
**Doctoral thesis submitted to
the Faculty of Behavioural and Cultural Studies
Heidelberg University
in partial fulfillment of the requirements of the degree of
Doctor of Philosophy (Dr. phil.)
in Psychology**

Title of the publication-based thesis
*Deep Learning Architectures for Amortized Bayesian Inference in
Cognitive Modeling*

presented by
Stefan T. Radev, M.Sc.

year of submission
2021

Dean: Prof. Dr. Dirk Hagemann
Advisors: Prof. Dr. Andreas Voss, Prof. Dr. Thorsten Meiser

CONTENTS

1	INTRODUCTION	3
1.1	Motivation and Scope	4
1.2	Contributions	5
1.3	List of Scientific Publications of the Publication-Based Dissertation	6
1.4	Notes on Notation	7
2	MODELS OF COGNITION	9
2.1	Cognition and Computation	9
2.2	Behavioral Simulators	12
2.3	The Likelihood	14
3	BAYESIAN INFERENCE	17
3.1	From Prior to Posterior	17
3.2	Uncertainty Reduction and Bayesian Surprise	19
3.3	Types of Uncertainty	21
3.4	Exchangeable Observations	23
3.5	Bayesian Model Comparison	24
3.5.1	The Prior Predictive and Occam’s Razor	25
3.5.2	The Posterior Predictive and Fortuna’s Knife	27
3.5.3	Bayesian Model Averaging and Wisdom of the Crowd	29
3.6	Bayesian Simulation-Based Inference	30
3.7	Samplers and Neural Samplers	31
4	RELATED WORK	33
5	AMORTIZED PARAMETER ESTIMATION WITH BAYESFLOW	37
5.1	Desiderata	37
5.2	Background	38
5.2.1	Deep Generative Modeling	38
5.2.2	Forward Inference	38
5.2.3	Normalizing Flows	39
5.2.4	Coupling Flows	40
5.2.5	Distribution Matching and Amortization	42
5.3	BayesFlow: Building Amortized Neural Samplers	43
5.3.1	Composing Invertible Networks	44
5.3.2	Summary Networks	45
5.3.3	Optimization Objective	49
5.4	Training Phase	51

Contents

5.5	Inference Phase	54
5.6	Sources of Error	55
5.7	A Bayesian Workflow with BayesFlow	58
5.7.1	Prior Consistency	58
5.7.2	Computational Faithfulness	58
5.7.3	Model Sensitivity	60
5.7.4	Posterior Predictive Checks	60
5.8	A Quick Demonstration	61
5.9	Concluding Remarks	64
6	AMORTIZED MODEL COMPARISON	65
6.1	Desiderata	65
6.2	Background	66
6.2.1	Bayesian Model Comparison	66
6.2.2	\mathcal{M} -Frameworks	67
6.2.3	Model Selection as Classification	68
6.2.4	Multi-Model Forward Inference	68
6.3	Training Amortized Evidence Approximators	69
6.3.1	Evidence Representation	69
6.3.2	Learning Evidence in an M-Closed Framework	72
6.3.3	Learning Absolute Evidence through Regularization	73
6.3.4	Implicit Preference for Simpler Models	74
6.3.5	Training and Inference	75
6.3.6	Sources of Error	76
6.4	A Simulated Experiment	76
6.4.1	Model Comparison Setting	76
6.4.2	Validation Results	78
6.5	Concluding Remarks	78
7	META-AMORTIZED INFERENCE	79
7.1	The Bayesian Hardships	79
7.2	Amortized Inference Revisited	81
7.3	Learning a Multi-Model Posterior	83
7.4	Use Cases and Challenges for Meta-Amortized Inference	86
8	APPLICATIONS	89
8.1	A Bayesian Brain Model of Adaptive Behavior	89
8.2	Jumping to Conclusion? A Lévy-Flight Model of Decision Making	91
8.3	Insights from Bayesian Modeling in a One Million Sample	92
8.4	OutbreakFlow: Model-Based Bayesian Inference of Disease Outbreak Dynamics	94
9	OUTLOOK	97
10	CONCLUSION	101

Contents

BIBLIOGRAPHY	105
APPENDIX A1 - MANUSCRIPT 1	117
APPENDIX A2 - MANUSCRIPT 2	137
APPENDIX A3 - MANUSCRIPT 3	159
APPENDIX A4 - MANUSCRIPT 4	167
DECLARATION IN ACCORDANCE TO § 8 (1) C) AND (D) OF THE DOCTORAL DEGREE REG- ULATION OF THE FACULTY	201

ACKNOWLEDGEMENTS

The intellectual wanderings transcribed in this thesis would not have been possible without the invaluable support and contributions of many people. First, I would like to thank my supervisor Andreas Voss for providing me with the unbounded freedom to choose and pursue my ideas and whims; for bearing with me and my diverging questions from the very start of my Masters'; for the indispensable discussions and for keeping an open mind for the future of our discipline.

Two other persons played a major supervisory role throughout my PhD life. Ullrich Köthe, whom I met during my computer science studies and whose teaching style and attitude towards students (and scientific life in general) I wholeheartedly admire, was always there for me (sometimes up at night) to discuss intellectual stuff, lift me up from reviewer and journal strikes, involve me in all sorts of crazy projects, and, most importantly, make me aware of my blind spots as a person and researcher. Paul-Christian Bürkner, whose workshop in Mannheim I accidentally attended, just as accidentally became my informal advisor in Bayesian statistics. During many discussions, he was able to clear up a whole lot of murkiness from my ideas, help me overcome a great deal of unknowledge, and, last but not least, repeatedly show me the need for having a Twitter account as a scientist (I still don't have one though, holding up...).

Perhaps I would not have commenced my PhD without the key influence of two very special psychologists. During my Masters', Ulf Mertens (together with the bold Ivan Marević) got me into deep learning and helped me discover my passion for programming languages, abstract problems, and long coding nights. Moreover, he and his wife, Alica Mertens, became good friends, beverage sharers, and idea tinkerers (they invited me to my first German wedding as well!). Also during my Masters', Veronika Lerche motivated me to fly to my first ever scientific conference in Coventry, where I, as an overdressed student assistant, presented a funny-looking graphical software to a bunch of stern researchers. Back then, she drew my attention to the bright sides of academia and got me thinking about staying for a bit longer than originally intended (unfortunately, we never started our very promising rock band).

I am thankful to my colleagues-turned-friends who accompanied me through this stage of life: Annika Stump, who represents the non-linearity of a creative life, who basically defined the taste of whiskey and always has an open ear for me; Marco D'Alessandro, the all-round Italian genius with an unceasing well of ideas in his pocket; Mischa von Krause, the master of work-life-balance; and Lukas Schumacher, my fellow investor, who reminded me of the unceasing enthusiasm and curiosity a scientist has to carry in his or her heart. Extra credits go to my friend-turned-colleague, Maximilian Knapp, and his amazing family.

In one way or another, I am indebted to all my international friends, collaborators, student assistants, and the entire SMiP community (especially David Izydorczyk and Martin Schnürch, who know how and what to drink during workshops and conferences).

Contents

At a whole different level, I am deeply grateful to my loving family, who never really let me feel away, and especially to my mother, who bravely struggled through all my boring papers and did not miss a single error.

Last but not least, this whole journey would have been unthinkable without Karin Prillinger, who was a bit of everything and everywhere.

1 INTRODUCTION

What I cannot create, I do not understand.

— Richard Feynman

Mathematical models are becoming increasingly important for describing, explaining, and predicting human behavior in terms of underlying mechanisms and systems of mechanisms. Although the ontology of such mechanisms remains largely unknown, their epistemic value and inferential power are now widely acknowledged throughout the behavioral sciences. Broadly speaking, whenever an assumed mechanism transforms information into behavior, it is referred to as a *cognitive process*. Cognitive processes are the conceptual fabric used to fill the explanatory gap between the mysterious firing of neurons and the mundane recognition of a long-forgotten acquaintance in the morning train. Consequently, modelers of cognitive processes earn their livelihood in an attempt to make the “ghost in a machine” tractable by replacing the ghost with *hidden parameters* embedded in an abstract functional framework.

The purpose of such parametric models is twofold. On the one hand, they can be viewed as formal expedients for understanding the messy and noisy human data in much the same way as the models physicists employ to make sense of the data coming from spiral galaxies and interstellar clouds. On the other hand, parametric models can be viewed as *behavioral simulators* and used to mimic the output of cognitive processes by generating synthetic behavior. Interestingly, there is a strange asymmetry in the challenges surrounding these two goals. Simulating behavior requires only specifying a cognitive model as a computer program and running the program with a desired parameter configuration. It is thus a generative process mainly constrained by the creativity and imagination of individual modelers. Differently, reverse engineering human data to recover hidden parameters is hampered by two external factors: the resolution and abundance of data and the availability of universal and efficient inferential methods. As for the latter, behavioral scientists have often sacrificed fidelity and complexity in order to adjust their models not to reality but to the limitations of existing inferential methods. Such a strategy is definitely viable in the early (often linear and beguilingly clear) stages of scientific inquiry, but it does not live up to the challenges and questions posed by later (often non-linear and disconcertingly fuzzy) stages.

The main argument of this thesis is that *questions of inferential tractability are of secondary importance for enhancing our understanding of the processes under study*. Accordingly, the core purpose of this thesis is to develop frameworks which leave such questions to specialized “black-box” artificial neural networks and enable researchers to focus on developing and validating faithful “white-box” models of cognition. Instead of a ready-made solution, the thesis explores a beginning of a solution. It presents a potentially fruitful coupling between human and artificial intelligence, an approach which is expected to gain more and more momentum as the world fills with artificial agents. Ultimately, this thesis strives to increase creativity by embracing complexity.

1.1 MOTIVATION AND SCOPE

This thesis is motivated by the question of how to offload *parameter estimation* and *model comparison* onto specialized neural network architectures. Undoubtedly, parameter estimation and model comparison are two of the most common and most challenging tasks in model-based behavioral sciences. Accordingly, for each of the two tasks, we will offer a general framework for composing neural networks capable of bootstrapping a wide range of inference tasks. The main principle will be to utilize prior domain expertise and guide these networks through simulations to become “experts” in inferring hidden parameters from data or selecting between plausible models of the data. The proposed frameworks are themselves embedded into the meta-framework of Bayesian inference which embraces probability theory as *the logic of science* [77].

Why probability theory? Simply put, a probabilistic approach to inference is appealing, as it provides consistent equations and principled methods for quantifying and communicating uncertainty [22]. Correspondingly, doing Bayesian inference and data analysis is nothing but applying the basic rules¹ of probability theory to amount of information gained through empirical inquiry. Curiously, it is primarily in the behavioral sciences that researchers applying probability theory are given the cultist label *Bayesians* and often seen as representatives of a statistical opposition against traditional (the cultist label being *frequentist*) methods. Accordingly, whenever the reader encounters the term Bayesian in this thesis, it should be read as *using the rules of probability theory to express uncertainty, update beliefs, revise knowledge, and inform scientific conclusions*.

Probabilistic reasoning is hard and time-consuming. It was not until general-purpose computers had shrunk considerably in size that Bayesian inference became useful for handling non-trivial practical problems. Until then, practitioners could leverage only a limited subset of the tools probability theory had to offer. Moreover, this restricted inference was further constrained by the ability to solve complicated integrals, or, as David MacKay puts it: “...a macho activity enjoyed by those who are fluent in definite integration” [103, p. 319].

With the advent of high-performance computing², Monte Carlo methods came to the rescue of probabilistic inference, the most prominent algorithmic family being Markov chain Monte Carlo (MCMC) methods [109]. Initially, MCMC proved instrumental in approximating the unimaginable integrals which had been thwarting the solution of relevant problems in chemical physics [141]. Today, the heirs of those rather crude grandfather sampling methods have become the Bayesian gold standard across the sciences, with novel and interesting modifications spawning in scientific journals on a regular basis. The behavioral sciences are no exception to this trend, being a field of both active application and development of novel Bayesian methods.

Notwithstanding the major contributions of MCMC methods to large-scale inference problems, they remain notoriously slow and sequential in nature [12]. These drawbacks can render inference with highly complex models practically infeasible, but they also transfer to applications of relatively simple models to big data where relevant computations need to be repeated for each observation. Consequently, it is not uncommon for researchers to wait a week or two for estimation algorithms to finish, only to notice afterwards that critical adjustments to the algorithm or the

¹Even though the rules of probability theory themselves are basic and intuitive, the application of these rules in practice can be anything but basic.

²The term *high-performance* having, of course, only a temporal meaning within the context of a computing generation.

model necessitate rerunning the entire loop from scratch. The inferential matters become even worse, when the model itself cannot be specified in a nicely closed, analytic form, but is only available as a Monte Carlo simulator program. The latter can often leave potentially relevant modeling territories vastly unexplored and confine certain model classes to a purely Platonic playground. Our goal is to provide an efficient and scalable framework for designing and testing solutions to precisely those challenging situations encountered frequently by behavioral scientists. Moreover, we believe that our ideas can positively impact computational modeling in research areas not solely confined to the behavioral sciences. At a high level, our framework is purely simulation-based and leverages the representational power of deep learning methods to build reusable estimators for two of the most important constructs in Bayesian inference: the *posterior distribution* and the *evidence*. Further, it utilizes the concept of *amortized inference* to increase the inferential efficiency of these estimators at every modeling step, from model development, to model selection. Importantly, our framework includes methods for self-diagnosis of miscalibrated inference due to algorithmic errors, which is an essential precondition for computational faithfulness in any Bayesian data analysis pipeline.

Beyond the development of a novel Bayesian framework for simulation-based inference, this thesis presents some concrete applications to relevant research questions in cognitive modeling and beyond. These applications demonstrate the utility of the framework for tackling challenging cognitive models dealing with both simple (typically independent and identically distributed, or *i.i.d.*) and more complex probabilistic structure (typically exhibiting temporal dependencies, or non-*i.i.d.*). Moreover, the models considered in these applications are themselves novel and serve the purpose to inspire further research and exploration in the corresponding areas.

Finally, the thesis includes a starter Python library for building own estimation or model comparison networks with minimal programming skills.

1.2 CONTRIBUTIONS

The main contributions of this thesis can be summarized as follows:

- Chapter 3 provides the necessary background and conceptual machinery for understanding uncertainty quantification in the context of Bayesian inference. It then discusses the challenges faced by standard Bayesian methods when applied to complex models and how these affect the field of cognitive modeling in particular. We introduce the concepts of simulation-based and amortized inference with neural samplers and set the stage for presenting our BayesFlow framework.
- Chapter 5 introduces our general BayesFlow framework for solving the task of amortized Bayesian parameter estimation. We demonstrate how to perform inference on data sets with different sizes and probabilistic structure by using specialized network architectures which preserve the probabilistic symmetry of the target Bayesian posterior. We formally derive a training procedure which ensures that neural networks in our framework recover the true target posteriors under perfect convergence of the optimization algorithm. We end the chapter with a simulation-study demonstrating the utility of our method.

- Chapter 6 introduces our *Dirichlet* evidence network for solving the task of amortized Bayesian model comparison. We explore a method to quantify absolute evidence as compared to relative evidence through a specific form of regularization in a meta-probabilistic framework. As in chapter 5, we show how to deal with variable numbers of observations and different model/data types. We also derive a simulation-based training method which ensures that evidential networks in our framework recover the true model probabilities under perfect convergence of the optimization algorithm. We end the chapter with two simulation studies using complex computational models from cognitive science and neuroscience.
- Chapter 7 introduces a visionary approach towards *meta-amortized inference*. It combines both parameter estimation and model comparison into a single unifying framework and presents initial conceptual results.
- Chapter 8 presents applications of the proposed Bayesian frameworks for model-based inference on real data. It starts with a direct estimation of an information-theoretic model of adaptive performance inspired by the Bayesian Brain Theory (BBT). Then, we describe a parameter estimation study concerned with a set of novel models of decision making. This is followed by an application of a custom diffusion model to a massive data set of human response times to disentangle questions of cognitive aging. Finally, we present an application to Covid-19 outbreak modeling with a version of BayesFlow for dynamic models.

1.3 LIST OF SCIENTIFIC PUBLICATIONS OF THE PUBLICATION-BASED DISSERTATION

The central ideas put forward in this thesis have been explored in the following publications by the author and his cooperators:

- S. T. Radev, U. K. Mertens, A. Voss, L. Ardizzone, and U. Köthe. “BayesFlow: Learning complex stochastic models with invertible neural networks”. *IEEE Transactions on Neural Networks and Learning Systems*, 2020, pp. 1–15. DOI: [10.1109/TNNLS.2020.3042395](https://doi.org/10.1109/TNNLS.2020.3042395)
- S. T. Radev, M. D’Alessandro, P.-C. Bürkner, U. K. Mertens, A. Voss, and U. Köthe. “Amortized Bayesian model comparison with evidential deep learning”, Manuscript submitted for publication in *IEEE Transactions on Neural Networks and Learning Systems*
- S. T. Radev, A. Voss, E. M. Wieschen, and P.-C. Bürkner. “Amortized Bayesian inference for models of cognition”. *International Conference on Cognitive Modelling (ICCM) Conference Proceedings*, 2020
- M. D’Alessandro, S. T. Radev, A. Voss, and L. Lombardi. “A Bayesian brain model of adaptive behavior: an application to the Wisconsin Card Sorting Task”. *PeerJ* 8, 2020, e10316

The author also contributed to the following publications which are related to the core topics of the current thesis:

- M. von Krause, S. T. Radev, and A. Voss. “Processing speed is high until age 60: insights from Bayesian modeling in a one million sample (with a little help of deep learning)”, Manuscript submitted for publication
- S. T. Radev, U. K. Mertens, A. Voss, and U. Köthe. “Towards end-to-end likelihood-free inference with convolutional neural networks”. *British Journal of Mathematical and Statistical Psychology* 73:1, 2020, pp. 23–43
- E. M. Wieschen, A. Voss, and S. Radev. “Jumping to conclusion? a lévy flight model of decision making”. *TQMP* 16:2, 2020, pp. 120–132
- S. T. Radev, F. Graw, S. Chen, N. Mutters, V. Eichel, T. Bärnighausen, and U. Köthe. “Model-based Bayesian inference of disease outbreak with invertible neural networks”. *arXiv preprint arXiv:2010.00300*, 2020
- S. Bieringer, A. Butter, T. Heimel, S. Höche, U. Köthe, T. Plehn, and S. T. Radev. “Measuring QCD splittings with invertible networks”. *arXiv preprint arXiv:2012.09873*, 2020
- L. Konicar, S. Radev, K. Prillinger, M. Klöbl, R. Diehm, N. Birbaumer, R. Lanzenberger, P. Plener, and L. Poustka. “Volitional modification of brain activity in adolescents with Autism Spectrum Disorder: A Bayesian analysis of Slow Cortical Potential neurofeedback”. *NeuroImage: Clinical*, 2021, p. 102557
- U. K. Mertens, A. Voss, and S. Radev. “ABrox—A user-friendly Python module for approximate Bayesian computation with a focus on model comparison”. *PLoS one* 13:3, 2018, e0193981

1.4 NOTES ON NOTATION

Throughout this thesis, we will follow some simple conventions for consistent mathematical notation. We will denote scalar variables by lowercase italic, e.g., x , y , z , vectors by lowercase bold italic, e.g., \mathbf{x} , \mathbf{y} , \mathbf{z} , and matrices by uppercase bold italic letters, e.g., \mathbf{X} , \mathbf{Y} , \mathbf{Z} . Data sets comprising multiple observations (e.g., multivariate responses of a single participant to a particular task) will be denoted as $\{\mathbf{x}_n\}_{n=1}^N = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \equiv \mathbf{x}_{1:N}$, where N indicates the number of observations. Occasionally, and when possible, we might stack all observations comprising a data set row-wise into a matrix, $\{\mathbf{x}_n\}_{n=1}^N \equiv \mathbf{X}$. Whenever the observations are assumed to be *time-dependent*, we will use T to denote the total number of observations in the resulting (multivariate) time-series $\mathbf{x}_{1:T}$. Occasionally, we will include the superscript *obs* to denote an actually observed data set, in contrast to a simulated one (i.e., $\mathbf{x}_{1:N}^{(obs)}$ vs. $\mathbf{x}_{1:N}$).

We will always collect the parameters of a mathematical model into a vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_D)$, and reserve the letter D for the dimensions of the parameter space. Finally, we will collectively refer to all trainable parameters of a neural network (e.g., weight matrices, biases, activation function parameters) as a vector (e.g., $\boldsymbol{\phi}$, $\boldsymbol{\psi}$, ...) even though these might be distributed across different functional components or layers of the network. Importantly, neural network parameters are *not to be confused* with the parameters of the mathematical model of interest, as the former are uninterpretable and high-dimensional, whereas the latter are carriers of theoretical value and usually low-dimensional.

Table 1.1: Table of important symbols and their corresponding description

Notation (Symbol)	Meaning
$\mathbf{x}, \mathbf{x}_{1:N}$	observed or simulated data point, data set
$\boldsymbol{\theta}, \boldsymbol{\omega}$	parameters of a mathematical model (simulator)
\mathbf{z}	latent variable learned by a (deep) generative model
$\mathcal{X}^N, \Theta, \Xi$	data space, parameter space, noise space
$\mathcal{M}, \mathcal{M}_j$	candidate model set, model index
g	generative (forward) model / simulator
p, q	probability density (mass) functions
ϕ, ψ	trainable parameters of neural networks
f_ϕ, h_ψ	functions parameterized via neural networks
$\mathbb{E}[\cdot]$	expected value of a random variable (vector)
$\mathbb{KL}[p q]$	Kullback-Leibler divergence between densities p and q
N	number of observations in a (simulated) data set
D	number of parameters / dimensions of the parameter space
B	number of simulations per training step / batch size

For the most part of this thesis, we will be concerned with (absolutely) continuous random vectors and their associated probability density functions (pdfs). For the sake of readability, the latter will be denoted by p even when they refer to pdfs of different random vectors defined on different spaces, which will be clear from the function arguments. For instance, we will write $p(\boldsymbol{\theta})$ for the (prior) probability density of the parameter vector $\boldsymbol{\theta} \in \Theta$ instead of $p_\Theta(\boldsymbol{\theta})$. Additionally, each pdf of interest will be implicitly associated with a corresponding probability measure P . Throughout the text, we will use density and distribution interchangeably.

By means of a slight abuse of notation, when a density function is approximated via a neural network with trainable parameters ϕ , we will often write $q_\phi(\boldsymbol{\theta}) \equiv q(\boldsymbol{\theta} | \phi)$, or, for a conditional density, $q_\phi(\boldsymbol{\theta} | \mathbf{x}) \equiv q(\boldsymbol{\theta} | \mathbf{x}, \phi)$. In this way, (i) we align our notation to the predominant notation in the literature on deep generative modeling; (ii) implicitly denote the dependence of the approximate density on the neural network parameters; (iii) make it immediately clear which is the density being approximated, e.g., $q_\phi(\boldsymbol{\theta} | \mathbf{x})$ approximates $p(\boldsymbol{\theta} | \mathbf{x})$ by means of neural network parameters ϕ .

The most important symbols and notation used throughout the text are summarized in Table 1.1.

2 MODELS OF COGNITION

Murky thoughts, like murky waters, can serve two purposes only: to hide what lies beneath, which is our ignorance, or to make the shallow seem deep.

— Giulio Tononi

Reasoning with models of empirical phenomena lies at the very heart of science. Models abstract away irrelevant details and focus attention on the theoretically relevant aspects of complex systems. Ideally, they simplify, but do not oversimplify reality. The importance of models for scientific progress is twofold. On the one hand, theories can be systematically instantiated and tested by specifying a mathematical model and inferring its hidden properties from data. On the other hand, competing theories can be tested against one another via formal model comparison. Thus, model-based reasoning complements verbal reasoning insofar as it reduces ambiguity and translates “murky” statements into precise and directly quantifiable hypotheses. Whether the latent properties of cognitive models represent faithful descriptors of the *unobservable causes* of behavior remains an open question whose surface we will only scratch here. The main purpose of this chapter is to establish the notion of a cognitive model, fix a useful notation, and introduce the concept of a likelihood function.

2.1 COGNITION AND COMPUTATION

Cognitive models exist to help cognitive scientists make sense of observed behavioral data in terms of unobservable (latent) cognitive processes, such as attention, memory decay, evidence accumulation, or belief updating, to name just a few [45]. Such models have been around for centuries (see Figure 2.1), with the most notable modern twists being a shift in contextualization (i.e., in the reference theoretical framework) and an increase in mathematical formalism. Whether cognitive processes actually exist as functional entities or simply represent useful metaphors psychologists live by, is a question that currently extends from philosophy down to single-cell neuroscience.

In the year 1994, Francis Crick, essentially reinventing naturalism, formulated his *astonishing hypothesis* in a rather flowery way:

The Astonishing Hypothesis is that “You”, your joys and your sorrows, your memories and your ambitions, your sense of personal identity and free will, are in fact no more than the behavior of a vast assembly of nerve cells and their associated molecules. As Lewis Carroll’s Alice might have phrased it: “You’re nothing but a pack of neurons”. [28, p. 3]

At the time of writing, hardly any cognitive scientist would find the relationship between brain, cognition, and behavior alien or astonishing. In fact, this relationship is the explicit or implicit

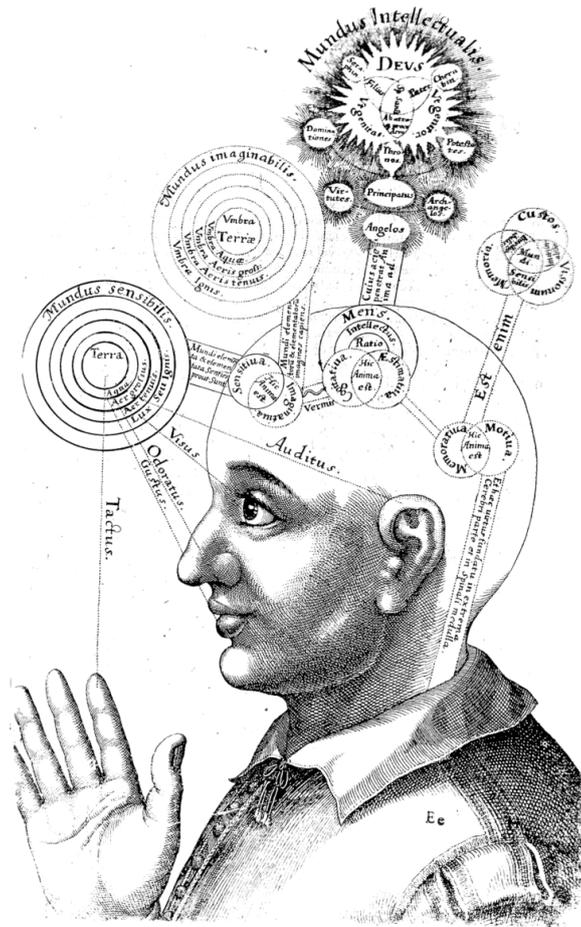


Figure 2.1: Robert Fludd’s “model” of a mind in the world from 1619 [15]. The diagram can be seen as an early predecessor of modern cognitive architectures incorporating cognitive functions such as perception, memory as well as an explicit flow of information between the mind and the world. Cognitive functions are localized in the head, even though neurons are yet to be discovered.

working hypothesis behind some of the most prominent cognitive architectures and model classes [91, 153, 162]. However, few would contend that even a complete description of all microtubules in each and every neuron in the brain will ever be sufficient to explain why a grandmaster under time pressure failed to spot an obvious checkmate in a crucial game of chess. Invoking such an explanation appears to require more than just a description of the behavior of molecules and neurons. In other words, the “no more” and the “nothing but” parts of the astonishing hypothesis are what provokes a certain theoretical dissatisfaction. Perhaps going up to a complete understanding of neural firing patterns would resolve the dissatisfaction. But how far up should one climb the reductionist ladder until a satisfactory level of analysis for cognition is reached? This is an example of the so called *bottom-up* problem, starting from the small and tractable constituents of a presumably complex process and building *up* towards its synthesis.

In contrast, one may start with a handy catalogue of cognitive processes inherited from the dictionary of psychology and look *down* for the “neural correlates” of these processes. Necessarily, such an approach rests on the assumption that there will be a serendipitous one-to-one correspondence between cognitive and neural processes. This is an example of the *top-down* problem, starting from the processes and analysing them in terms of their constituents.

Undoubtedly, both approaches bear obvious risks. Following a purely bottom-up research agenda, we might run into the problem of being unable to reconstruct the entire list of assumed cognitive processes. Following a purely top-down approach, we might not end up discovering the neatly corresponding constituents we are looking for, leaving us with a list of substance-free metaphors.

It is perhaps no coincidence, that neuroscientists favor a predominantly bottom-up approach. For instance, György Buzsáki [18] outlines the three missing pieces for a purportedly complete understanding of how the brain generates behavior. These are the understanding of (i) the dynamical structural organization of the brain; (ii) the physiological functions of its constituents; (iii) and the computational mode of operation that enables the neurons in a given anatomical hardware to execute actions [18, p. 24]. The eventual synthesis of this knowledge is expected to provide a satisfactory explanation for all kinds of behavior in terms of underlying neural mechanisms. In the process of studying and systematizing these mechanisms, only the neurophysiologically plausible cognitive constructs would therefore stand the test of rigorous empirical verification.

Cognitive scientists, on the other hand, prefer Marr’s interpretative framework for talking about cognitive processes [106]. Its starting point is the basic concept of *information processing*, that is, the detectable (i.e., the difference that makes a difference) transformation of information (e.g., wavelengths becoming cone cell responses or the ringing of the phone becoming an increase in heart-rate variability). It then distinguishes three levels of analysis for understanding any information processing system: (i) *computational theory* - concerning the goals (the why) of the system and its computational logic; (ii) *representation and algorithm* - concerning the representation of input and output as well as the step-by-step instructions for carrying out the input-output transformation (the how); (iii) *hardware implementation* - concerning the physical realization of the algorithm [106, p. 25].

Accordingly, we can equate cognitive processes with the algorithms transforming neural representations into observed behavior (or into further representations), without committing to a particular physical ontology. Thus, computational models of cognition (cognitive models, for short) are our best instruments to learn something about these algorithms from behavioral data alone, without resorting to the expensive methods of neuroscience. Ideally, the functional form and parameters of cognitive models would capture the most important algorithmic and representational characteristics of the system under study. At the same time, the assumed hardware underpinnings of a cognitive process should impose a number of parametric and functional constraints (e.g., the maximal speed of information processing), if a model of the process is to be taken seriously by substantive neuroscientists.

Cognitive models are occasionally termed *computational models*, highlighting their algorithmic essence and their conceptual relatedness with notions borrowed from computer science. Even though neurophysiological plausibility seems to be a prerequisite for the ultimate validation of any cognitive model, it is often ignored in favor of a strong embedding in the nomological nexus of psychology. Accordingly, as long as the parameters of a cognitive model are interpretable with

a reference to an overarching psychological theory, the actual neuroanatomical hardware becomes of secondary importance. In addition, neurophysiological plausibility becomes even less important if the presumed cognitive processes are eventually to be implemented by goal-directed artificial agents. In the latter case, an artificial agent may perfectly mimic human or animal behavior without the need or physical possibility to invoke any neural representations. As a consequence, the concept of a cognitive process may become decoupled from its neural realization and thus become synonymous to an algorithm, which is per definition independent of its implementation. Nevertheless, for those striving to infer something about real brains from real behavioral data, the hope remains that their models are at least able to tap into the distant echo of neural system parameters within a coherent “neurocognitive” framework.

Future theoretical developments may also see *phenomenological plausibility* come as an additional criterion for models representing cognitive processes which are *experienced* in some way (e.g., the experience of coming to a decision as compared to all non-experienced factors that contributed to the decision). However, the fact that information processing is experienced in a particular way seems less important for cognitive modelers than the functional task of describing how information is transformed in a way that ultimately leads to manifest behavior. And even though simulating experience appears to be unimaginably hard (except perhaps for the brain), we can easily “build” abstract machines which simulate behavior in various experimental and observational contexts by executing a series of well-defined computational steps. It is this property which makes model building a creative process and cognitive models *generative* in nature.

2.2 BEHAVIORAL SIMULATORS

Formally, we can represent a cognitive model as a function $g : \Theta \times \Xi \rightarrow \mathcal{X}$ which generates observable quantities $\mathbf{x} \in \mathcal{X}$ from a particular configuration of the hidden parameters $\boldsymbol{\theta} \in \Theta$ and an independent source of noise $\boldsymbol{\xi} \in \Xi$. The function is typically realized as a Monte Carlo computer simulation which mimics quantifiable manifestations of actual behavior. Simulation programs involving random number generation are also known as Monte Carlo engines or probabilistic programs.

Since humans rarely behave the same way even when provided with the same information, the stochastic component $\boldsymbol{\xi}$ in the model formulation ensures that g is non-deterministic given the same time-invariant parameter configuration $\boldsymbol{\theta}$. Ideally, the noise in a cognitive model should capture all non-modeled factors which nevertheless influence the generation of behavior, but it may also reflect inherent randomness of the cognitive system under study. In the latter case, the noise distribution might itself contain learnable parameters which are part of $\boldsymbol{\theta}$, so it should be denoted as $p(\boldsymbol{\xi} | \boldsymbol{\theta})$. Thus, the general functional form of a (stochastic) cognitive model is:

$$\mathbf{x}_n = g(\boldsymbol{\theta}, \boldsymbol{\xi}_n) \quad \text{with} \quad \boldsymbol{\xi}_n \sim p^*(\boldsymbol{\xi}) \quad (2.1)$$

where the subscript n indicates that the simulator can be run repeatedly with the same parameter vector to yield an entire data set $\{\mathbf{x}_n\}_{n=1}^N \in \mathcal{X}^N$ and p^* denotes the true underlying noise distribution. Usually, this distribution is unknown and its form needs to be either theoretically deduced or approximated *ad hoc* via a simpler distribution p from which random samples can easily be obtained (e.g., Gaussian, Poisson).

The model form in Equation 2.1 represents a *memoryless* or a stateless process: each run of the simulator with a fixed parameter combination is independent of the previous runs and does not influence future runs. As we will show later, such models generate data sets consisting of independent and identically distributed (*iid*) observations. More complex models involving some form of memory can be formulated by a recursive dependence of the simulator on the previous observation:

$$\mathbf{x}_n = g(\mathbf{x}_{n-1}, \boldsymbol{\theta}, \boldsymbol{\xi}_n) \quad \text{with} \quad \boldsymbol{\xi}_n \sim p^*(\boldsymbol{\xi}) \quad (2.2)$$

or by introducing a persistent memory variable $\boldsymbol{\omega}_n$ which is updated after each simulation run and may itself be unobserved and treated as a time-varying parameter:

$$(\mathbf{x}_n, \boldsymbol{\omega}_n) = g(\boldsymbol{\omega}_{n-1}, \boldsymbol{\theta}, \boldsymbol{\xi}_n) \quad \text{with} \quad \boldsymbol{\xi}_n \sim p^*(\boldsymbol{\xi}) \quad (2.3)$$

Importantly, such *stateful* models give rise to more complex sets of observations which are no longer *i.i.d.* and thus need to be tackled differently than *i.i.d.* observations. Accordingly, a major aim of this thesis is to develop and validate a framework for performing model-based inference on both *i.i.d.* and non-*i.i.d.* observations. In anticipation of future modeling challenges including joint models of neural and behavioral data as well as more unstructured data such as graphs, text, or motion time-series, we will ensure that our framework is extendable to incorporate these future challenges. We further anticipate that our ideas will be potentially useful in different fields having embraced model-based reasoning and inference.

Scientifically useful cognitive models should work both forward and backward. Given a set of parameters, researchers should be able to generate artificial observations and data sets which are indistinguishable from their real counterparts even when scrutinized by expert observers or subjected to rigorous statistical tests. Conversely, given only a set of observations, researchers should be able to recover the hidden data-generating parameters which are seen as epistemically valuable *explanans* of the data. We call this the *inverse inference* problem.

As already briefly mentioned in the introduction, the two tasks are notably asymmetric with respect to the computational and epistemic burden they carry. To further appreciate this asymmetry, consider the example of an ice cube left at room temperature to eventually become a puddle of water [154, p. 196]. Having a model of the process of fusion, one can feasibly simulate how an arbitrary ice cube transitions into a puddle over time (forward inference). However, suppose that you are presented with only a puddle and tasked to recover an unobserved cube (inverse inference). Even though you can infer something about the volume or the density of the ice cube, the task as a whole appears insurmountable, since there are different cubes that could have melted into the same puddle, which presents a case of inherent non-identifiability. What is more, there is uncertainty about the model of the ice cube itself, since there is a multitude of ice forms, not necessarily cubic, that could have resulted in the same puddle.

Researchers striving to understand and model how “the mind can occur in a physical universe” [2] face a similarly tough challenge. When it comes to modeling cognition and behavior, this challenge is further aggravated by the possibility that *there could have been no ice cube to begin with*, but an entirely different puddle-generating process.

In short, forward inference is easy, whereas inverse inference is hard. Forward inference requires “only” the ability to write down the model as a simulator program in a programming language and

run the simulation with the desired set of parameters. It also requires creativity and theoretical insight (as well as acceptance by the community of other modelers). Inverse inference, on the other hand, requires invoking additional estimation procedures, largely external to the model or the modelers themselves. Further, it implies facing a large portion of epistemic uncertainty. In addition, estimation procedures are impeded by the following properties of complex simulators:

1. The simulator is typically non-deterministic, so that there is intrinsic uncertainty about the true value of θ .
2. The simulator is typically not information-preserving, so that there is ambiguity among possible values of θ .
3. The simulator is typically too complex to admit a closed-form mathematical expression for evaluating the likelihood of θ .

A multitude of inference frameworks with a varying degree of generality have been proposed for addressing problems of inverse inference. The most prominent among these are manual tuning, parameter search methods, kernel methods, optimization methods, maximum likelihood, maximum-a-posteriori (MAP inference), variational inference, fully Bayesian inference, approximate Bayesian computation (ABC), machine learning approaches as well as various hybrid methods [12, 29, 52, 105, 123, 136, 140, 142, 150]. A comprehensive review of the multiverse of methods for inverse inference is beyond the scope of this work (but see Chapter 4 or the excellent review by [27] for a broad classification). Essentially, all approaches for inverse inference optimize a trade-off between efficiency, scalability, accuracy, and practical utility. As there is no free lunch in capital markets, there is no silver bullet in statistical inference. Put differently, there is no method or framework that simultaneously maximizes all of the above criteria and the *best method* will be application-dependent.

This thesis focuses on scaling and utilizing fully Bayesian inference for very complex (often deemed intractable) models exhibiting all three inferential predicaments. In the next chapter, we will see why such models necessitate a *simulation-based* approach to Bayesian inference. Before we conclude this chapter, however, a word on the concept of *likelihood* seems warranted.

2.3 THE LIKELIHOOD

When looked upon through a probabilistic lens, the outputs of a cognitive model can be associated with some (potentially very complex) probability distribution. This distribution is referred to as the *likelihood* and denoted as $p(\mathbf{x} | \theta)$. Loosely speaking, when evaluated, the likelihood returns the relative probability of an observation \mathbf{x} given a set of parameters θ . When the parameters are systematically varied, the likelihood can be used to quantify how well each model instantiation fits the data.

If the likelihood has a known distributional form (e.g., Gaussian, Laplace, Dirichlet), the model in Equation 2.1 can be formulated entirely in terms of the likelihood. Moreover, in these simple cases, the likelihood can be evaluated analytically or numerically for any pair (\mathbf{x}, θ) . In a sense, all stochastic models have a *dual generative representation*:

$$\mathbf{x}_n \sim p(\mathbf{x} | \theta) \iff \mathbf{x}_n = g(\theta, \xi_n) \text{ with } \xi_n \sim p(\xi) \quad (2.4)$$

where the likelihood function $p(\mathbf{x} | \boldsymbol{\theta})$, being a probability distribution itself, naturally captures the effects of the extrinsic stochastic factor $\boldsymbol{\xi}$. Note, that this representation can be trivially extended to incorporate models with memory (Equations 2.2 and 2.3). In fact, each stochastic model viewed as a Monte Carlo simulator defines an implicit likelihood given by the relationship:

$$p(\mathbf{x} | \boldsymbol{\theta}) = \int_{\Xi} \delta(\mathbf{x} - g(\boldsymbol{\theta}, \boldsymbol{\xi})) p(\boldsymbol{\xi} | \boldsymbol{\theta}) d\boldsymbol{\xi} \quad (2.5)$$

where $\delta(\cdot)$ is the Dirac delta function and the integral runs over all possible execution paths of the stochastic simulation for a fixed $\boldsymbol{\theta}$. For most complex models, this integral is analytically intractable or too expensive to approximate numerically, so it is much easier to specify the model directly in terms of the simulation program g instead of using the likelihood. However, according to Equation 2.4, we can still *sample* from the likelihood by running the simulator with different Monte Carlo realizations of $\boldsymbol{\xi}$.

It is at this point where likelihood-based and simulation-based inference methods diverge. The former require the ability to evaluate the likelihood for any pair $(\mathbf{x}, \boldsymbol{\theta})$. The latter require only the ability to sample (simulate) arbitrary pairs $(\mathbf{x}, \boldsymbol{\theta})$ from the likelihood¹. In the next chapter, we will see how the inability to evaluate or derive the likelihood prohibits standard Bayesian methods. We will then briefly peruse the frontier of simulation-based inference before introducing our frameworks for parameter estimation and model comparison.

¹Although the literature has adopted the term *likelihood-free* inference, we will completely avoid it in this text, since it incorrectly implies the absence of a likelihood, even though the likelihood is implicitly defined via the action of the simulation program. It is the calculation of its actual numerical value for simulated or real observations which is impossible. We will therefore prefer the term *simulation-based* inference, since it unambiguously captures the essence of the task.

3 BAYESIAN INFERENCE

It is an extraordinary feature of science that the most diverse, seemingly unrelated, phenomena can be described with the same mathematical tools.
— Benoit Mandelbrot

Probability theory is the language for describing uncertainty with mathematical objects. Uncertainty is not merely a philosophical phantasm or a convenient excuse for the erroneous forecasts of meteorologists, but rather a fundamental epistemic basis for decision making in an incomplete universe. Scientists face uncertainty both in their academic and private lives. Finite data, approximation errors, noisy measurements, and inherently stochastic models are the common culprits for cultivating the habit of reporting some confidence measures alongside point estimates. But also everyday questions such as "Will it rain tomorrow?", "Which party is likely to win the election?", or "What is the likelihood of encountering a dragon in the park tonight?" call for reasoning with a varying degree of confidence. Since the core topic of this thesis is model-based inference, we will see in this chapter how Bayesian methods provide a self-consistent framework for uncertainty quantification and communication when trying to extract cognitive models from behavioral data. Apart from the parlor of behavioral sciences, Bayesian methods have been employed for tasks as diverse as predicting global equity indices [73], inferring latent infectious disease dynamics [42], and evaluating forensic DNA profiles [9]. Henceforth, we will assume the utility of Bayesian methods as given and proceed to the conceptual and mathematical details of Bayesian inference.

3.1 FROM PRIOR TO POSTERIOR

The core mathematical workhorse in Bayesian probabilistic reasoning is the famous Bayes' rule [52], which specifies how to update prior knowledge about a given quantity upon making an informative observation. In cognitive modeling, the quantities of interest are the parameters of a cognitive model, which capture relevant computational characteristics of the process under scientific scrutiny. Thus, when collecting behavioral data, we ultimately strive to increase our knowledge about these parameters as a proxy for understanding the assumed cognitive processes.

The point of departure in Bayesian inference is the *prior distribution* (or just *prior*¹, for short), denoted as $p(\theta)$. Ideally, the prior is supposed to capture both what we already (believe to) know about the parameters, but also what we still do not (believe to) know. Knowledge is expressed through a reasonable domain and concentration of probability density (i.e., our "best" guess). Lack of knowledge is expressed through the spread of probability density over the domain (i.e.,

¹Whenever a model has a multi-dimensional parameter space, one speaks of a *joint prior* highlighting the fact that the the distribution p extends over multiple parameters (dimensions). Correspondingly, one speaks of a *joint posterior* when referring to the updated distribution of multiple parameters.

our uncertainty about the “best” guess). To make this idea concrete, Figure 3.1 depicts three one-dimensional priors with the same average value of zero but different allocations of uncertainty and belief.

Even though the prior is typically defined as an unconditional density, this definition is merely an aesthetic consideration, since knowledge and uncertainty are fundamentally contextual. Thus, $p(\boldsymbol{\theta})$ is really a shorthand for $p(\boldsymbol{\theta} | C)$, where C is a placeholder for all contextual information, such as previous research, theoretical constraints, model assumptions, and cultural preferences. Indeed, one of the greatest appeals (and, perhaps, at the same time, greatest deterrents) of Bayesian inference is that it allows to systematically and explicitly incorporate contextual information.

In a sense, Bayesian inference does away with the illusion of objectivity by allowing subjectivity to be expressed explicitly and transparently. Accordingly, one might question the specification of $p(\boldsymbol{\theta} | C)$ on empirical or theoretical grounds and propose $p(\boldsymbol{\theta} | C')$ as an alternative. In this case, the difference in contextual information leads to different epistemic states *prior* to observing any new data and potentially different conclusions *after* observing some data².

The process of knowledge updating in Bayesian inference essentially consists in transforming the prior into the *posterior* according to the well known Bayes’ rule:

$$p(\boldsymbol{\theta} | \mathbf{x}) = \frac{p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathbf{x})} = \frac{p(\boldsymbol{\theta}, \mathbf{x})}{\mathbb{E}_{p(\boldsymbol{\theta})}[p(\mathbf{x} | \boldsymbol{\theta})]} \quad (3.1)$$

where $p(\mathbf{x} | \boldsymbol{\theta})$ denotes the likelihood function, as discussed in the last chapter, and $p(\boldsymbol{\theta} | \mathbf{x})$ denotes the posterior given some observation \mathbf{x} . In the context of model-based inference, attaining the posterior corresponds to Bayesian parameter estimation and is sometimes referred to as *inverting the likelihood*, due to the reflection of the arguments across the $|$ symbol. The denominator $p(\mathbf{x})$ in the second term of Equation 3.1 is known as the *evidence* and represents a normalizing constant for the posterior. The equivalent denominator in the third term rephrases the evidence as the expectation of the likelihood with respect to the prior, that is, $\mathbb{E}_{p(\boldsymbol{\theta})}[p(\mathbf{x} | \boldsymbol{\theta})] = \int_{\Theta} p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$. The latter definition underlines the useful interpretation of the evidence as an *average* over all possible $\boldsymbol{\theta}$ or, equivalently, as a *marginal likelihood*. Even though the marginal likelihood is usually bypassed in parameter estimation tasks due to the obvious proportionality

$$p(\boldsymbol{\theta} | \mathbf{x}) \propto p(\mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta}), \quad (3.2)$$

it becomes a key object in the context of Bayesian model comparison.

Importantly, all distributions in Equation 3.1 are implicitly conditioned on a given generative model g , in addition to unspecified contextual information C . Looking at the computational definition of the posterior, another consequence of Bayesian inference stands out, namely, that different observations will lead to different updated epistemic states. Thus, observing \mathbf{x}' will generally elicit a different change in knowledge, that is, $p(\boldsymbol{\theta} | \mathbf{x}) \neq p(\boldsymbol{\theta} | \mathbf{x}')$, and, again, a potentially different conclusion based on that modified knowledge³. This property seems desirable, as

²The influence of different prior specifications on model-based inference and model-derived decisions can be systematically investigated and quantified via *prior sensitivity analysis* [7]. More on this in Chapter 7.

³However, if the observations are equally informative or equally uninformative, the two posteriors may be the same. In general, the information gained by observing different states of the world will be different.

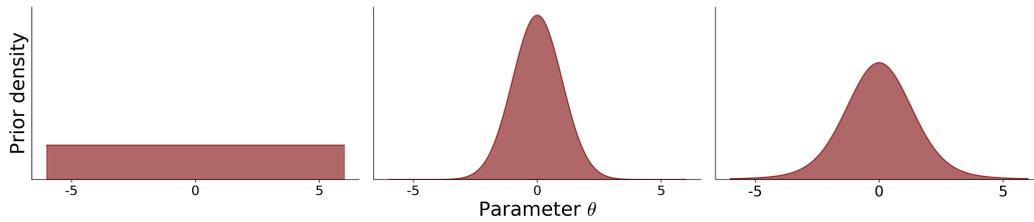


Figure 3.1: Three different ways to express prior knowledge and uncertainty about a single parameter θ in a probabilistic way. The left panel depicts a uniform prior around 0 which assigns constant density to values within the interval $[-6, 6]$ and treats all values outside this interval as literally impossible. The middle and right panel depict symmetric stable priors [179] around 0 which assign highest density to the central value of 0 and exponentially decaying density to values diverging from the center. Note, that the distribution in the right panel appears *shorter* than its *normal* counterpart in the middle panel, since it has thicker tails and treats values farther from 0 as less improbable.

it captures the (rather trivial) intuition that two learners presented with different facts will learn different things.

A further desirable property of Bayesian inference is that it allows the simultaneous or sequential integration of all available information (i.e., yesterday’s posteriors become today’s priors). Provided that information does not decay between sequential updates, both simultaneous and sequential updating should lead to the same endpoint posterior. The latter method is especially useful when data are collected at different points in time or arrive as a stream of observations. Curiously, even though often cited as a crucial asset of Bayesian inference, sequential updating is largely underutilized in behavioral research⁴.

3.2 UNCERTAINTY REDUCTION AND BAYESIAN SURPRISE

The transition from prior to posterior essentially conveys a *reduction in uncertainty* brought about by observing reality. Equivalently, it can be seen as communicating the *gain in information* achieved by querying nature through an empirical endeavor. Accordingly, we expect the posterior to be *narrower* or *sharper* than the prior, as the opposite would imply a loss of information through observation - a scenario which appears rather paradoxical. However, it is reasonable to expect circumstances when the posterior will exactly equal the prior, namely, whenever the data carry no information about the parameters of interest.

The concept of *posterior contraction* formalizes the idea that the posterior should get sharper as the number N of available observations increases. In the simplest case, the posterior variance should decrease at rate $1/N$, but a more nuanced behavior can occur for more complex models inducing, for example, multi-modal or asymmetric posteriors. In general, the posterior contraction (PC) for multivariate posteriors can be formally expressed as a ratio between generalized variances:

⁴And for good reasons, since updating turns out to be technically hard if the posterior is represented only as random draws (e.g., as provided by MCMC algorithms).

$$\text{PC}[\boldsymbol{\theta} | \boldsymbol{x}] = 1 - \frac{\det(\text{Cov}[\boldsymbol{\theta} | \boldsymbol{x}])}{\det(\text{Cov}[\boldsymbol{\theta}])} \quad (3.3)$$

where $\text{Cov}[\boldsymbol{\theta} | \boldsymbol{x}]$ and $\text{Cov}[\boldsymbol{\theta}]$ denote the posterior and prior covariance matrices, respectively. Posterior contraction near zero indicates that the data contribute little to nothing to reducing prior uncertainty about $\boldsymbol{\theta}$. Posterior contraction near one indicates a large reduction in prior uncertainty after accounting for the data [144]. Note, that Equation 3.3 provides *local* (i.e., per-observation) information about information gain. If one is interested about *global* (i.e., in expectation over all possible observations) information gain for a particular model, then the expected posterior contraction (EPC) should be considered:

$$\text{EPC}[\boldsymbol{\theta} | \boldsymbol{x}] = 1 - \mathbb{E}_{p(\boldsymbol{x})} \left[\frac{\det(\text{Cov}[\boldsymbol{\theta} | \boldsymbol{x}])}{\det(\text{Cov}[\boldsymbol{\theta}])} \right] \quad (3.4)$$

Accordingly, different model parameterizations can be compared with respect to their expected information gain by computing the corresponding EPCs. However, the EPC will usually be intractable for complex models, so it needs to be approximated via its empirical mean over a sufficiently large number of different observations. For most non-trivial models, this will only be feasible in an *amortized inference* setting (to be discussed in the next chapter) which makes inference globally efficient.

Notably, posterior contraction only considers the variance, that is, the second moments of the prior and posterior distributions. However, other differences between prior and posterior may manifest themselves in differences between higher-moments of the distributions. A concept for quantifying arbitrary differences between prior and posterior is the *Bayesian surprise*, which can be defined as the Kullback-Leibler (KL) divergence between the two densities:

$$\mathcal{B} = \mathbb{KL}[p(\boldsymbol{\theta} | \boldsymbol{x}) || p(\boldsymbol{\theta})] \quad (3.5)$$

$$= \int_{\Theta} p(\boldsymbol{\theta} | \boldsymbol{x}) \log \left(\frac{p(\boldsymbol{\theta} | \boldsymbol{x})}{p(\boldsymbol{\theta})} \right) d\boldsymbol{\theta} \quad (3.6)$$

Importantly, the Bayesian surprise is non-negative and equals zero if and only if the two densities are equal. In information theory, this quantity is termed the *relative entropy*, which, in a Bayesian context, represents the information gained by replacing the prior with the posterior. The units of information are then determined by the base of the logarithm. The *expected Bayesian surprise* (EBS) then encodes the average information gained from applying a particular model to the entire data space. If evaluation of the prior and posterior densities is analytically tractable, the empirical approximation of the EBS can therefore be used to quantify the impact of the data on Bayesian updating⁵.

Global information gain and uncertainty reduction are especially useful in the model development phase, since they can inform researchers about the general utility of a computational model. Accordingly, a researcher may experiment with different theoretically plausible parameterizations of a model and prune those which result in a poor information gain. Such an approach involves

⁵In cases where even approximating the (expected) Bayesian surprise is infeasible, an integral metric such as the maximum mean discrepancy (MMD, [61]) can be used to define Bayesian surprise as well.

performing repeated simulations from each model and then obtaining a posterior for each single simulation and model⁶. Notably, this approach can quickly become practically infeasible even for a few simulations if estimating the posterior is computationally demanding. Again, we will see in later sections why efficient amortized inference is particularly helpful for such undertakings.

3.3 TYPES OF UNCERTAINTY

The notion of uncertainty is central to scientific and daily life. However, in most research contexts dealing with probabilistic matters, it is often useful to introduce a *taxonomy of uncertainties*. A canonical approach in predictive modeling draws a distinction between *aleatoric* (sometimes called *ontic*) and *epistemic uncertainty* [74, 82].

Broadly speaking, aleatoric uncertainty is caused by intrinsic randomness, whereas epistemic uncertainty is brought about by ignorance. To illustrate this distinction with a (rather trivial) example [74], consider a coin-flipping game taking place in Bulgaria. A probabilistic model of the coin-flipping process could inform us about the likelihood of obtaining head or tail on any given toss. Further, assuming our probabilistic model is calibrated to reality, we can derive from it the expected value of the coin-flipping game and use decision theory to select certain actions (i.e., whether to bet my watch on the next toss). However, our model cannot foretell the *concrete* future outcome due to aleatoric reasons (matters of statistical physics aside), so there is uncertainty regarding the future possession of my watch. Differently, one might be equally uncertain about the meaning of the word “ezi” in Bulgarian. Yet the possible answers (and corresponding probabilities) are the same as before: head or tail. In this case, the ensuing uncertainty is caused purely by our lack of linguistic knowledge (i.e., epistemic uncertainty).

It has been argued that distinguishing between aleatoric and epistemic uncertainty constitutes a so-called *distinction without a difference* [154], that is, one without any practical consequences for decision makers. However, since ignorance (as opposed to intrinsic randomness) can in principle be reduced or even removed with additional information, it appears that pinpointing the source of uncertainty implies different handling of uncertainty. Thus, in our coin-flipping game, one can easily get rid of the epistemic uncertainty surrounding the word “ezi” by means of a dictionary. No need to set up a probabilistic model. Accordingly, epistemic uncertainty is generally treated as reducible; aleatoric uncertainty is generally treated as irreducible.

Note, however, that the irreducibility of aleatoric uncertainty is, in most cases, rather a practical decision than an ontological necessity. For instance, even in the coin-flipping game (the all-time favorite example of probability textbooks), knowledge of the initial conditions and of all forces acting on the coin at all times would, in principle, render the system deterministic. To further illustrate this point, consider the following (deterministic) linear model with two covariates, x_1 , x_2 , and one outcome y :

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} \quad (3.7)$$

⁶Despite being self-consistent, such an approach does not guarantee the utility of a model when applied to real data. A potential *simulation-gap* between simulation and reality can be detected in later modeling stages with different tools.

3 Bayesian Inference

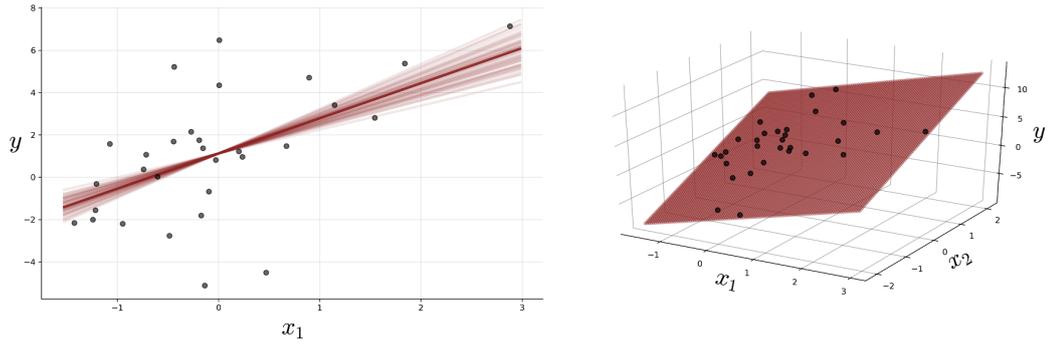


Figure 3.2: An illustrative example of the contextual irreducibility of aleatoric uncertainty. Left panel: having limited access to the relevant data (only x_1) induces aleatoric variability which cannot be reduced with further observations of the same kind. The deterministic model appears stochastic. Right panel: having access to the full data (both x_1 and x_2) makes the uncertainty disappear, since the true model is deterministic.

with all $x \sim \mathcal{N}(0, 1)$. If one were to fit a proverbial simple linear regression on a reduced data set $\mathbf{D} = \{x_{1i}, y_i\}_{i=1}^{N=30}$ generated from the model, with x_2 treated as unknown to the modeler, the data would behave as if they had been generated via the well-known:

$$y_i = \beta_0 + \beta_1 x_{1i} + \xi_i \quad (3.8)$$

with $\xi \sim \mathcal{N}(0, \sigma)$, where σ is aleatoric Gaussian noise dependent on the unknown β_2 and x_2 .

The data set and the resulting best-line fit are depicted in the left panel of Figure 3.2. In this simple regression scenario, epistemic and aleatoric uncertainty are seemingly separable. Due to the finite N , there is epistemic uncertainty about the precise values of β_0 and β_1 , which might change if a different set of 30 observations was generated from the model. In a Bayesian setting this uncertainty would be captured via the posterior distribution $p(\beta_0, \beta_1 | \mathbf{D})$, which will get narrower as N increases. Note also, that the epistemic uncertainty in β_0 and β_1 results in uncertainty in the best-line fit, as depicted by the multiple shaded lines. However, regardless of how large N gets, the aleatoric factor ξ will cause the points to vary unsystematically around the line, and, correspondingly, there will be irreducible predictive uncertainty about the true outcome of an upcoming observation $x_{1i'}$. On the other hand, if we could augment our original data set with (the previously hidden) x_2 , that is, if we use $\mathbf{D}' = \{x_{1i}, x_{2i}, y_i\}_{i=1}^{N=30}$ in a multiple regression model, all previous aleatoric uncertainty disappears, since the additional information of x_2 renders all points to perfectly lie on the best-fitting plane (cf. Figure 3.2, right panel). Thus, what we ultimately treat as irreducible uncertainty might depend on the particular modeling context instead of referring to an absolute notion.

When performing Bayesian inverse inference with Monte Carlo simulators, we typically want to mimic the real-world randomness, regardless of its source, via the stochastic component $p(\boldsymbol{\xi})$ involved in forward inference:

$$\mathbf{x}_n = g(\boldsymbol{\theta}, \boldsymbol{\xi}_n) \quad \text{with} \quad \boldsymbol{\xi}_n \sim p(\boldsymbol{\xi}) \quad (3.9)$$

Since we are dependent on computer-based random number generation, we typically need to assume a parametric form for $p(\xi)$ in order to make forward inference tractable in the first place. Even though in some applications (mostly experiments in physics), the precise source and form of the data-corruption process might be known in advance, most cognitive models are confined to a convenient form of uncertainty (i.e., one which we can easily simulate or write down mathematically). Thus, there is uncertainty about the chosen form of randomness, which, along with the uncertainty about the generative capabilities of the simulator, would fall under the rubric of *model uncertainty*, generally considered as another (sub-)type of epistemic uncertainty [74]. However, since the true data generator in areas reserved for the behavioral sciences is almost never transparent, model uncertainty is almost always practically irreducible. Accordingly, the reducible uncertainty encoded by the Bayesian posterior $p(\theta | \mathbf{x}_{1:N})$ is immersed into an ocean of potentially irreducible uncertainty and the sought-after Bayesian surprise is only guaranteed when the observed data are informative for the target parameters θ .

3.4 EXCHANGEABLE OBSERVATIONS

In the previous chapter, we briefly touched upon the idea of a memoryless probabilistic program. Running a memoryless simulation means that each run of the simulator is not affected by any of the previous runs and will, in turn, not affect any of the future runs. But on what grounds can we assume that such models can even remotely do justice to the noisy vicissitude of reality?

The answer is simplicity. Consider an observed sequence of N data points $\mathbf{x}_{1:N}$, that is, a data set obtained from an experiment or in an observational study, waiting to be analyzed by an eager scientist. Without a model in mind, the data set can be assumed to arise from an unknown random process $p^*(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, also known to statisticians as the *data-generating* process. Generally, it will be beyond our intellectual reach to provide a complete description of how even a single observation in a given behavioral data set has come about⁷. Thus, we often need to leverage some probabilistic symmetry imposed on p^* which renders the data describable with compact and interpretable models. One such symmetry, tirelessly assumed in the Bayesian modeling literature, is *exchangeability*.

To illustrate exchangeability, imagine a scenario, in which researcher A, for whatever reasons, conspired to shuffle the observations in a data set collected by their colleague B. Under the assumption of exchangeability, the mischievous researcher A would be wasting their time - the ordering of the data points does not matter at all to researcher B. Formally speaking, researcher B views the sequence of data points as *invariant under all permutations of the data points* and plans to build a memoryless model of the data conforming to this view.

Exchangeable observations may come in various forms, for instance, patients entering a hospital, psychology students participating in a response time experiment, or raw response times recorded from a single participant. Exchangeable observations are not only conceptually simple to work with, but also admit a particularly useful probabilistic treatment. According to de Finetti's

⁷Just try to think of all possible factors affecting a single response time of a drunk participant in a reaction time study.

representation theorem [33, 116], an exchangeable data-generating distribution p^* has the following integral decomposition:

$$p^*(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \int_{\Upsilon} p^*(\mathbf{v}) \prod_{n=1}^N p^*(\mathbf{x}_n | \mathbf{v}) d\mathbf{v} \quad (3.10)$$

with $\mathbf{v} \in \Upsilon$ being an (potentially infinitely dimensional) absolutely continuous random vector. Even though the representation theorem does not provide any clues on how to actually find the correct integral decomposition, it motivates the formulation of Bayesian models of the form

$$p(\boldsymbol{\theta}, \mathbf{x}_1, \dots, \mathbf{x}_N) = p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta}) \quad (3.11)$$

as useful approximations of the unknown “parameter vector” \mathbf{v} and densities $p^*(\mathbf{v})$ and $p^*(\mathbf{x} | \mathbf{v})$. Note also, that Equation 3.11 represents the numerator of Bayes’ rule for multiple *i.i.d.* observations. Essentially, this equation represents a statement of *conditional independence* and specifies a generative recipe for simulating synthetic observations: first, obtain a random sample from the prior $p(\boldsymbol{\theta})$ and run the simulator N times with the corresponding sample to obtain a simulated tuple $(\boldsymbol{\theta}, \mathbf{x}_1, \dots, \mathbf{x}_N)$.

Exchangeable observations impose a symmetry on the posterior as well. As a consequence of the assumption, the resulting posterior $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ should also be invariant with respect to the ordering of the individual observations \mathbf{x}_n . In other words, if the function $\mathbb{S}_N(\mathbf{x}_{1:N}) = (\mathbf{x}_{\pi(1)}, \dots, \mathbf{x}_{\pi(N)})$ represents an arbitrary permutation of N elements, the following should hold for the posterior:

$$p(\boldsymbol{\theta} | \mathbf{x}_{1:N}) = p(\boldsymbol{\theta} | \mathbb{S}_N(\mathbf{x}_{1:N})) \quad (3.12)$$

The same applies to any (sufficient) summary statistic of the data $h(\mathbf{x}_{1:N})$:

$$p(\boldsymbol{\theta} | h(\mathbf{x}_{1:N})) = p(\boldsymbol{\theta} | h(\mathbb{S}_N(\mathbf{x}_{1:N}))) \quad (3.13)$$

and, by extension, to any estimator of the posterior which is a function of the raw data. In later chapters, we will elaborate on why the preservation of posterior symmetry poses a challenge to standard neural network estimators. We will further discuss and describe how to employ specialized symmetry-preserving networks for addressing this challenge [13]. Albeit common, memoryless models are not ubiquitous throughout the behavioral sciences. Thus, our inference frameworks will incorporate different choices of neural architectures for stateful models and non-exchangeable distributions.

3.5 BAYESIAN MODEL COMPARISON

Our discussion on Bayesian inference so far has concentrated on a single abstract model. In fact, we have even treated the model as an invisible part of the unspecified context C implicit to all distributions involved in the computation of Bayes’ rule (Equation 3.1). However, this simplified setting is rather rare in the behavioral sciences. In most model-rich research areas, such as decision

making, attention, judgement formation, risk-taking, and memory, there is a multitude of models and model parameterizations competing to account for a given behavior. Moreover, in a Bayesian context, different models can be defined not only via different generative mechanisms, but also by different priors (e.g., as in prior sensitivity analysis, [163]), or by different data processing pipelines (e.g., as in multiverse analysis, [149]). Thus, researchers often find themselves in a scenario calling for formal model comparison, model selection, or model averaging.

In order to extend our discussion and formal notation to multiple models, consider a collection of J candidate models $\mathcal{M} = \{\mathcal{M}_1, \dots, \mathcal{M}_J\}$ and corresponding parameter spaces $\mathcal{H} = \{\Theta_1, \dots, \Theta_J\}$. We assume that each model is realizable via a generative algorithm and a simulation program, such that $g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi}_j)$ realizes \mathcal{M}_j for each j . Thus, in order to perform Bayesian parameter estimation for each model given an observation \boldsymbol{x} , we need to compute J posteriors of the form

$$p(\boldsymbol{\theta}_j | \boldsymbol{x}, \mathcal{M}_j) = \frac{p(\boldsymbol{x} | \boldsymbol{\theta}_j, \mathcal{M}_j) p(\boldsymbol{\theta}_j | \mathcal{M}_j)}{p(\boldsymbol{x} | \mathcal{M}_j)} \quad (3.14)$$

with the prior, posterior, likelihood, and marginal likelihood explicitly written as dependent on the particular model \mathcal{M}_j .

How does one assign preferences to competing models using a probabilistic toolkit? Assuming that all quantities in Equation 3.14 are tractably computable, Bayesian methods for model selection revolve around two key concepts: *prior predictive measures* and *posterior predictive measures* [52]. Prior predictive and posterior predictive approaches to model comparison answer somewhat different questions, so asking “which one is better” for a specific modeling problem is rarely expedient. Moreover, their answers can occasionally diverge, so oftentimes, it is informative to explore *both* approaches in order to obtain a more nuanced picture of the candidate models’ performance. The most commonly faced obstacle in practice is feasibility: both approaches are computationally expensive (sometimes intractable) for complex models and can thus benefit from our proposed frameworks. We now briefly discuss each approach.

3.5.1 THE PRIOR PREDICTIVE AND OCCAM’S RAZOR

The canonical measure of prior predictive performance is the already encountered *evidence*, which is the denominator in Bayes’ rule:

$$p(\boldsymbol{x} | \mathcal{M}_j) = \int_{\Theta_j} p(\boldsymbol{x} | \boldsymbol{\theta}_j, \mathcal{M}_j) p(\boldsymbol{\theta}_j | \mathcal{M}_j) d\boldsymbol{\theta}_j \quad (3.15)$$

Note, that the integral runs over the prior space of each model \mathcal{M}_j and thus represents the expected likelihood with respect to each model’s prior. The evidence thus penalizes prior complexity, that is, the inelegance, of a model, since the prior acts as a weight on the likelihood. It also induces a well-known and widely appreciated source of intractability in Bayesian inference, since it typically involves a multi-dimensional integral over potentially unbounded parameter spaces. For most non-trivial models, this integral cannot be computed in closed-form or approximated numerically. Sophisticated algorithms for efficiently approximating the evidence have been proposed in the Bayesian universe, such as *bridge sampling* and *path sampling* [54, 62]. However, these methods still depend on the ability to evaluate the likelihood $p(\boldsymbol{x} | \boldsymbol{\theta}_j, \mathcal{M}_j)$ for each candi-

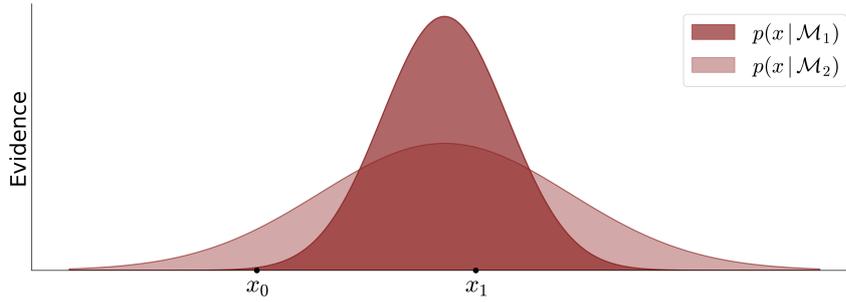


Figure 3.3: An illustrative example of Bayesian Occam's razor. The figure depicts the hypothetical evidence (marginal likelihood) of a simple model \mathcal{M}_1 vs. that of a more complex model \mathcal{M}_2 . The complex model has a larger generative scope and thus accounts for a broader range of observations by spreading its marginal likelihood to cover the whole range. It does so at the cost of diminished sharpness. Even though observation x_1 is well within its generative scope, the simpler model \mathcal{M}_1 yields a higher evidence and is therefore favored. In contrast, observation x_0 has a higher evidence under model \mathcal{M}_2 , as it is very unlikely to be generated by the simpler model.

date model. If the likelihood itself is intractable, as is the case with complex simulators, the task becomes increasingly hopeless even with the most efficient current methods.

Provided that the marginal likelihood can be reliably approximated, one can compute the ratio of marginal likelihoods for two models \mathcal{M}_j and \mathcal{M}_k

$$\text{BF}_{jk} = \frac{p(\mathbf{x} | \mathcal{M}_j)}{p(\mathbf{x} | \mathcal{M}_k)} \quad (3.16)$$

This famous ratio is called a Bayes factor (BF) and is used in Bayesian settings for quantifying relative model preference. Thus, a $\text{BF} > 1$ indicates preference for model j over model k , taking observation \mathbf{x} into account. Alternatively, one can directly focus on the (marginal) posterior probability of a model \mathcal{M}_j

$$p(\mathcal{M}_j | \mathbf{x}) \propto p(\mathbf{x} | \mathcal{M}_j) p(\mathcal{M}_j) \quad (3.17)$$

which equips the model space itself with a multinomial prior distribution $p(\mathcal{M})$ encoding potential prior beliefs on the plausibility of each model before collecting any data. Such a prior might be useful if a model embodies extraordinary claims (e.g., psychokinesis) and thus requires extraordinary evidence supporting it. However, if no prior reasons can be given for favoring some models over others (i.e., one prefers not to prefer), a uniform model prior $p(\mathcal{M}_j) = 1/J$ can be assigned.

The ratio of posterior model probabilities is called the *posterior odds* and is connected to the Bayes factor via the corresponding model priors:

$$\frac{p(\mathcal{M}_j | \mathbf{x})}{p(\mathcal{M}_k | \mathbf{x})} = \frac{p(\mathbf{x} | \mathcal{M}_j)}{p(\mathbf{x} | \mathcal{M}_k)} \times \frac{p(\mathcal{M}_j)}{p(\mathcal{M}_k)} \quad (3.18)$$

If both models are deemed equally likely *a priori*, that is, $p(\mathcal{M}_j) = p(\mathcal{M}_k)$, the posterior odds are identical to the Bayes factor. In this case, if the Bayes factor, or, equivalently, the posterior odds equal one, the observed data provide no evidence for favoring one of the models over the other. Importantly, a relative evidence of one does not distinguish whether the data are equally likely or equally unlikely under both models, as this is a question of absolute evidence. It simply means that the observations are not informative to the question of model comparison⁸.

It follows from our discussion so far, that, from a prior predictive perspective, model complexity is determined by (i) the prior of a model and (ii) the likelihood function of a model. Together, these two quantities establish the *generative scope* of a model and dictate how to select between two hypothetical cognitive models which both provide a reasonable account of some behavioral manifestations.

Thus, a complex model g_j with a large generative scope can generate a larger variety of behaviors, and so the density of $p(\mathbf{x} | \mathcal{M}_j)$ must spread over a larger portion of the observation space \mathcal{X} . On the other hand, a simpler model g_k with a smaller generative score can generate a limited range of behaviors, so its density will be restricted to a smaller portion of \mathcal{X} , and thus more peaked (cf. Figure 3.3). This is the reason why a marginal likelihood is said to automatically embody a fundamental trade-off between a model’s complexity (an antonym of scientific elegance) and its ability to convincingly account for a wide variety of empirical phenomena. Indeed, this trade-off is sometimes embraced as a probabilistic version of the famous Occam’s razor.

From a generative perspective, simpler simulators will tend to produce synthetic observations which are more similar to each other compared to those generated via more complex simulators. Indeed, this is the most important property of Bayesian models that we will later leverage in order to perform efficient model comparison between complex cognitive simulators. Essentially, by approximating Equation 3.17 directly, we will be circumventing two common sources of intractability: the marginal likelihood and the likelihood function itself.

3.5.2 THE POSTERIOR PREDICTIVE AND FORTUNA’S KNIFE

Prior predictive measures such as the Bayes factor do not utilize the posterior when comparing models (cf. Equation 3.15). In contrast, posterior predictive methods quantify the ability of models to forecast *new observations* which have not been used for Bayesian updating. In other words, if the posterior of each model represents a modified state of knowledge upon integrating observation \mathbf{x} , one tests which form of modified knowledge can “best” predict an unseen observation \mathbf{x}' . The most straightforward way to perform such an operation consists in computing the *posterior predictive* distribution:

$$p(\mathbf{x}' | \mathbf{x}, \mathcal{M}_j) = \int_{\Theta_j} p(\mathbf{x}' | \boldsymbol{\theta}_j, \mathbf{x}, \mathcal{M}_j) p(\boldsymbol{\theta}_j | \mathbf{x}, \mathcal{M}_j) d\boldsymbol{\theta}_j \quad (3.19)$$

where the likelihood $p(\mathbf{x}' | \boldsymbol{\theta}_j, \mathbf{x}, \mathcal{M}_j)$ simplifies to $p(\mathbf{x}' | \boldsymbol{\theta}_j, \mathcal{M}_j)$ when working with memoryless models. The main difference now is that integration is performed with respect to the posterior distribution of each model. Intuitively, the posterior predictive uses all information gained

⁸Unless one assumes that the model set \mathcal{M} provides an exhaustive description of reality, in which case relative evidence is absolute evidence.

through previous observation(s) to provide an uncertainty-aware forecast for the new query \mathbf{x}' , where the parameter posterior $p(\boldsymbol{\theta}_j | \mathbf{x}, \mathcal{M}_j)$ acts as a weight on the likelihood. Thus, the probability of an upcoming observation under a model \mathcal{M}_j is a weighted average over its probabilities under all plausible parameter configurations $\boldsymbol{\theta}_j$. Thereby, epistemic uncertainty, as captured by the parameter posterior, is essentially “averaged out” in the posterior predictive.

The canonical approach for Bayesian posterior predictive comparisons are cross-validation (CV) methods [166]. Examples for widely applied methods that fall into this category are approximate cross-validation methods using Pareto-smoothed importance sampling (PSIS-CV) [17, 165], information criterion measures, such as the widely applicable information criterion (WAIC; [171]), or stacking of posterior predictive distributions [176]. All of these methods require not only the ability to evaluate the likelihood of each model for each observation during parameter estimation, but also for new observations during prediction.

What is more, if application of exact CV methods is required because approximations are insufficient or unavailable, models need to be estimated several times based on different data sets or subsets of the original data set. This renders such methods practically infeasible when working with complex simulators for which posterior inference is already computationally demanding. Thus, even a single intractable model in the candidate model set suffices to disproportionately increase the difficulty of performing posterior predictive model comparison.

Posterior predictive measures based on evaluating the likelihood of new data points provide information about the *relative* performance of models. For certain applications, one can replace the likelihood in Equation 3.19 with a scoring function $S : \mathcal{X} \rightarrow \mathbb{R}$ and arrive at a measure of model predictive performance with respect to the scoring function:

$$s_j = \int_{\Theta_j} \int_{\Xi_j} S(\mathbf{x}', g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi}_j)) p(\boldsymbol{\xi}_j | \boldsymbol{\theta}_j) p(\boldsymbol{\theta}_j | \mathbf{x}, \mathcal{M}_j) d\boldsymbol{\xi}_j d\boldsymbol{\theta}_j \quad (3.20)$$

where the integrals are typically approximated via Monte Carlo samples from the corresponding posterior and noise distributions. If the scoring function is well aligned with the particular goals of an inference task, it can serve as a useful proxy for quantifying both absolute and relative performance of the models under scrutiny. Moreover, it replaces the dependence on the likelihood, at least during prediction. Still, the estimation of each posterior and its repeated re-estimation for different data sets when using variants of exact LOO-CV remain as potential bottlenecks. These bottlenecks become even more challenging to overcome when doing simulation-based inference even with a single model of interest.

Note, that oftentimes researchers use the term *predict* when they are actually referring to a model’s ability to *reproduce* the data used to inform the estimation of model parameters [177]. Formally, this amounts to replacing the new observation \mathbf{x}' in the previous two equations with the observation \mathbf{x} (or set of observations $\mathbf{x}_{1:N}$) used to condition the posterior. Such a procedure is indeed useful, since it can be indicative of a model’s *generative performance* [120] and can thus help in diagnosing model misspecification or a simulation gap. However, generative performance is generally *not indicative* of a model’s predictive performance [53, 177], so we find it important to highlight and keep the distinction as clear and explicit as possible.

At this point, there is an important distinction to be made when it comes to predicting *new* observations. Suppose that each observation in a data set was generated by the unknown process

p^* of which we formulate a parametric model p_θ . If we now make the model “blind” to certain observations in the original data set (as in CV) and use these observations to assess predictive performance, we are essentially testing the model’s ability to perform *induction* about the statistical regularities of the process. In such a scenario, however, we are not assessing the model’s ability to faithfully forecast the future, since the observations are new only from the relative perspective of model fitting. An attempt to forecast future behavior with a static cognitive model will only be meaningful, if the process p^* is stationary (i.e., its regularities are invariant with respect to time) or if the model somehow explicitly incorporates a potential time-dependent change inherent to the generator. Since such changes are extremely hard to know in advance (otherwise they would have been predicted and incorporated into the model’s equations), a model which claims time-invariant performance should regularly be subjected to the falsification of time.

Finally, note that the probabilistic Occam’s razor from prior predictive approaches does not automatically show up in posterior predictive approaches. Differently, from a posterior predictive perspective, a model’s quality should be judged based on how well it *generalizes* to unseen instances it is supposed to accurately predict. In other words, a model has to withstand the sharp challenges of its future, and models which fail to do so, are discarded. However, one might still anticipate that overly complex models would *overfit* the data at hand and fail dramatically at predicting new data, again being implicitly subjected to some form of Occam’s razor. Indeed, such an anticipation has been formally framed under the so called *bias-variance dilemma* [47] which bounds the generalization error of supervised learning algorithms. However, the practical consequences of such a dilemma have recently been called into question by the achievements of “black-box” neural networks having billions of parameters and still being able to generalize beyond their training data [111]. It remains therefore unclear to what extent Occam’s razor is implicitly encoded in posterior predictive measures when applied to “white-box” models of cognition.

3.5.3 BAYESIAN MODEL AVERAGING AND WISDOM OF THE CROWD

A rather disparate approach to model selection is that of *not selecting* a single model but instead *averaging* across the predictions of all candidate models. However, in a Bayesian model averaging (BMA) setting, models are not created equal, and thus not weighted with indifference to their performance or elegance. Instead, posterior model probabilities are used as weights in the aggregate prediction:

$$p(\mathbf{x}' | \mathbf{x}) = \sum_{j=1}^J p(\mathbf{x}' | \mathbf{x}, \mathcal{M}_j) p(\mathcal{M}_j | \mathbf{x}) \quad (3.21)$$

Crucially, BMA depends on the marginal likelihoods of the different models, as well as on their likelihood functions. Naturally, BMA can also be used for averaging over scoring functions instead of posterior predictive distributions, when a proxy of absolute performance is needed. BMA is particularly useful when predictive performance matters and predictions of a single model are expected to be unstable. In fact, model averaging generally yields superior predictive results in expectation compared to those obtained by model selection [126]. Broadly speaking, it can be regarded as the Bayesian embodiment of the well-known *wisdom of the crowd*, or *vox populi*, statistical phenomenon [50].

Clearly, BMA is of limited use when the goal of inference is to compare competing theories instantiated by formal models and subsequently choose the most plausible among all (in a probabilistic sense). However, BMA might still come in handy for selecting between *model classes*, where the performance of single model instances is not crucial or when competing models might admit different plausible parameterizations.

Consider, for example, a class of mathematical models $\mathcal{A} = \{\mathcal{M}_j^A\}_{j=1}^J$, having property A in common, and another class of models $\mathcal{B} = \{\mathcal{M}_j^B\}_{j=1}^J$, having property B in common. Applying Equation 3.21 to the models in each class, we can obtain two model-averaged posterior predictive distributions $p(\mathbf{x}' | \mathbf{x}, \mathcal{A})$ and $p(\mathbf{x}' | \mathbf{x}, \mathcal{B})$. Subsequently, we can use these to quantify and compare the bulk predictive performance of the two model classes.

Performing BMA for selecting between model classes featuring complex models (e.g., memoryless vs. stateful models) is, however, even more computationally demanding, which makes it a highly underutilized approach in the behavioral sciences. It also presents a challenge we consider worthy of future investigation in the context of frameworks for simulation-based Bayesian inference.

3.6 BAYESIAN SIMULATION-BASED INFERENCE

We have seen that all central objects in Bayesian inference, from the Bayesian parameter posterior to the Bayesian model-averaged predictive distribution, depend on the likelihood function. Thus, we reiterate, when the likelihood cannot be efficiently evaluated or is not available in closed-form, standard Bayesian methods relying on the central proportionality $p(\boldsymbol{\theta} | \mathbf{x}) \propto p(\mathbf{x} | \boldsymbol{\theta})p(\boldsymbol{\theta})$ do not apply. This includes both the less efficient but asymptotically appealing MCMC and the more efficient but asymptotically less winsome variational methods. Moreover, a potential intractable likelihood can also appear on top of the well-known intractability of the *marginal likelihood*, which is necessary for model comparison in a prior predictive context. In fact, this is not just a theoretical hardship, but a real practical predicament faced by researchers working with various complex models [27]. Complexity is usually a direct consequence of the desire to build high-fidelity models of cognitive processes, sometimes coupled with a modicum of neurophysiological realism.

In such cases, not all is lost. Instead of simplifying the model *ad hoc*, one can still retain the advantages (and disadvantages) of Bayesian inference by “simply” performing repeated simulations and using them to guide the process of inference. Such an approach is known as *simulation-based inference*. When the resulting quantities still represent the reduction of prior uncertainty by conditioning on data, we are essentially doing the same Bayesian inference as before, only with the likelihood implicitly involved in the process. All Bayesian simulation-based methods repeatedly go through the following three steps, given informally by:

1. Obtain a random sample from the prior.
2. Simulate an artificial data set with the sampled parameters.
3. Do something with the simulated pair of parameters and data.

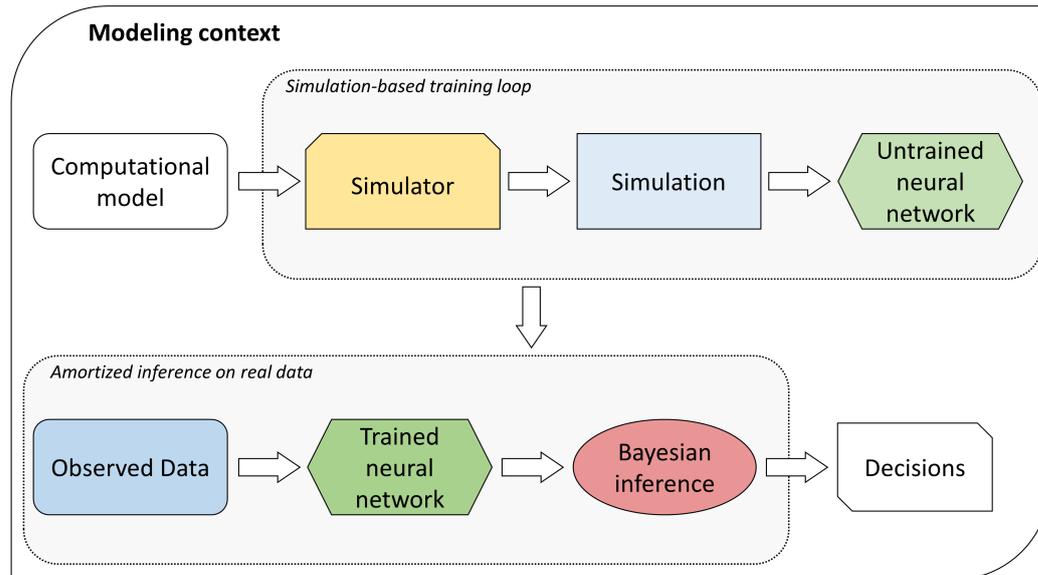


Figure 3.4: An abstract overview of the central idea underlying our proposed solutions to Bayesian inverse inference.

The essential difference between most Bayesian simulation-based methods lies in the particular implementation of step number three. Our frameworks are no exception to this pattern. A common theme will be to use the cognitive model as a trainer for a specialized neural network, which is driven through a number of simulations towards the Bayesian answer to an estimation or a model comparison problem (cf. Figure 3.4). The next chapter will provide a brief and very broad survey of the related work on simulation-based inference. Before we move on, however, we shortly discuss the idea of *samplers*, as it is central to our developed methods and probabilistic modeling as a whole.

3.7 SAMPLERS AND NEURAL SAMPLERS

It is important to note, that the posterior distribution itself is hardly ever available as a known density function which can be analytically calculated. In the typical textbook cases where the likelihood belongs to a known family of probability distributions and the prior is chosen to be *conjugate* to the likelihood, then the posterior must belong to the same distribution family as the prior. However, such mathematical convenience has proven insufficient for addressing most unidealized real-world problems, as it leaves little room for flexible or high-fidelity modeling. Thus, researchers and statistical software developers have resorted to MCMC methods.

MCMC methods, such as the Metropolis-Hastings algorithm [66], Gibbs sampling [51], Hamiltonian Monte Carlo [112] or its extension to No-U-Turn sampling [71], approximate the posterior in the form of *random samples from the target posterior*. MCMC methods belong to a family of stateful algorithms which generate a sequence of correlated samples that converge in distribution to a target equilibrium distribution (i.e., the posterior, in a Bayesian setting).

The idea of approximating a complicated distribution via dependent random samples, albeit rather straightforward in hindsight, has gradually transformed and shaped the field of Bayesian inference. Moreover, it forms the main logic behind major probabilistic programming languages such as JAGS [127] or Stan [19], which return posterior estimates in the form of random samples. A *sampler* is thus a program which uses computer-generated randomness to “draw” samples from a distribution instead of deriving or estimating its algebraic form.

Recently, the idea of random sampling has entered the rapidly expanding field of deep learning under the umbrella term *deep generative modeling*. As a consequence, the concept of *neural samplers* has emerged [72, 87, 113, 114]. Neural samplers emulate sampling from a distribution via neural networks that transform a random input vector into a sample from a target probability distribution defined by the network weights. The random input vector is typically drawn from a simple distribution (e.g., uniform or Gaussian) which is computationally cheap and easy to sample from. The expressive transformation of simple inputs creates diversity and flexibility. Accordingly, the network weights defining the transformation are optimized in a way to ensure that the subsequently generated samples are actually representative of the target distribution.

Neural samplers have so far shown tremendous success in computer vision and natural language processing. Our inference frameworks also make extensive use of neural samplers. However, our neural networks are trained to approximate the Bayesian posterior induced by a particular model given a set of observations and thus generate posterior samples in a way similar to MCMC. Moreover, we will utilize the important fact that neural samplers distil global information about a particular model family into their trainable parameters (e.g., network weights). Thus, once trained, neural samplers are easy to store on a computer or a server (a couple of megabytes memory demand) and re-use across multiple applications of the same model family to many data sets of potentially variable size. This property gives rise to *amortized inference*, a concept we will repeatedly encounter in later chapters.

4 RELATED WORK

As with any new and rapidly expanding field, it is nearly impossible to manage a comprehensive review of the existing literature, since new methods would have emerged upon the review’s completion. Indeed, this is exactly what happened during the (rather expeditious) completion of this thesis. Thus, this chapter attempts to provide a cursory glance upon the landscape of simulation-based inference with a special focus on deep learning methods. For a more detailed review of recent developments, see, for example, [27]. For a collection of classical methods with a focus on cognitive science, see, for example, [119].

There are multiple ways to devise a taxonomy for members of the zoo of simulation-based inference. For the purpose of this thesis, we can classify methods as *amortized* vs. *non-amortized*, with different degrees of amortization possible (see Chapters 5 and 7). The main difference between the two methodological endpoints is this. Non-amortized methods require a repetition of the same computations for each model application, that is, estimation starts from scratch for each observed data set $\mathbf{x}_{1:N}^{(obs)}$. In contrast, amortized methods involve an upfront optimization/training phase which ensures that subsequent applications of the model to any observed data set are very cheap (i.e., the cost of the optimization phase *amortizes* over multiple inferences).

The standard non-amortized solution to intractable modeling problems is offered by *approximate Bayesian computation* (ABC) methods [29, 150]. ABC methods approximate the posterior by repeatedly sampling parameters from a proposal (prior) distribution $\boldsymbol{\theta}^{(l)} \sim p(\boldsymbol{\theta})$ and then simulating a synthetic data set by running forward inference, $\mathbf{x}_n \sim p(\mathbf{x} | \boldsymbol{\theta}^{(l)})$ for $n = 1, \dots, N$, with the sampled parameters. If the simulated data set is sufficiently similar to an actually observed data set (as measured by a user-defined distance function), the corresponding parameter configuration $\boldsymbol{\theta}^{(l)}$ is retained as a random draw from the desired posterior, otherwise rejected. However, in practice, ABC methods are notoriously inefficient and suffer from various problems, such as the *curse of dimensionality* or *curse of inefficiency* [104], to name the most severe.

More efficient and sophisticated methods, such as sequential Monte Carlo algorithms (ABC-SMC) or Markov chain Monte Carlo with implicit likelihoods (ABC-MCMC), employ different creative techniques to optimize sampling or correct potential biases [65, 90, 107, 119, 123, 157, 158, 160]. Currently, the gold-standard in cognitive science and mathematical psychology appears to be non-amortized ABC-MCMC with kernel density estimation (KDE) [44, 160, 161]. This method has the advantage of doing away with hand-crafted summary statistics and distance functions, since each (simulated or actual) observation \mathbf{x}_n enters the KDE computation at each MCMC step. However, the application of KDE-based MCMC to stateful models yielding non-*i.i.d.* observations is not at all straightforward, since KDE methods typically assume *i.i.d.* observations. Furthermore, every non-amortized method becomes infeasible in data-rich settings which require the estimation of the same model hundreds or even thousands of times (or even more than one million times, as was the case in [94]), unless a researcher has access to a high-end computing cluster with countless cores.

Recently, machine learning and deep learning innovations have permeated the field of simulation-based inference with the goal of scaling up standard ABC methods. Most of these innovations yield *amortized methods*, since they involve an expressive machine learning approximator (e.g., random forests or neural networks) trained on simulations from the Monte Carlo engine which emulates the behavior of the analytically intractable model. These approximators are either only able to return summaries of the full posterior (e.g., posterior means, variances, or quantiles) or capable of performing fully Bayesian inference. We now discuss each of the two approaches in turn.

Perhaps the most straightforward inference method has been to cast the problem of parameter estimation as a supervised regression task [14, 79, 136, 140]. In this setting, the simulator is run repeatedly to create a large dataset of the form $\mathcal{D} = \{(h(\mathbf{x}_{1:N}^{(m)}), \boldsymbol{\theta}^{(m)})\}_{m=1}^M$, also referred to as a *reference table* in the ABC literature [150]. Typically, the dimensionality of the simulated data is reduced by computing summary statistics with a fixed summary function $h(\mathbf{x}_{1:N})$ (but see [136]). Then, the reference table \mathcal{D} is used as training data for a supervised learning algorithm (e.g., random forest [140], or a convolutional neural network [136]). The learning algorithm is trained to output an estimate of the true data-generating parameters and an optional uncertainty estimate (e.g., the posterior variance [14, 136] or quantiles [140]). Thus, supervised methods attempt to approximate the intractable inverse model directly and globally via non-linear regression $\hat{\boldsymbol{\theta}} = f(h(\mathbf{x}_{1:N}))$. Importantly, the trained algorithms can be cheaply stored and re-used for estimation on an arbitrary number of observed data sets or integrated into an ABC routine [79].

A severe shortcoming of supervised approaches is that they provide only limited information about the full posterior or impose overly restrictive assumptions on its distributional form (e.g., Gaussian). This is especially problematic when the true posterior is acutely skewed or multimodal, in which case the mean or the variance are not particularly representative of its relevant characteristics.

To address this shortcoming, neural density estimation (NDE) methods employ specialized neural networks capable of performing fully Bayesian inference (i.e., returning full posteriors). NDE methods approximate different components of the intractable joint Bayesian model, that is, $p(\boldsymbol{\theta}, \mathbf{x}_{1:N}) = p(\boldsymbol{\theta})p(\mathbf{x}_{1:N} | \boldsymbol{\theta})$.

Sequential neural posterior estimation (SNPE) methods iteratively refine a proposal distribution via specialized neural networks (e.g., mixture density networks, autoregressive or normalizing flows) to generate parameter samples which closely match a particular observed data set [60, 101, 121]. Even though these methods also entail a relatively expensive learning phase and a cheap inference phase, they are capable of amortized inference only when operating in a non-sequential manner (i.e., the prior is used as a proposal throughout every optimization step). Otherwise, a separate neural density estimator has to be trained for each observed data set, which quickly becomes infeasible when working with many data sets. The main feature of SNPE methods is that they avoid MCMC sampling altogether and are able to sample from the *true posterior* given perfect convergence. This is in contrast to variational methods which optimize a lower-bound on the posterior [87, 89], and oftentimes need to assume Gaussian approximate posteriors through the reconstruction error.

The sequential neural likelihood (SNL, [122]) method and the method of emulated likelihoods [100] propose to learn the likelihood instead of the posterior. In this way, the trained neural ap-

proximators can be integrated into standard Bayesian pipelines and can be re-used across changes in the prior model $p(\boldsymbol{\theta})$, as long as the likelihood remains invariant. Sequential neural ratio (SNR) methods [40, 68], on the other hand, propose to train a classifier to approximate density ratios. These density ratios can then be used to sample from the posterior via standard MCMC methods. Of course, the computational time of SNL and SNR methods will still be dominated by the non-amortized components (e.g., MCMC or alternative sampling schemes) when faced with more than a few data sets. The same will be true for *inference compilation* approaches, which calibrate specialized neural networks through simulations in order to improve proposal distributions within (non-amortized) sequential Monte Carlo [117].

An interesting approach for amortized inference which does not rely on neural networks is the pre-paid estimation method without likelihoods [108]. This method memorizes a large database of pre-computed summary statistics for fast nearest-neighbor inference, aided by advanced interpolation methods. Even though the pre-paid method is very powerful and applicable to all kinds of models, it still crucially depends on the ability to (heuristically) select good summary statistics. Thus, in a future work, it seems worthwhile to explore the possibility of combining the pre-paid method with a neural network capable of learning maximally informative summary vectors (as proposed in our frameworks).

Ideas for direct posterior estimation via NDE are closely related to the concept of *optimal transport maps* and its application in Bayesian inference [11, 23, 36, 125]. A transport map defines a transformation between (probability) measures which can be constructed in a way to *warp* a simple probability distribution into a more complex one. In the context of Bayesian inference, transport maps have been applied to accelerate MCMC sampling [125], to perform sequential inference [36], and to solve inference problems via direct optimization [11]. In fact, our BayesFlow framework can be loosely viewed as a parameterization of invertible transport maps via invertible neural networks. An important distinction from this line of research is that NDE methods do not require an explicit likelihood function for approximating the target posteriors and are capable of amortized inference.

Curiously, throughout the advancement of amortized NDE methods, the task of simulation-based Bayesian model comparison appears to have taken a backseat. With certain caveats, neural density estimators can be adapted for posterior/prior predictive Bayesian model comparison by post-processing the samples from an approximate posterior/likelihood over each model’s parameters. However, such an approach will involve training a separate neural estimator for each model in the candidate set \mathcal{M} and has not yet been systematically investigated. In addition, most NDE methods also rely on fixed summary statistics [121, 122]) as inputs to the networks and few applications using raw data directly exist [60].

Alongside advancements in simulation-based inference, there has been an upsurge in the development of methods for uncertainty quantification in deep learning applications [74]. For instance, much work has been done on the efficient estimation of Bayesian neural networks [69, 96, 99] since the pioneering work of [102]. Parallel to the establishment of novel variational methods [85, 86], the drive for representing uncertainty has paved the way towards more interpretable and trustworthy neural network applications. Moreover, the need for distinguishing between different sources of uncertainty and the overconfidence of deep neural networks in both classification and regression tasks has been demonstrated quite effectively in recent works [82, 145].

Our frameworks are inspired by recent methods for deep probabilistic modeling [38, 88] and uncertainty representation in classification tasks [145, 156]. However, our goal is to efficiently approximate Bayesian posteriors and posterior odds between competing mechanistic models using *non-Bayesian neural networks*, not to estimate neural network parameters (e.g., weights) via Bayesian methods. Indeed, the incorporation of Bayesian neural networks into our frameworks appears to be an interesting avenue for future research. Moreover, our frameworks combine some of the latest ideas from simulation-based inference and uncertainty quantification for training efficient and uncertainty-aware estimators capable of amortized Bayesian parameter estimation and model comparison.

Accordingly, we propose to solve each task *globally*, that is, over the entire range of plausible parameters, data sets, and models. For parameter estimation, we will employ invertible neural networks (INN, [3, 4, 37, 38]). Previously, INNs have been successfully employed to tackle inverse problems in astrophysics and medicine [3]. We will adapt these flow-based INN architectures to suit the task of Bayesian parameter estimation in the context of various intractable model and data types. As for model comparison, we will employ *evidential neural networks*, which have previously been convincingly employed for uncertainty-aware classification [145].

Finally, by introducing our frameworks, we will further discuss many open questions, such as end-to-end estimation of various model classes (e.g., memoryless models, stateful models, joint models) via algorithmic alignment, model validation, Bayes factor approximation, and ideas for meta-amortized inference.

5 AMORTIZED PARAMETER ESTIMATION WITH BAYESFLOW

Let reality be reality. Let things flow naturally forward in whatever way they like.

— Laozi

Estimating the parameters of cognitive models is a crucial task in behavioral and cognitive modeling. However, the task can prove notably difficult or even impossible when a model can faithfully simulate behavior but the probabilistic form (i.e., the likelihood) of its outputs cannot be described analytically. In this chapter, we introduce our framework for amortized Bayesian parameter estimation which we coined *BayesFlow*. It comprises a new Bayesian solution to the simulation-only setting in terms of *invertible neural networks* and the theory of *normalizing flows*. The main idea behind BayesFlow is to split Bayesian analysis into a potentially expensive upfront training phase, followed by a much cheaper inference phase. The goal of the upfront training phase is to train a neural sampler that yields well-calibrated posteriors for *any* observed data set from the generative scope of a model. Subsequently, applying the neural sampler to specific observations during inference is very fast and can easily be performed in parallel, so that the training effort amortizes over repeated evaluations. The following chapter describes the mathematical details behind BayesFlow and discusses its strengths, limitations and future enhancements.

5.1 DESIDERATA

We have seen in the previous chapter that simulation-based methods need to optimize multiple, often conflicting, requirements concerning their performance. We therefore commence this chapter by stating the concrete desiderata for the utility of our framework:

1. Fully Bayesian estimation without framework-imposed restrictions on the type of priors, simulators, and posteriors amenable for inverse inference
2. Automatic extraction of maximally informative data representations instead of reliance on potentially suboptimal hand-crafted summary statistics
3. Scalability to high-dimensional problems (regarding both the data space \mathcal{X} and the parameter space Θ) via algorithmic alignment and considerations on probabilistic symmetry
4. Full amortization over multiple empirical data sets and data set sizes
5. Parallel computing and GPU acceleration applicable to forward inference (simulations), training, and inverse inference

6. Low memory demands, both during training (i.e., through online learning and on-the-fly simulations) and after training (i.e., no need to store large grids, reference tables or parameter databases)
7. A theoretical guarantee for convergence of the approximate posterior to the true posterior under certain optimal conditions

Throughout this chapter, we will gradually introduce our BayesFlow framework and discuss how it addresses each desideratum. Along the way, we will point out unexplored conceptual or empirical territories and lay out ideas for potential future improvements and applications.

5.2 BACKGROUND

5.2.1 DEEP GENERATIVE MODELING

As previously mentioned, BayesFlow draws on major advances in modern deep generative modeling, also referred to as deep probabilistic modeling. The core idea behind deep generative modeling is to represent a complicated target distribution as a non-linear transformation of some simpler latent distribution (e.g., Gaussian or uniform), a so called pushforward. Density estimation of the target distribution, a very complex task, is thus reduced to learning a non-linear transformation, a task that is ideally suited for gradient-based neural network training via standard backpropagation. Typically, deep generative models approximate the target distribution by sampling from the simpler latent distribution and applying the (inverse) transformation learned during gradient-based optimization. Consequently, we will train our neural networks to sample from intractable posterior distribution over the parameters of complex (behavioral) simulators.

Deep generative methods have demonstrated tremendous successes in applications dealing with very high-dimensional data, such as images, texts, or videos [88, 170, 175]. To draw an equivalent between these applications and simulation-based inference, consider a prototypical generative task in computer vision. In this context, the target distribution runs over the pixels of an image, and estimating a generative model of this distribution poses a major challenge. Conditional image generation, an even more challenging task, involves modeling the distribution over pixels contingent on additional information, such as image categories and descriptions (captions). Notably, the same ideas can be seamlessly transferred to model-based Bayesian inference, where the associated challenge is estimating the distribution over model parameters contingent on some observed data. Moreover, recent work demonstrating excellent generative performance with high-resolution images [88] suggests that deep generative models might be excellent candidates for overcoming the curse of dimensionality from which standard simulation-based methods notoriously suffer [27].

5.2.2 FORWARD INFERENCE

BayesFlow depends on the ability to (efficiently) perform forward inference, that is, simulate artificial data sets given possible parameter configurations. Moreover, in order to use the simulator as a calibrator for the neural networks involved in BayesFlow, we need to simulate multiple data sets from the Bayesian joint distribution $p(\boldsymbol{\theta}, \mathbf{x}_{1:N}) = p(\boldsymbol{\theta}) p(\mathbf{x}_{1:N} | \boldsymbol{\theta})$, which, for memoryless

models, decomposes further into $p(\boldsymbol{\theta}) \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\theta})$. As already mentioned, the decomposition of the joint distribution is essentially a generative receipt, which can be carried out with or without an explicit likelihood (due to Equation 2.4).

Corresponding to the previous description, Algorithm 1 describes the steps for generating a batch of synthetic data sets using simulations from a memoryless model and a randomly drawn data size N for each batch. The data-set size will typically be drawn from a discrete uniform distribution, $p(N) = \mathcal{U}(N_{\min}, N_{\max})$, but fixed sizes or different distributions are possible (and reasonable for certain applications). Note also, that the algorithm is trivially extendable to stateful models by including memory variables or an explicit dependence of the simulator on previous outputs.

Algorithm 1 Monte Carlo generation of B synthetic data sets

Require: $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ - stochastic model simulator, $p(\boldsymbol{\theta})$ - prior over model parameters, $p(\boldsymbol{\xi})$ - noise distribution, $p(N)$ - distribution over data set sizes, B - number of data sets to generate.

- 1: Draw data set size: $N \sim p(N)$.
 - 2: **for** $b = 1, \dots, B$ **do**
 - 3: Draw model parameters from prior: $\boldsymbol{\theta}^{(b)} \sim p(\boldsymbol{\theta})$.
 - 4: **for** $n = 1, \dots, N$ **do**
 - 5: Sample noise instance: $\boldsymbol{\xi}_n \sim p(\boldsymbol{\xi})$.
 - 6: Run simulator to obtain n -th synthetic observation: $\mathbf{x}_n = g(\boldsymbol{\theta}^{(b)}, \boldsymbol{\xi}_n)$.
 - 7: **end for**
 - 8: Store pair $(\boldsymbol{\theta}^{(b)}, \mathbf{x}_{1:N}^{(b)})$.
 - 9: **end for**
 - 10: **Return** mini-batch $\mathcal{D}_N^{(B)} := \{\boldsymbol{\theta}^{(b)}, \mathbf{x}_{1:N}^{(b)}\}_{b=1}^B$.
-

Importantly, the efficiency of Algorithm 1 depends highly on its actual implementation. The naive complexity of the data-generation algorithm is at least $\mathcal{O}(N * B * G)$, where G denotes the cost of executing the simulator once to obtain a single synthetic observation \mathbf{x}_n . Thus, the algorithm can benefit from three levels of parallelism: i) over the number of data sets (B); ii) over the number of observations in each data set (N); iii) and over the computational steps of the simulator itself (G). In the ideal case where all levels can be executed in parallel, the computational complexity reduces to $\mathcal{O}(1)$. For some applications, even parallelizing the most costly level can bring about a significant speedup in practice. Notably, the parallelization of data generation with stateful models is generally not immediately obvious due to the sequential dependence of each output on previous outputs (the loop over N), but multiple data sets from a stateful model can still be readily generated in parallel (the loop over B).

5.2.3 NORMALIZING FLOWS

In order to perform neural density estimation, we will implement a *normalizing flow* via an invertible neural network (INN, [37, 38]). A normalizing flow represents a transformation of a simple probability density (e.g., Gaussian) into a more complex (unknown) density by a sequence of invertible and differentiable mappings [38]. In contrast to variational methods [12], flow-based methods can perform asymptotically exact inference by using lossless compression. Additionally,

they scale favourably from simple low-dimensional problems to (potentially intractable) high-dimensional distributions with complex dependencies [4, 88].

To set the stage, let $\mathbf{z} \in \mathbb{R}^D$ be a random variable with a known (simple) probability density and $\boldsymbol{\theta} \in \mathbb{R}^D$ a random variable with an unknown (complicated) probability density. Let $f : \mathbb{R}^D \rightarrow \mathbb{R}^D$ be an invertible, differentiable function such that $\mathbf{z} = f(\boldsymbol{\theta})$ and $\boldsymbol{\theta} = f^{-1}(\mathbf{z})$. By using the change of variable rule of probability theory, the density of the variable $\boldsymbol{\theta}$ can be computed as:

$$p(\boldsymbol{\theta}) = p(\mathbf{z} = f(\boldsymbol{\theta})) \left| \det \left(\frac{\partial f(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) \right| \quad (5.1)$$

In our framework, we use a unit Gaussian as a base distribution, $p(\mathbf{z}) = \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I})$, and the pushforward density will be the posterior $p(\boldsymbol{\theta} | \mathbf{x})$ over model parameters $\boldsymbol{\theta}$ given a single (for now) observation \mathbf{x} . Thus, our aim is to learn an approximate posterior q which matches the pushforward posterior. Accordingly, we reparameterize the approximate posterior in terms of a conditional invertible neural network (cINN) estimator f_ϕ which implements a normalizing flow between $\boldsymbol{\theta}$ and \mathbf{z} given observation \mathbf{x} :

$$q_\phi(\boldsymbol{\theta} | \mathbf{x}) = p(\mathbf{z} = f_\phi(\boldsymbol{\theta}; \mathbf{x})) \left| \det \left(\frac{\partial f_\phi(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} \right) \right| \quad (5.2)$$

Accordingly, sampling from the approximate posterior involves sampling from the base density and transforming the sample via the inverse operation of the cINN into the pushforward posterior:

$$\boldsymbol{\theta} \sim q_\phi(\boldsymbol{\theta} | \mathbf{x}) \iff \boldsymbol{\theta} = f_\phi^{-1}(\mathbf{z}; \mathbf{x}) \text{ with } \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \quad (5.3)$$

Thus, our solution to the task of simulation-based parameter estimation is to train a cINN which approximates the true posterior for all possible observations \mathbf{x} arising from a given model \mathcal{M}_j as accurately as possible. There are many ways to implement a cINN in practice, and we will illustrate our preferred architecture below based on coupling flows. Note, however, that the field of deep generative modeling is rapidly expanding, with novel architectures emerging almost on a daily basis. Thus, it is likely that the concrete cINN architecture proposed in this work might be supplanted by a better candidate in the not-so-distant future. In any case, the problem of inverse inference in science is here to stay, so the BayesFlow framework could easily be adapted to incorporate a different neural sampler f_ϕ .

5.2.4 COUPLING FLOWS

Coupling flows are one of the most widely used invertible architectures [92] because they are i) conceptually simple; ii) easily invertible; and iii) able to represent highly expressive transformation with tractable Jacobian determinants [88]. A coupling flow $C_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D$ with trainable parameters ϕ between $\boldsymbol{\theta} \in \mathbb{R}^D$ and $\mathbf{z} \in \mathbb{R}^D$ can be realized as follows. Consider a disjoint partition of the input parameters $\boldsymbol{\theta} \in \mathbb{R}^D$ into two subspaces: $(\boldsymbol{\theta}^A, \boldsymbol{\theta}^B) \in \mathbb{R}^d \times \mathbb{R}^{D-d}$. Input

partitioning is required by most coupling flows to ensure invertibility [92]. The invertible forward transformation of a coupling flow can be defined as:

$$\mathbf{z}^A = C_1(\boldsymbol{\theta}^A; \Omega_1(\boldsymbol{\theta}^B)) \quad (5.4)$$

$$\mathbf{z}^B = C_2(\boldsymbol{\theta}^B; \Omega_2(\mathbf{z}^A)) \quad (5.5)$$

$$\mathbf{z} = (\mathbf{z}^A, \mathbf{z}^B) \quad (5.6)$$

where C_1 and C_2 are invertible functions and Ω_1 and Ω_2 are called *conditioners*. The conditioners can be realized via arbitrarily complex functions (e.g., deep neural networks) which themselves need not be invertible, as long as C_1 and C_2 are (easily) invertible. Correspondingly, the inverse transformation C_ϕ^{-1} of a coupling flow is defined as:

$$\boldsymbol{\theta}^B = C_2^{-1}(\mathbf{z}^B; \Omega_2(\mathbf{z}^A)) \quad (5.7)$$

$$\boldsymbol{\theta}^A = C_1^{-1}(\mathbf{z}^A; \Omega_1(\boldsymbol{\theta}^B)) \quad (5.8)$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}^A, \boldsymbol{\theta}^B) \quad (5.9)$$

Note, that there are many viable ways to parameterize a coupling flow [92]. Our BayesFlow method uses a composition of conditional affine coupling layers (cACLs). A single cACL performs the following bijective mapping on its split input

$$\mathbf{z}^A = \boldsymbol{\theta}^A \odot \exp(S_1(\boldsymbol{\theta}^B; \mathbf{x})) + T_1(\boldsymbol{\theta}^B; \mathbf{x}) \quad (5.10)$$

$$\mathbf{z}^B = \boldsymbol{\theta}^B \odot \exp(S_2(\mathbf{z}^A; \mathbf{x})) + T_2(\mathbf{z}^A; \mathbf{x}) \quad (5.11)$$

$$\mathbf{z} = (\mathbf{z}^A, \mathbf{z}^B) \quad (5.12)$$

where \odot denotes element-wise multiplication and the functions S_1, S_2, T_1, T_2 are implemented as fully connected (FC) neural networks with \mathbf{x} passed through an additional input head. By construction, this bijection works independently of the form of the functions s and t , which themselves are never inverted. The inverse transformation of the cACL is thus given by:

$$\boldsymbol{\theta}^B = (\mathbf{z}^B - T_2(\mathbf{z}^A; \mathbf{x})) \odot \exp(-S_2(\mathbf{z}^A; \mathbf{x})) \quad (5.13)$$

$$\boldsymbol{\theta}^A = (\mathbf{z}^A - T_1(\boldsymbol{\theta}^B; \mathbf{x})) \odot \exp(-S_1(\boldsymbol{\theta}^B; \mathbf{x})) \quad (5.14)$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}^A, \boldsymbol{\theta}^B) \quad (5.15)$$

The Jacobian of the forward coupling transformation is a product of two triangular matrices

$$\frac{\partial C_\phi(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \text{diag}(\exp(S_1(\boldsymbol{\theta}^B; \mathbf{x}))) & \text{finite} \\ 0 & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mathbf{1} & 0 \\ \text{finite} & \text{diag}(\exp(S_2(\mathbf{z}^A; \mathbf{x}))) \end{bmatrix}, \quad (5.16)$$

so its determinant is easy to compute and given by

$$\det \mathbf{J}_{C_\phi} = \exp\left(\sum_{i=1}^d S_1(\boldsymbol{\theta}^B; \mathbf{x})_i + \sum_{i=1}^{D-d} S_2(\mathbf{z}^A; \mathbf{x})_i\right), \quad (5.17)$$

where we have abbreviated the Jacobian of the cACL as \mathbf{J}_{C_ϕ} . In section 5.3.1, we will show how to compose multiple cACLs into an invertible network and discuss some additional features for improving the basic design introduced here. Before we introduce our complete Bayesian framework for amortized inference, we briefly peruse the concepts of distribution matching and amortized inference.

5.2.5 DISTRIBUTION MATCHING AND AMORTIZATION

The term amortized inference refers to an approach which reduces the cost of inference by casting some or all inferential phases as an optimization task. In particular, for a given simulator, one can approximate an unknown ground-truth posterior distribution $p(\boldsymbol{\theta} | \mathbf{x})$ via a parameterized distribution $q_\phi(\boldsymbol{\theta} | \mathbf{x})$ by minimizing some f -divergence between the two distributions:

$$\phi^* = \arg \min_{\phi} D_f(p(\boldsymbol{\theta} | \mathbf{x}) || q_\phi(\boldsymbol{\theta} | \mathbf{x})) \quad (5.18)$$

$$= \arg \min_{\phi} \int_{\Theta} q_\phi(\boldsymbol{\theta} | \mathbf{x}) f\left(\frac{p(\boldsymbol{\theta} | \mathbf{x})}{q_\phi(\boldsymbol{\theta} | \mathbf{x})}\right) d\boldsymbol{\theta} \quad (5.19)$$

where f is a convex function. Usually, the Kullback-Leibler (\mathbb{KL}) divergence is chosen, so the objective becomes:

$$\phi^* = \arg \min_{\phi} \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{x}) \log \frac{p(\boldsymbol{\theta} | \mathbf{x})}{q_\phi(\boldsymbol{\theta} | \mathbf{x})} d\boldsymbol{\theta} \quad (5.20)$$

$$= \arg \min_{\phi} \mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})}[-\log q_\phi(\boldsymbol{\theta} | \mathbf{x})] - \mathbb{H}[p(\boldsymbol{\theta} | \mathbf{x})] \quad (5.21)$$

$$= \arg \min_{\phi} \mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})}[-\log q_\phi(\boldsymbol{\theta} | \mathbf{x})] \quad (5.22)$$

where $\mathbb{H}[p(\boldsymbol{\theta} | \mathbf{x})]$ in Equation 5.21 is the Shannon entropy of the true posterior and can be dropped from the optimization objective since it does not depend on ϕ . We will now differentiate between three types of amortized inference, which all leverage the fact that we can generate synthetic datasets via a scientific simulator. We believe that such an explicit distinction is important, given the current abundance of neural network methods for Bayesian inference.

Case-wise amortization In case-wise amortized inference, we perform an optimization loop for each individual observation \mathbf{x} (e.g., as in sequential neural posterior estimation, [40, 60]). Thus, in expectation over all possible observations, the criterion can be expressed as:

$$\mathbb{E}_{p^*(\mathbf{x})} \left[\min_{\phi_{\mathbf{x}}} \mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})}[-\log q_\phi(\boldsymbol{\theta} | \mathbf{x})] \right] \quad (5.23)$$

and inference is only amortized in the context of individual observations and models. This implies that different neural network parameters $\phi_{\mathbf{x}}^*$ are found for each approximate posterior defined by a particular observation \mathbf{x} and a particular model \mathcal{M}_j . Crucially, the fact that the minimization objective is inside the outer expectation can render inference infeasible when multiple observations and models are available.

Model-wise amortization In the model-wise amortized scenario, optimization is performed globally for the entire range of plausible observations, which involves pulling the minimum operator (\min) out of the outer expectation:

$$\min_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [-\log q_{\phi}(\boldsymbol{\theta} | \mathbf{x})]] \quad (5.24)$$

In this case, the training effort amortizes over the entire data range for a given model. This is the main approach taken in the *BayesFlow* method. The training (optimization) phase in the model-wise amortized setting is considerably longer than the training phase in the case-wise amortized setting. However, once optimization has converged to an approximator of ϕ^* , the resulting neural estimator f_{ϕ^*} can be reused for arbitrarily many observations assumed to arise from a given model \mathcal{M}_j . In some cases, the break-even in terms of efficiency between case-wise and model-wise amortized inference occurs even after a few observations, without noticeable accuracy degradation [133]. However, when multiple candidate models should be estimated and compared, the training effort can become prohibitively large, since a separate set of neural network parameters needs to be learned for each model.

Meta-amortization In the meta-amortized setting, optimization is performed over all possible models simultaneously, introducing one more expectation into the objective:

$$\min_{\phi_{\mathcal{M}}} \mathbb{E}_{p(\mathcal{M})} [\mathbb{E}_{p^*(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x}, \mathcal{M})} [-\log q_{\phi}(\boldsymbol{\theta} | \mathbf{x}, \mathcal{M})]]] \quad (5.25)$$

In this way, an estimator with parameters $\phi_{\mathcal{M}}^*$ (the subscript denoting amortization over the entire model set \mathcal{M}) can be reused for inference on multiple observations with an arbitrary number of competing models from a particular research domain or model class. Importantly, such setting can only be useful if the latent parameter spaces are allowed to vary across the models, as competing models in various domains can have widely different parameterizations. Note also, that both the model-wise and meta-amortized setting can employ amortization over different dataset sizes N (to be discussed shortly) as long as the dimensionality of the data summary statistic stays the same, a property which can be of great utility in practice.

Nevertheless, each further amortization step might introduce an *amortization gap*. The issue has been discussed in the context of variational inference [25] and refers to a potential drop in performance as a consequence of optimizing neural network parameters in expectation as opposed to optimizing for each individual observation. However, a potential amortization gap has not been investigated outside the context of variational inference and warrants an empirical assessment in a meta-amortized context.

In a later chapter, we will show how to make the meta-amortized objective tractable by using ideas from the literature on normalizing flows and multi-task learning. For now, however, we focus on model-wise amortization with BayesFlow.

5.3 BAYESFLOW: BUILDING AMORTIZED NEURAL SAMPLERS

At a high level, *BayesFlow* [133] incorporates a summary network h and an inference network f to jointly invert a generative Bayesian model. The *summary network* $h(\mathbf{x}_{1:N})$ reduces data sets

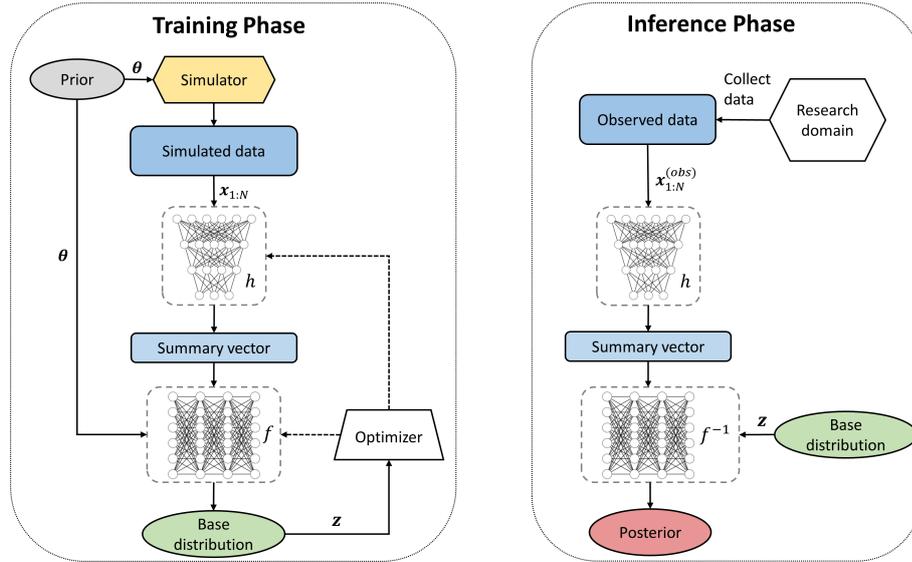


Figure 5.1: Both phases of the BayesFlow framework. Left panel: During the training phase of our BayesFlow method, a summary (h) and an inference network (f) are trained jointly with random draws from the prior and synthetic data from the simulator; Right panel: During the inference phase, BayesFlow works entirely in a feed-forward manner, that is, no training or optimization happens in this phase. The upfront training effort amortizes over arbitrary many observations and data sets from a research domain working on the same model family.

of arbitrary size to fixed-size vector representations. The *inference network* samples from an approximate posterior q via a conditional invertible neural network (cINN) f which implements a normalizing flow between θ and a normally distributed \mathbf{z} given the outputs of the summary network:

$$q(\theta | \mathbf{x}_{1:N}) = p(\mathbf{z} = f(\theta; h(\mathbf{x}_{1:N}))) \left| \det \left(\frac{\partial f(\theta; h(\mathbf{x}_{1:N}))}{\partial \theta} \right) \right| \quad (5.26)$$

where the dependence on all neural network parameters is implicit and has been omitted for clarity. The introduction of a summary network whose structure is aligned to the structure of the simulator (i.e., stateless vs. stateful) frees our framework from a restriction to a particular model class or data type. Moreover, the summary network itself does not have to be invertible, since its output is concatenated with θ and fed to each coupling layer, but not directly mapped to \mathbf{z} . Figure 5.1 illustrates the different components and phases of our BayesFlow framework.

5.3.1 COMPOSING INVERTIBLE NETWORKS

In this section, we describe how to stack multiple coupling layers to obtain a deep invertible network. For now, consider the case when raw simulated data $\mathbf{x}_{1:N}$ of size $N = 1$ is entered directly into the invertible network without using a summary network. In order to ensure that our architecture is expressive enough to encode complex posterior distributions, we chain multiple ACLs,

so that the output of each ACL becomes the input to the next one. In this way, the whole network remains invertible from the first input to the last output and can be viewed as a single bijective function. Such chaining of operations is possible, since a composition of invertible functions is itself invertible and its Jacobian determinant is the product of the Jacobian determinants of the individual coupling blocks. Therefore, we refer to a cINN as a composition of K conditional ACLs.

$$\mathbf{z} = f_{\phi}(\boldsymbol{\theta}; \mathbf{x}) \equiv C_{\phi_K} \circ C_{\phi_{K-1}} \circ \dots \circ C_{\phi_1}(\boldsymbol{\theta}; \mathbf{x}) \quad (5.27)$$

with trainable parameters $\phi = (\phi_1, \dots, \phi_K)$ and inverse:

$$\boldsymbol{\theta} = f_{\phi}^{-1}(\mathbf{z}; \mathbf{x}) \equiv C_{\phi_1}^{-1} \circ C_{\phi_2}^{-1} \circ \dots \circ C_{\phi_K}^{-1}(\mathbf{z}; \mathbf{x}) \quad (5.28)$$

Note, that the observation \mathbf{x} (or a transformation thereof) is fed unchanged to each coupling layer C_{ϕ_k} . In our applications, the input to the first ACL is the parameter vector $\boldsymbol{\theta}$, and the output of the final ACL is a D -dimensional vector \mathbf{z} representing the non-linear transformation of the parameters into \mathbf{z} -space. Shortly, we will show how to ensure that \mathbf{z} follows a unit Gaussian distribution through optimization, that is, we will enforce $p(\mathbf{z}) = \mathcal{N}_D(\mathbf{z} | 0, \mathbb{I})$. We also use fixed permutation matrices before each ACL to ensure that each axis of the transformed parameter space \mathbf{z} encodes information from all components of $\boldsymbol{\theta}$, in order to capture posterior dependencies (e.g., posterior covariance). In addition, we apply soft clamping of the exponential outputs in each ACL for numerical stability.

Intuitively, our cINN realizes the following process: the forward pass maps data-generating parameters $\boldsymbol{\theta}$ to \mathbf{z} -space using conditional information from the observation \mathbf{x} , while the inverse pass maps data points from \mathbf{z} -space to the data-generating parameters of interest using the same conditional information.

5.3.2 SUMMARY NETWORKS

Since the number of observations might vary in practical scenarios (e.g., different number of trials or time points) or measurements might arrive in streams, we need to perform some form of dimensionality reduction on the data before feeding it to the cINN. As previously mentioned, we want to avoid information loss through restrictive hand-crafted summary statistics and, instead, learn the most informative finite summary vectors directly from data. Therefore, instead of feeding the raw simulated or observed data to the cINN, we pass the data through an auxiliary summary network to obtain a fixed-sized vector representation $\tilde{\mathbf{x}} = h_{\psi}(\mathbf{x}_{1:N})$.

As already alluded to in previous sections, the architecture of the summary network should match the probabilistic symmetry of the observed data (which, in turn, is dictated by the simulator), a property we refer to as *algorithmic alignment* [174]. In other words, different network architectures are needed for exchangeable (generated by memoryless models) and non-exchangeable (generated by stateful models) data. In the following, we illustrate three common scenarios in model-based inference.

Memoryless Models Memoryless models typically generate *i.i.d.* observations, which imply exchangeability and induce permutation invariant posteriors. In other words, changing (permuting) the order of individual elements should not change the associated likelihood or posterior (see

Section 3.4). Memoryless models abide in the cognitive sciences [39, 43, 138, 162], mainly due to their convenient simplicity, but also due to the computational limitations of existing methods for Bayesian estimation. Following [13], we encode probabilistic permutation invariance through functional permutation invariance realized by a deep invariant network. Such a network is capable of learning expressive permutation invariant functions through a combination of *equivariant* and *invariant* transformations.

First, we can obtain a permutation invariant function via an *invariant module* Σ_I which performs an equivariant non-linear transformation h_1 followed by a pooling operator (e.g., sum or max) and another non-linear transformation h_2 :

$$\tilde{\mathbf{x}} = \Sigma_I(\mathbf{x}_{1:N}) = h_1 \left(\sum_{n=1}^N h_2(\mathbf{x}_n) \right) \quad (5.29)$$

where h_1 and h_2 can be arbitrarily complex neural networks (cf. Figure 5.2, lower left panel). Second, in order to increase the capacity of the invariant transformation, we can stack together multiple *equivariant modules* Σ_E . Each equivariant module implements a learnable equivariant transformation by performing the following operations for each input element \mathbf{x}_n :

$$\mathbf{y}_n = \Sigma_E(\mathbf{x}_n, \tilde{\mathbf{x}}) = h_3(\mathbf{x}_n, \tilde{\mathbf{x}}) \quad \text{for } n = 1, \dots, N, \quad (5.30)$$

so that Σ_E is a combination of element-wise (equivariant) and invariant transforms (cf. Figure 5.1, lower middle panel). Again, the internal function h_3 can be parameterized via an arbitrary feed-forward neural network. Importantly, each equivariant module also contains a separate invariant model whose output is concatenated with each observation in order to increase the expressiveness of the learned transformation.

Finally, we can stack multiple equivariant modules followed by an invariant module, in order to obtain a deep invariant summary network $h_\psi : \mathcal{X}^N \rightarrow \mathbb{R}^S$:

$$\tilde{\mathbf{x}} = h_\psi(\mathbf{x}_{1:N}) = (\Sigma_I \circ \Sigma_E^{(K)} \circ \Sigma_E^{(K-1)} \circ \dots \circ \Sigma_E^{(1)})(\mathbf{x}_{1:N}), \quad (5.31)$$

where ψ denotes the vector of all learnable neural network parameters and S denotes the dimensionality of the output layer of the last invariant module Σ_I . The complete inference phase of BayesFlow using a deep invariant summary network is depicted in Figure 5.1.

Stateful Models Stateful models incorporate some form of memory and are thus capable of generating observations with complex dependencies. A prime example are dynamic models, which typically describe the evolution trajectory of a system or a process, such as an infectious disease, over time [81]. Observations generated from such models are usually the solution of a stochastic differential equation (SDE) and imply a more complex probabilistic symmetry than those generated from memoryless models.

In a recent application of the BayesFlow framework for estimating key epidemiological parameters [135], we have proposed a set-up specifically designed to tackle dynamic models with simulation-based inference. Our BayesFlow architecture comprises three sub-networks: (i) a convolutional *filtering* network performing noise reduction and feature extraction on the raw time-series data; (ii) a recurrent *summary* network reducing pre-processed time-series of *varying* length

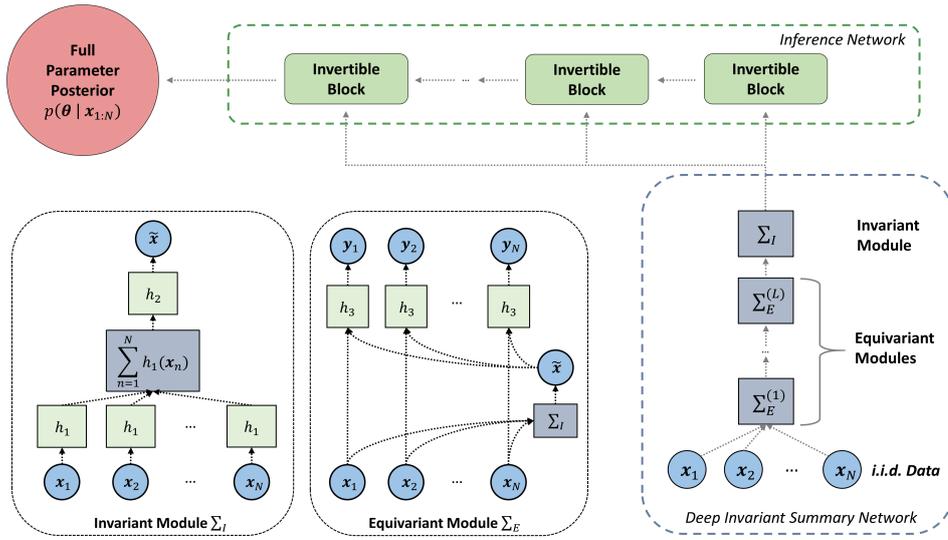


Figure 5.2: Inference with BayesFlow on *i.i.d.* data from a memoryless model using a deep permutation invariant summary network. The summary network is composed of a sequence of flexible equivariant neural modules followed by an invariant neural module. In this way, *i.i.d.* data (sets) of varying length are embedded into fixed-size vector representations which carry maximal information for posterior inference with a memoryless model.

to statistical summaries of *fixed* size; (iii) a cINN inference network performing Bayesian parameter inference given the learned summary vectors of the observations. Figure 5.2 depicts the architecture of this composite network.

The design of the convolutional network is inspired by that of the *Inception* neural architecture which has demonstrated tremendous success in a wide variety of computer vision tasks [152]. In particular, our network is implemented as a deep fully convolutional network which applies adjustable one-dimensional filters of different size at each level (cf. Figure 5.3, lower left panel). The intuition behind this design is that filters of different size might capture patterns at different temporal scales. For instance, if $t = 1, \dots, T$ is measured in days, a filter of size one will capture daily fluctuations whereas a filter of size seven will capture weekly dynamics. This, in turn, should ease the task of extracting informative temporal features for Bayesian updating.

The output of the convolutional network is a multivariate sequence containing the filtered time-series $\mathbf{x}_{1:T}$. In order to reduce the filtered sequence to a fixed-size vector, we pass it through a long-short term memory (LSTM) recurrent network [57]. Importantly, the LSTM network (see Figure 5.3, lower right panel) can deal with sequences of varying length, which enables online learning (i.e., Bayesian updating when new observations become available) and makes the same inference network applicable to settings with different degrees of data availability. Compared to a fixed pooling operation (e.g., mean or max), our many-to-one recurrent network performs a learnable pooling operation which respects the sequential probabilistic symmetry of the data. In this way, the composite summary network learns to filter and extract the most informative fea-

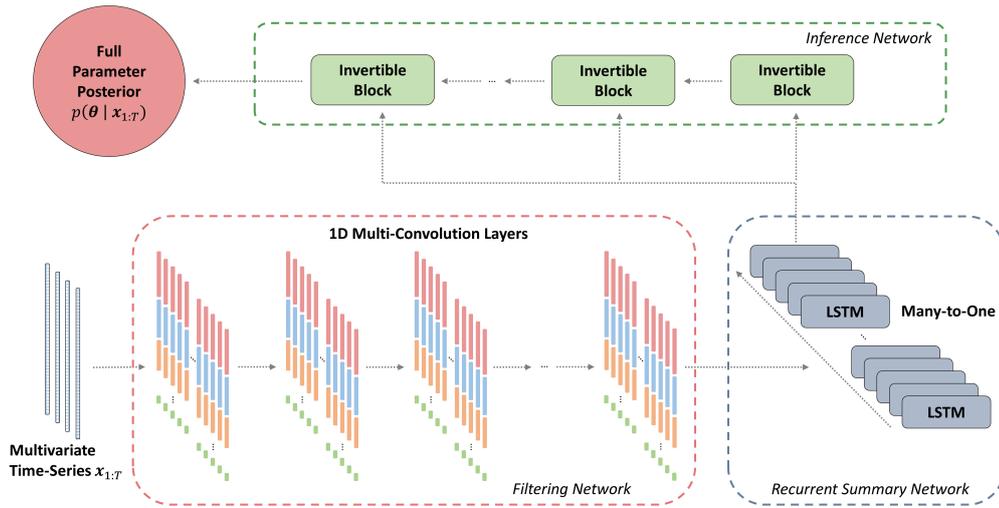


Figure 5.3: Inference on multivariate time-series data arising from a stateful model using a composite summary network architecture. The summary network is composed of an inception-like 1D fully convolutional network, followed by a many-to-one recurrent LSTM network. In this way, time-series of varying length are embedded into fixed-size vector representations which carry maximal information for posterior inference with a stateful model.

tures from the noisy observations in an end-to-end manner, such that no manual (and potentially suboptimal) selection of hand-crafted data features is required from the user at any point.

More formally, let us denote the functions represented by the filtering and summary networks as h_c and h_r . Then, the convolutional filtering network yields a filtered time-series $\tilde{x}_{1:T'} = h_c(x_{1:T})$ from observed data $x_{1:T}$, where the number of time steps T' may vary according to data availability. The recurrent summary network turns the outputs of the filtering network into fixed-size vectors $\tilde{x} = h_r(\tilde{x}_{1:T'})$. The cINN thus generates samples $\theta \sim q_\phi(\theta | \tilde{x})$ from the parameter posterior by computing $\theta = f_\phi^{-1}(z, \tilde{x})$ with normally distributed random vectors $z \sim \mathcal{N}_D(0, \mathbb{I})$. The complete inference phase of BayesFlow using a deep sequence network is depicted in Figure 5.3.

Joint Models Joint models are an attempt to account for different processes (e.g., neural and cognitive) within a single composite model [30, 118, 159]. Thus, joint models integrate different sources and types of data and require more complex summary architectures. A hypothetical scenario with three data sources (e.g., behavioral data, neural data, and eye-tracking data) is depicted in Figure 5.4. In this case, a separate processor network, each aligned to the particular data type, reduces a separate set of observations from a given source. The outputs of the individual processor networks are then concatenated and fed through an *integrator network*, which combines the information from all processor networks into a single vector representation. In this way, the main cINN architecture can remain the same as in the previous examples and utilize the modularity of neural network. Since the application of integrative joint models is still in its infancy, such composite BayesFlow architectures are yet to prove their usefulness.

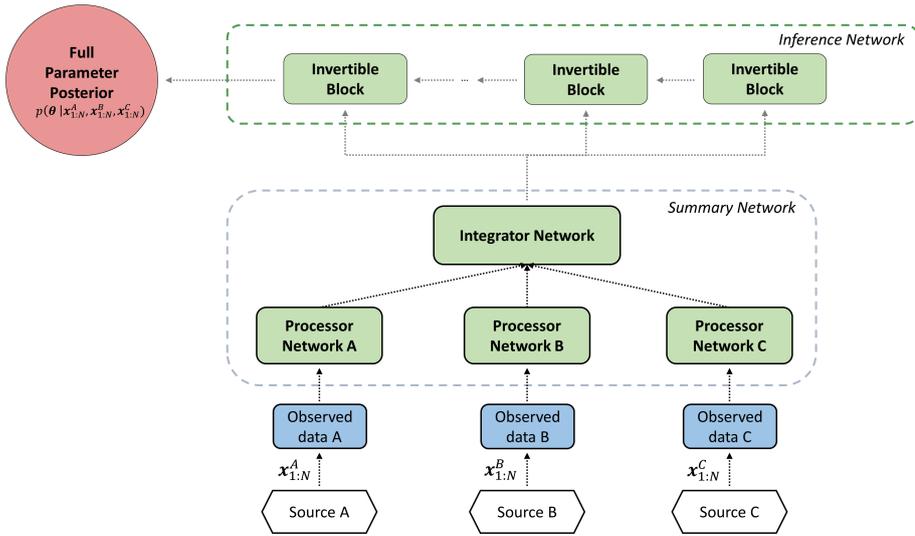


Figure 5.4: Inference on different data sources from a joint model using an ensemble of summary networks and an integrator network. The different data sets, potentially of varying size and structure, are processed by separate algorithmically aligned summary networks. The outputs of all summary networks are then combined into a fixed-size vector representations by the integrator network, which informs the inference network about the joint posterior over all model parameters.

Regardless of the summary network’s concrete design, its parameters ψ are optimized jointly with those of the cINN via backpropagation. Thus, the training phase remains completely end-to-end, and BayesFlow learns to generalize to data sets of different sizes by suitably varying N during training (see Algorithm 1).

5.3.3 OPTIMIZATION OBJECTIVE

For *any* given (simulated or observed) dataset $\mathbf{x}_{1:N}$, our framework needs to ensure that the inverse transformation of the trained cINN, $\theta = f_{\phi}^{-1}(z; h_{\psi}(\mathbf{x}_{1:N}))$ with $z \sim \mathcal{N}_D(z | 0, \mathbb{I})$, yields samples from the true posterior $p(\theta | \mathbf{x}_{1:N})$. To achieve this, we resort to the concept of distribution matching introduced earlier and minimize the expected KL divergence between the true and the approximate posterior for all possible observations within the generative scope of a

Monte Carlo simulator. For clarity, we first derive our optimization objective with $N = 1$ and no summary network,

$$\phi^* = \arg \min_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}) || q_{\phi}(\boldsymbol{\theta} | \mathbf{x}))] \quad (5.32)$$

$$= \arg \min_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [\log p(\boldsymbol{\theta} | \mathbf{x}) - \log q_{\phi}(\boldsymbol{\theta} | \mathbf{x})]] \quad (5.33)$$

$$= \arg \min_{\phi} \mathbb{E}_{p^*(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [-\log q_{\phi}(\boldsymbol{\theta} | \mathbf{x})]] \quad (5.34)$$

$$= \arg \min_{\phi} - \int_{\mathcal{X}} p^*(\mathbf{x}) \int_{\Theta} p(\boldsymbol{\theta} | \mathbf{x}) \log q_{\phi}(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} d\mathbf{x} \quad (5.35)$$

which corresponds to model-wise amortization, as defined earlier. To render optimization of this criterion tractable, we first apply the change of variable rule to $q_{\phi}(\boldsymbol{\theta} | \mathbf{x})$ as given in Equation 5.2 to obtain:

$$\phi^* = \arg \min_{\phi} - \int_{\mathcal{X}} \int_{\Theta} p^*(\mathbf{x}) p(\boldsymbol{\theta} | \mathbf{x}) (\log p(f_{\phi}(\boldsymbol{\theta}; \mathbf{x})) + \log |\det \mathbf{J}_{f_{\phi}}|) d\boldsymbol{\theta} d\mathbf{x} \quad (5.36)$$

where we have abbreviated $\partial f_{\phi}(\boldsymbol{\theta}; \mathbf{x}) / \partial \boldsymbol{\theta}$ (the Jacobian of the entire cINN f_{ϕ} evaluated at $\boldsymbol{\theta}$ and \mathbf{x}) as $\mathbf{J}_{f_{\phi}}$ and moved $p^*(\mathbf{x})$ inside the inner expectation, as it does not depend on $\boldsymbol{\theta}$. Since Equation 5.36 defines an expectation over the true and unknown data-generating distribution, we replace it with the Bayesian joint model $p(\boldsymbol{\theta}, \mathbf{x})$ from which we can obtain Monte Carlo samples (e.g., by using Algorithm 1). Accordingly, for a batch of B parameters and corresponding synthetic data sets $\mathcal{D}^{(B)} = \{(\boldsymbol{\theta}^{(b)}, \mathbf{x}^{(b)})\}_{b=1}^B$, we can define the following loss function

$$\mathcal{L}(\phi) = \frac{1}{B} \sum_{b=1}^B \left(-\log p(f_{\phi}(\boldsymbol{\theta}^{(b)}; \mathbf{x}^{(b)})) - \log |\det \mathbf{J}_{f_{\phi}}^{(b)}| \right) \quad (5.37)$$

$$= \frac{1}{B} \sum_{b=1}^B \left(\frac{\|f_{\phi}(\boldsymbol{\theta}^{(b)}; \mathbf{x}^{(b)})\|_2^2}{2} - \sum_{k=1}^K \log |\det \mathbf{J}_{C_{\phi_k}}^{(b)}| \right), \quad (5.38)$$

which we minimize using standard backpropagation to arrive at an unbiased estimate $\hat{\phi}$ of ϕ^* . The first term follows from Equation 5.36 due to the fact that we have prescribed a unit Gaussian distribution to \mathbf{z} . It represents the negative log of $\mathcal{N}_D(\mathbf{z} | 0, \mathbb{I}) \propto \exp(-\frac{1}{2}\|\mathbf{z}\|_2^2)$. The second term follows from Equation 5.27 and controls the rate of volume change induced by the non-linear transformation from $\boldsymbol{\theta}$ to \mathbf{z} learned by f_{ϕ} . Thus, minimizing Equation 5.38 ensures that \mathbf{z} follows the prescribed unit Gaussian and that f_{ϕ}^* is a model-wise amortized neural sampler which yields independent samples from the true posterior under perfect convergence [133].

When the number of observations varies during inference, we need to vary it during training as well, in order to achieve amortization over data sets $\mathbf{x}_{1:N}$ with different sizes (if required by the application). Thus, we introduce a suitable summary network h_{ψ} which renders the cINN

independent of N and learns to extract maximally informative statistics from the (raw) simulated data in an end-to-end manner. Our modified criterion then becomes:

$$\phi^*, \psi^* = \arg \min_{\phi, \psi} \mathbb{E}_{p(N, \theta, \mathbf{x})} [-\log q_{\phi}(\theta | h_{\psi}(\mathbf{x}_{1:N}))] \quad (5.39)$$

Accordingly, our modified loss function for a batch $\mathcal{D}_N^{(B)} = \{(\theta^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ simulated from the Bayesian model $p(N, \theta, \mathbf{x})$ becomes:

$$\mathcal{L}(\phi, \psi) = \frac{1}{B} \sum_{b=1}^B \left(\frac{\|f_{\phi}(\theta^{(b)}; h_{\psi}(\mathbf{x}_{1:N}^{(b)}))\|_2^2}{2} - \sum_{k=1}^K \log |\det \mathbf{J}_{C_{\phi_k}^{(b)}}| \right), \quad (5.40)$$

which corresponds to a trivial change that simply sets the conditioning vector of the cINN to the output of the summary network. Again, we can use backpropagation with any gradient-based optimization method to obtain unbiased estimates $\hat{\phi}, \hat{\psi}$ of the optimal neural network parameters ϕ^*, ψ^* from Equation 5.39. Note, that minimizing the above loss function leads to a self-consistent criterion which recovers the true posterior $p(\theta | \mathbf{x}_{1:N})$ over all \mathbf{x} and N under perfect convergence of both networks [133]. However, perfect convergence is often a chimera in practice, so, in a later section, we will discuss the potential sources of errors and respective remedies in detail. Having formulated our optimization criterion, we now describe the different training regimes of BayesFlow.

5.4 TRAINING PHASE

The training phase of the BayesFlow framework (left panel of Figure 5.1) can be implemented in different ways, depending on the modeling scenario and the modelers' computational budget. The starting point of all Bayesian analysis is the observed data itself. If a single observed data set $\mathbf{x}_{1:N}$ should be analyzed with a complex model that is custom-tailored for this and only this data set, it is worth considering a case-wise amortized approach, such as SNPE [60, 122]. The speed break-even point between case-wise and model-wise amortized inference is application-dependent and currently being investigated [133]. However, at present, a systematic quantitative comparison between different model classes and network architectures is missing from the literature, so modelers need to base their decisions on empirical considerations or pilot simulation studies. Be that as it may, we now present and discuss three viable simulation-based training approaches in the context of model-wise amortization with BayesFlow.

Offline learning The starting point of traditional simulation-based approaches has been the so-called *reference table* $\mathcal{D}^{(S)}$, which is simply a large data structure containing S pairs of parameters and summary statistics of synthetic observations [29, 150]. Indeed, initial machine learning approaches have already recognized the potential of using the reference table as training data for learning algorithms, such as quantile random forests [140] or deep neural networks [79], [136]. In this way, the problem of inverse inference becomes a supervised learning task which can easily be tackled with expressive learning algorithms. With BayesFlow, we can take a similar approach, as outlined in Algorithm 2.

Algorithm 2 BayesFlow training phase using offline learning

Require: f_ϕ - invertible inference network, h_ψ - algorithmically aligned summary network, S - total number of simulations, B - number of simulations per batch (batch size).

- 1: Generate a large reference table $\mathcal{D}_N^{(S)} := \{\boldsymbol{\theta}^{(s)}, \mathbf{x}_{1:N}^{(s)}\}_{s=1}^S$ using Algorithm 1.
- 2: **repeat**
- 3: Sample a mini-batch: $\mathcal{D}_N^{(B)} \sim \mathcal{D}_N^{(S)}$.
- 4: Pass each synthetic data set through the summary network: $\tilde{\mathbf{x}}^{(b)} = h_\psi(\mathbf{x}_{1:N}^{(b)})$.
- 5: Pass each pair $(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$ through the inference network: $\mathbf{z}^{(b)} = f_\phi(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$.
- 6: Compute loss according to Equation 5.40 from the training batch.
- 7: Update neural network parameters ϕ, ψ via backpropagation.
- 8: **until** convergence to $\hat{\phi}, \hat{\psi}$
- 9: **Return** trained inference and summary networks $f_{\hat{\phi}}, h_{\hat{\psi}}$.

A few points regarding Algorithm 2 are worth mentioning. First, it involves a single call to Algorithm 1 to generate the entire reference table (step 1), which will return data sets with the same size N if called only once¹. Thus, if we want to vary N during offline learning, we need to create the reference table via multiple calls to Algorithm 1 and make sure that we have an efficient data structure to store entries with different sizes. Second, steps 3 – 7 can be executed with GPU parallelization leading to a considerable speed-up in convergence. Third, the convergence criterion can be chosen as in standard deep learning application. For instance, we can establish a pre-defined number of epochs (i.e., loops through the entire training data) or an early stopping condition (i.e., if the loss does not improve in some number of consecutive epochs).

The offline learning regime is particularly useful when active calls to the simulator are computationally expensive, since data generation and training are clearly separated. It also has the advantage of reusing the simulated data and being closest to standard applications of deep learning. Obvious drawbacks of the offline learning regime are the memory demands for storing potentially large and heterogeneous data structures as well as the need to address potential overfitting.

Online learning An alternative to the offline learning regime utilizes the possibility to generate a theoretically limitless number of synthetic data sets on-the-fly. In this way, the networks never “experience” the same inputs (simulated parameters and data sets) twice, since simulations are discarded after each backpropagation update. Moreover, since classical overfitting is nearly impossible in an online learning regime, training can continue as long as the networks keep improving (i.e., the loss keeps decreasing).

Algorithm 3 outlines the online learning regime with BayesFlow. Note, that the key difference to offline training is the fact that learning and data generation are tightly intertwined when performing online learning. The most prominent advantage of online learning is also its most notable disadvantage: since simulations are not reused, the simulator needs to work actively and presents a potential bottleneck. Note also, that this detail presents less of a problem if simulations are computationally cheap or implemented efficiently (e.g., by utilizing different forms of parallelism as discussed in section 5.2.2).

¹Mathematically, the fixed N scenario is simply a special case where $p(N)$ reduces to a point mass distribution.

Algorithm 3 BayesFlow training phase using online learning

Require: f_ϕ - invertible inference network, h_ψ - algorithmically aligned summary network, B - number of simulations per iteration (batch size).

- 1: **repeat**
 - 2: Generate a mini-batch $\mathcal{D}_N^{(B)} := \{\boldsymbol{\theta}^{(b)}, \mathbf{x}_{1:N}^{(b)}\}_{b=1}^B$ using Algorithm 1.
 - 3: Pass each simulated data set through the summary network: $\tilde{\mathbf{x}}^{(b)} = h_\psi(\mathbf{x}_{1:N}^{(b)})$.
 - 4: Pass each pair $(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$ through the inference network: $\mathbf{z}^{(b)} = f_\phi(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$.
 - 5: Compute loss according to Equation 5.40 from the training batch.
 - 6: Update neural network parameters ϕ, ψ via backpropagation.
 - 7: **until** convergence to $\hat{\phi}, \hat{\psi}$
 - 8: **Return** trained inference and summary networks $f_{\hat{\phi}}, h_{\hat{\psi}}$.
-

Hybrid learning Offline and online learning represent the two endpoints on a hypothetical continuum of training strategies. However, various hybrid learning approaches appear viable for optimizing the total simulation budget available for a given modeling problem. For instance, we can use a technique used widely in reinforcement learning called *experience replay* [98, 148]. Experience replay is a hybrid learning approach aimed at balancing data usage and computational efficiency. It uses a data structure called a circular buffer which keeps past simulations in main memory and discards the oldest once its capacity has been exceeded. We outline this type of hybrid learning in Algorithm 4.

Algorithm 4 BayesFlow training phase using hybrid learning with experience replay

Require: f_ϕ - invertible inference network, h_ψ - algorithmically aligned summary network, S - memory capacity, \mathcal{F} - replay memory buffer, B - number of simulations per iteration (batch size).

- 1: Initialize replay memory buffer \mathcal{F} with capacity S .
 - 2: **repeat**
 - 3: Generate a mini-batch $\mathcal{D}_N^{(B)} := \{\boldsymbol{\theta}^{(b)}, \mathbf{x}_{1:N}^{(b)}\}_{b=1}^B$ using Algorithm 1.
 - 4: Store mini-batch $\mathcal{D}_N^{(B)}$ in memory buffer \mathcal{F} .
 - 5: Sample a mini-batch $\tilde{\mathcal{D}}_N^{(B)}$ randomly from \mathcal{F} .
 - 6: Pass each pair $(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$ through the inference network: $\mathbf{z}^{(b)} = f_\phi(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$.
 - 7: Compute loss according to Equation 5.40 from the sampled batch.
 - 8: Update neural network parameters ϕ, ψ via backpropagation.
 - 9: **until** convergence to $\hat{\phi}, \hat{\psi}$
 - 10: **Return** trained inference and summary networks $f_{\hat{\phi}}, h_{\hat{\psi}}$.
-

To further increase the efficiency when using experience replay, we can introduce a dummy parameter $\alpha \in [0, 1]$ which controls the probability of creating new simulations by executing Algorithm 1. In other words, if $\alpha = 0.5$, new parameters and synthetic observations will be generated in roughly every other pass through lines 3 – 8, thus reducing the overall number of simulations by a half.

Another hybrid learning approach utilizes a round-based strategy inspired from SNPE methods [40, 60]. Accordingly, the training phase moves through a progression of rounds and each round introduces its own simulation phase. In this way, we keep a reference table in main memory and augment it in a step-wise manner for a pre-defined number of rounds R . Thereby, each round becomes potentially longer but also reuses simulations from all previous rounds. This approach appears preferable to pure offline learning, especially when it is difficult to estimate the required number of simulations in advance. Moreover, an early stopping criterion can be grafted in-between rounds, in case further training is not conducive to the networks' performance. Algorithm 5 lays out the essential steps of the round-based approach.

Algorithm 5 BayesFlow training phase using round-based hybrid learning

Require: f_ϕ - invertible inference network, h_ψ - algorithmically aligned summary network, R - number of rounds, S - number of simulations per round, B - batch size.

- 1: Initialize reference table $\mathcal{D}^{(R \times S)} := \{\}$.
- 2: **for** $r = 1, \dots, R$ **do**
- 3: Generate synthetic data $\mathcal{D}_r^{(S)} := \{\boldsymbol{\theta}^{(s)}, \mathbf{x}_{1:N}^{(s)}\}_{s=1}^S$ using Algorithm 1.
- 4: Aggregate data: $\mathcal{D}^{(R \times S)} := \mathcal{D}^{(R \times S)} \cup \mathcal{D}_r^{(S)}$.
- 5: **repeat**
- 6: Sample a mini-batch: $\mathcal{D}_N^{(B)} \sim \mathcal{D}^{(R \times S)}$.
- 7: Pass each synthetic data set through the summary network: $\tilde{\mathbf{x}}^{(b)} = h_\psi(\mathbf{x}_{1:N}^{(b)})$.
- 8: Pass each $(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$ through the inference network: $\mathbf{z}^{(b)} = f_\phi(\boldsymbol{\theta}^{(b)}, \tilde{\mathbf{x}}^{(b)})$.
- 9: Compute loss according to Equation 5.40 from the sampled batch.
- 10: Update neural network parameters ϕ, ψ via backpropagation.
- 11: **until** convergence to $\hat{\phi}_r, \hat{\psi}_r$
- 12: **end for**
- 13: **Return** trained inference and summary networks $f_{\hat{\phi}_R}, h_{\hat{\psi}_R}$.

To sum up, one should keep an open mind regarding alternative training regimes which go beyond the ones discussed in this section. The field of neural Bayesian inference is new and, despite being an area of active research, systematic analyses of key practical issues are currently missing from the literature. As we saw in this section, there are various ways to implement the training phase of BayesFlow in practice, each coming with its own advantages and disadvantages. Ideally, the training phase should be structured so as to maximize the performance of the networks while minimizing the number of simulations. Albeit not always easy to achieve in practice, the attainment of this (informal) criterion can greatly benefit from prior considerations on computational resources and domain knowledge of the modeling problem.

5.5 INFERENCE PHASE

Once the training phase has completed, the converged BayesFlow networks can be stored on any computer and used for efficient amortized inference on any upcoming data set from the generative scope of the simulator. In other words, the summary and the inference networks have become

“domain experts” for Bayesian inference with a particular model family. Moreover, since the price of inference has been pre-paid during the upfront training phase, uncertainty-aware model inversion is now extremely efficient using the pre-trained networks. Indeed, we have extensively demonstrated the efficiency benefits of amortized inference with BayesFlow in our main paper [133]. Throughout the examples considered there, we have shown that we can obtain thousands of samples on hundreds of data sets for a couple of seconds. Algorithm 6 describes the inference phase of BayesFlow (see also Figure 5.1, right panel) on a list of I observed data sets.

Algorithm 6 BayesFlow inference phase with pre-trained networks

Require: $f_{\hat{\phi}}$ - pre-trained invertible inference network, $h_{\hat{\psi}}$ - pre-trained summary network, $\{\mathbf{x}_{1:N_i}^{(i,obs)}\}_{i=1}^I$ - list of observed data sets for inference, L - number of posterior samples.

- 1: **for** $i = 1, \dots, I$ **do**
- 2: Pass the i -th data set through the summary network: $\tilde{\mathbf{x}}^{(i,obs)} = h_{\hat{\psi}}(\mathbf{x}_{1:N_i}^{(i,obs)})$.
- 3: **for** $l = 1, \dots, L$ **do**
- 4: Sample a latent variable instance: $\mathbf{z}_l^{(i)} \sim \mathcal{N}_D(\mathbf{z} | 0, \mathbb{I})$.
- 5: Evaluate the inference network in reverse: $\boldsymbol{\theta}_l^{(i)} = f_{\hat{\phi}}^{-1}(\mathbf{z}_l^{(i)}; \tilde{\mathbf{x}}^{(i,obs)})$.
- 6: **end for**
- 7: Store $\{\boldsymbol{\theta}_l^{(i)}\}_{l=1}^L$ as samples from the i -th posterior $p(\boldsymbol{\theta} | \mathbf{x}_{1:N} = \mathbf{x}_{1:N_i}^{(i,obs)})$.
- 8: **end for**

Note, that all components of Algorithm 6 can also benefit from a tremendous speed-up with the aid of GPU acceleration. In particular, both loops over I and L can be performed in parallel using a GPU. Thus, it seems evident that every step of a BayesFlow analysis pipeline is amenable to modern parallel computing, from Monte Carlo simulations to inference on real data (and also validation, as we will discuss shortly). The correctness of Algorithm 6 is guaranteed under the conditions of perfect convergence and self-consistency, that is:

$$f_{\phi^*}^{-1}(\mathbf{z}; h_{\psi^*}(\mathbf{x}_{1:N})) \sim p(\boldsymbol{\theta} | \mathbf{x}_{1:N}) \quad \text{with} \quad \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | 0, \mathbb{I}) \quad (5.41)$$

where ϕ^* , ψ^* are global minimizers of the modified criterion (Equation 5.39) and the simulation gap induced by modeling $p^*(\mathbf{x}_1, \dots, \mathbf{x}_N)$ via $p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \int_{\Theta} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$ is negligible (see [133] for a detailed proof). Moreover, the samples obtained by perfectly converged BayesFlow networks are fully independent, in contrast to MCMC and other stateful Bayesian methods which sometimes induce severe auto-correlation among successive samples. In practice, however, it is important to be aware of potential deficiencies in computational faithfulness, to which we turn next.

5.6 SOURCES OF ERROR

Computational faithfulness refers to the adequacy or the ability of a Bayesian method to recover the correct target posterior in a particular (simulated or real-world) modeling scenario. Thus, computational faithfulness is not just a nice-to-have extra, but a crucial prerequisite for trustworthy model-based inference. No Bayesian method is exempt from the privilege of occasionally leading

modelers and decision makers astray. Therefore, even though each method will eventually err in some (inevitably unexpected) situation, it seems important to at least have a handy catalogue of errors, listed together with their potential causes and fixes. Such a catalogue does not have to be static, but can dynamically grow as a particular method is continuously used in novel applications or integrated in existing analysis pipelines. In the following, we discuss five prominent sources of error which can potentially compromise faithful Bayesian inference with BayesFlow.

The first source of error is the *simulation gap* which can occur under model misspecification or when the observed data are contaminated in ways not covered by the stochastic component ξ of the simulator. Despite being an issue which needs to be addressed via *prior predictive checks*, that is, *before doing inference*, errors due to model misspecification will result in incorrect posteriors that might be hard to detect in practice. In some cases, model misspecification might manifest itself in posteriors which are incompatible with the prior (e.g., posterior samples having 0 density under the prior), but more complex misbehavior is also possible. In other cases, researchers might anticipate how data will be contaminated (e.g., inattention by participants in an experiment or guesswork during a performance test) and explicitly model the contaminants². However, in most cases, model misspecification will be far from obvious (otherwise one would have taken steps to eliminate it), so its potential to bias subsequent inference remains a real issue. This underlines the importance of domain expertise consistency when setting up a model and highlights the fact, that all steps in a Bayesian workflow are inter-dependent, with errors inherent in initial phases tacitly propagating to further phases of data analysis.

The second source is the Monte Carlo error introduced by necessarily using a finite number of simulations from the joint model $p(N, \theta, \mathbf{x})$ to approximate the expectation in Equation 5.39. This source is also referred to as *approximation error* and is a widely accepted concomitant of all Monte-Carlo methods. It is also relatively easy to mitigate in an online learning regime, since, in principle, we can run the simulator as long as we can afford and thus generate a potentially infinite amount of training data. In this respect, neural simulation-based inference is in a better position to exploit the capacity of data-hungry deep neural networks than more prototypical deep learning applications operating a limited-data regime.

The third source is the *amortization gap* which refers to a potential deficiency in the inference phase due to the use of a single set of summary and inference networks parameters $(\hat{\psi}, \hat{\phi})$ to perform inverse inference globally (i.e., to obtain model-wise amortization). An amortization gap can be elusive and non-trivial to detect with certainty in practical scenarios unless one performs case-wise inference alongside (which would be wasteful in practice) and quantifies the quality of both analyses. Sometimes, an amortization gap can be detected via probabilistic calibration methods (e.g., simulation-based calibration, SBC, [155]), although the reasons for miscalibrated inference might be obscure at first. The more severe problem with this approach, however, is that miscalibrated inference might have different and overlapping causes, and thus not be directly attributable to an amortization gap. To make matters worse, perfectly calibrated inference on the basis of simulations might still be perfectly miscalibrated when transferred to real data, so, detrimentally, an amortization gap can manifest effects similar to a simulation gap. Thus, proper posterior model checking is needed to ensure a model's generative and predictive performance meet the modeler's needs. In the case of a determinable amortization gap, moving to case-wise amortization might be

²In fact, this is what we did in our application of BayesFlow to Covid-19 modeling [135]

a viable option when dealing with complex models. Alternatively, increasing the expressiveness of the summary and inference networks could also ameliorate amortization-related problems (to be discussed shortly).

The fourth source is due to a summary network which may not fully capture the relevant information in the data or when sufficient summary statistics do not exist. All things being equal, not capitalizing on the information contained in the data will result in incorrect inference, usually in the form of overdispersed or otherwise miscalibrated posteriors [133]. Thus, the choice of summary network is a crucial proviso for the overall performance of a BayesFlow application. Indeed, the design and architecture of optimal summary networks is a subject of ongoing research. And even though concrete guidelines for optimal summary network design are currently lacking, there are at least two wells of guidance. On the one hand, recent work on *probabilistic symmetry* [13] and *algorithmic alignment* [174] can provide theoretical ideas on how to select a suitable summary architecture for a particular problem. On the other hand, recent simulation-based applications using the BayesFlow framework to tackle complex stochastic models in different research domains can provide viable empirical hints for aligning the summary network to the data at hand. Currently, the BayesFlow method has been employed to perform inference on complex models from psychology [172], cognitive science [137], computational psychiatry [31], epidemiology [135], mathematical finance [147], and physics [10]. Nevertheless, more theoretical and empirical work is needed for definite recommendations at the current stage of development.

The fifth source is due to an inference network which does not accurately transform the true posterior into the prescribed (Gaussian) latent space. This error can be easily detected by passing multiple simulations through the networks and exploring the structure of the latent space $p(\mathbf{z})$. Industrious modelers might even consider computing a formal metric between the desired latent space (e.g., Gaussian) and the one obtained by the networks. In the presence of a mismatch, increasing the capacity of the inference network should be the first step to take before further investigations into the problem. Accordingly, both the depth (number of coupling layers) of the cINN as well as the design of the coupling layers themselves could be tuned to increase the expressiveness of the learned transformation from θ -space to \mathbf{z} -space. The benefits of neural network depth have been confirmed both in theory and in practice [5, 97], so one should expect better performance in complex settings with increasing network depth. However, one should also bear in mind, that an underexpressive summary network could also be responsible for a deficient transformation, since summary and inference network are optimized jointly during the training phase of BayesFlow. Thus, an exclusive focus on the inference network might not be conducive to solving all possible transformation errors. In any case, visualizing the learned latent space and inspecting it for deviations from the desired one (i.e., as prescribed by the optimization criterion) is integral to any application of BayesFlow.

To sum up, as in any Bayesian framework, care should be taken to ensure computational faithfulness as a basis for reliable amortized inference with BayesFlow. Fortunately, we can address model misspecification (error 1) with standard Bayesian prior/posterior predictive checks [52, 56, 144]. Moreover, we can establish deficiencies in self-consistency (errors 2-5) by simply visualizing the latent space obtained in any application of BayesFlow, which provides us with a self-diagnostic method. Naturally, using this method does not help us pinpoint the exact source of error, but only indicates its potential presence. As previously discussed, certain heuristics can be applied for a more detailed error checking in particular applications. In addition, future research should

take steps towards a more fine-grained theoretical error analysis, elucidating the consequences of imperfect convergence and investigating error bounds.

5.7 A BAYESIAN WORKFLOW WITH BAYESFLOW

We will now briefly discuss the place of BayesFlow in a principled Bayesian workflow with a focus on cognitive modeling [144]. In the context of a single cognitive model, a principled Bayesian workflow proposed by [144] goes through the following steps:

1. Prior predictive checks
2. Computational faithfulness checks
3. Model sensitivity checks
4. Posterior predictive checks

Prior predictive checks are designed to test whether a model is consistent with the relevant domain expertise. Computational faithfulness refers to the accuracy of the estimation method. Model sensitivity asks whether the parameters of a model can be recovered given the model’s prior specification, generative scope, and algorithmic form. Finally, posterior predictive checks assess whether the model captures the relevant structure of the assumed true data generating process. Needless to say, these steps are all computationally intensive and associated with their own specific challenges. In the following, we describe the significant role of amortized inference with BayesFlow at each step of the Bayesian workflow.

5.7.1 PRIOR CONSISTENCY

Since no inference happens at the (pre-data) stage of ensuring consistency with domain expertise, there is little room for amortized inference either. However, prior predictive checks should be an integral part of any Bayesian (simulation-based) analysis. Inconsistent models can require either a modification of the prior $p(\boldsymbol{\theta})$ or/and the simulator $g(\boldsymbol{\theta}, \boldsymbol{\xi})$, in order to resolve conflicts with self-evident domain expertise. Ideally, cognitive models should (re-)produce meaningful patterns of human behavior and not harness pathological patterns in their generative scope (e.g., superhuman reaction times or flawless memory). The easiest way to control for inconsistent model behavior is to constrain the priors to meaningful domains (numerical ranges). In other cases, incorporating certain constraints into the simulator and extensive exploration of the data space (e.g., via *prior pushforward checks*) might be necessary.

5.7.2 COMPUTATIONAL FAITHFULNESS

Computational faithfulness is best ensured when performing Bayesian inference with methods capable of self-diagnosis. For example, convergence issues in MCMC sampling methods in general can be detected by inspecting the Gelman-Rubin (\hat{R}) metric [55] or specific problems with Hamiltonian Monte Carlo (HMC) can be indicated by divergent transitions [8].

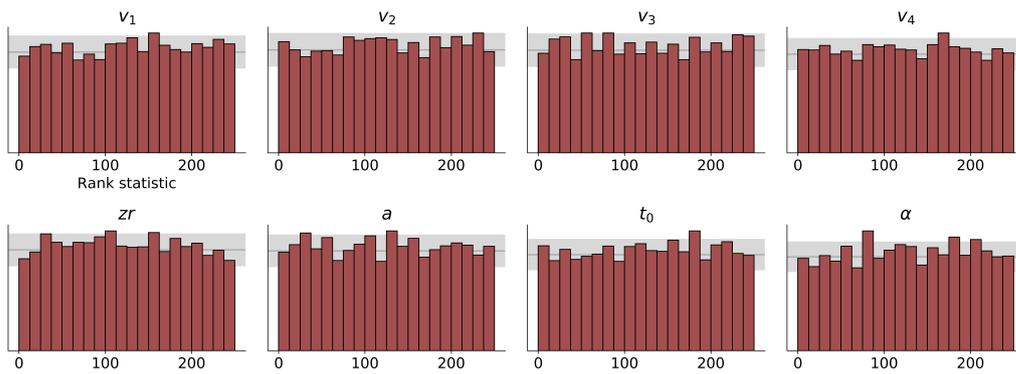


Figure 5.5: Simulation-based calibration (SBC) results for a Lévy flight model with 8 parameters at $N = 800$ trials as a validation check for computational faithfulness. The histograms indicate no systematic deviations from uniformity across marginal posteriors.

A natural self-diagnostic of BayesFlow can be derived by inspecting its ability to correctly transform $p(\boldsymbol{\theta} | \boldsymbol{x})$ into $p(\boldsymbol{z})$ for any \boldsymbol{x} . To ensure this, one can simulate a set of pairs $(\boldsymbol{\theta}, \boldsymbol{x})$, pass them through a converged BayesFlow configuration and inspect the resulting latent space for deviations from the prescribed latent space (a spherical Gaussian in our case). This can be done either visually, or numerically, by computing, for instance, the maximum mean discrepancy (MMD, [61]). Note, that this procedure is very fast, since it requires only simulations and forward evaluations of the network, which can all be performed in parallel and furthered through GPU acceleration.

Alternatively, one can resort to calibration algorithms, which can reveal systematic biases in the approximate posteriors. One such approach is simulation-based calibration (SBC, [155]), which is a variant of probabilistic calibration [58] specifically tailored for generative Bayesian models. SBC can be used to validate the inferential correctness of a Bayesian sampling method without knowing the true posterior distribution, which makes it a very powerful diagnostic tool.

However, SBC is extremely time-intensive with standard Bayesian methods, since the computational model needs to be estimated repeatedly, potentially hundreds of times, on different simulated data sets. In addition, the obtained posterior samples should be independent for SBC to yield reliable results, which further increases the required computing time to eliminate auto-correlation via thinning while still retaining enough posterior samples afterwards [155]. These requirements often render SBC practically infeasible for non-amortized Bayesian methods.

Within the BayesFlow framework, SBC can be performed with extreme efficiency once the training phase is over. It simply requires running Algorithm 6 repeatedly with simulated data sets instead of actual observations. Amortized inference ensures that these runs are very efficient. In addition, a perfectly converged BayesFlow configuration yields independent samples from the posterior. Using GPU acceleration, SBC with BayesFlow typically takes a couple of seconds, assuming that the synthetic observations have already been simulated. Thus, we advise the routine and automated use of SBC when doing amortized Bayesian inference.

5.7.3 MODEL SENSITIVITY

Model sensitivity, or model adequacy, refers to the feasibility of inverse inference. In other words, it asks about the amount of information gained through Bayesian updating, assuming computational faithfulness of the inferential method and self-consistency of the Bayesian simulator. A straightforward way to obtain a measure of model sensitivity is to compute the expected Bayesian surprise (see Section 3.2), which can also be used for model comparison in a pre-data stage. However, since Bayesian surprise could be hard to interpret in practical terms and without reference to information-theoretic notions, one can resort to other proxies of information gain, such as posterior contraction or posterior z -score [144].

Posterior contraction is a measure of sharpness achieved by Bayesian updating and can be computed for both marginal as well as joint distributions (see Section 3.2). Higher values indicate a high degree of uncertainty reduction and, equivalently, a noticeable posterior sharpness. Likewise, the posterior z -score is a measure of accuracy computed as the difference between the posterior mean (expected value) and the true parameter configuration of a simulated data set, standardized by the posterior variance. Accordingly, smaller values suggest that the posterior concentrates strongly around the true parameter (i.e., the posterior mean is a reasonable representation of the full posterior) while larger values suggest a posterior that concentrates in other parts of the prior domain.

In order to avail themselves of posterior contraction and posterior z -score as useful measures of model sensitivity, modelers need to simulate multiple data sets from the generative model, perform inverse inference on all of them, and compute the corresponding metrics. Similarly to SBC, the feasibility of this procedure depends heavily on the efficiency of the Bayesian estimation method. Thus, evaluating model sensitivity with non-amortized methods might turn out to be prohibitively slow, whereas it becomes trivial when doing amortized inference with a pre-trained BayesFlow configuration. The same would be true for *any* measure of model sensitivity requiring repeated inverse inference on multiple simulated data sets, so model sensitivity is another step of a principled Bayesian workflow which can massively profit from amortized inference.

5.7.4 POSTERIOR PREDICTIVE CHECKS

Posterior predictive checks are vital for evaluating a computational model on actually observed data with respect to the model’s generative and predictive performance. Moreover, posterior predictive metrics, such as cross-validation or Bayesian information criteria, can be used for subsequent model comparison and selection in a multi-model setting. As already discussed in Section 3.5, posterior predictive checks comprise a serious computational bottleneck in Bayesian pipelines, even more so when dealing with intractable models.

For instance, k -fold or leave-one-out (LOO) cross-validation (CV) require re-estimating the same model on multiple sub-sets of the original data set in order to approximate out-of-sample predictive performance. When multiple data sets are to be modeled, the computational load increases in a multiplicative manner with the number of data sets B , so extensive posterior predictive checks with standard Bayesian methods quickly become too costly to perform. Once again, amortized inference with BayesFlow offers considerable efficiency gains, since repeated applications of the same model simply involve running the pre-trained networks in a feed-forward mode with different (sub)-sets of the full data.

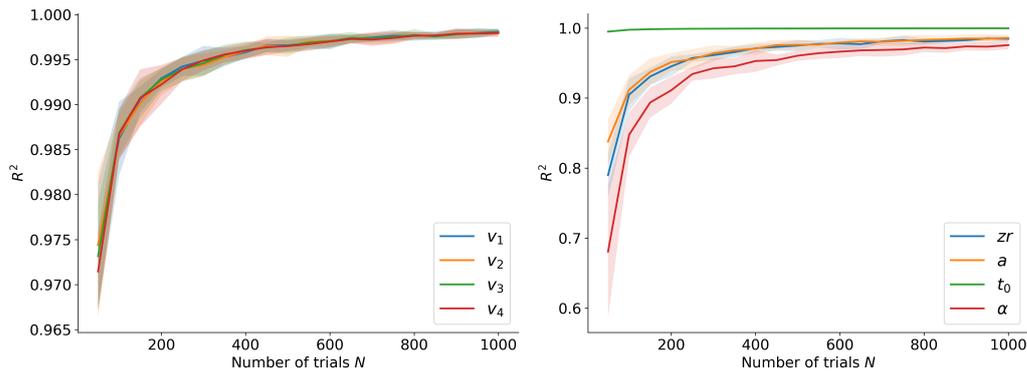


Figure 5.6: The left panel depicts parameter recovery of the four drift rate parameters as a function of trial numbers N using the R^2 metric between true and estimated values. The right panel depicts recovery of the other four parameters. Posterior means are used as summaries of the full posteriors and shaded regions represent bootstrap 95% confidence intervals.

5.8 A QUICK DEMONSTRATION

As an illustrative example, we present an application of BayesFlow to a recent intractable evidence accumulation model (EAM). Further applications to models from different research domains are described in Chapter 8 or in applied works [10, 31, 135, 147]. EAMs are a popular class of mechanistic models in psychology and cognitive science, since they enable a principled model-based analysis of human response time (RT) data obtainable in controlled experimental environments.

For this example, we focus on a Lévy flight model (LFM) with a non-Gaussian noise assumption [169, 172]. The Lévy flight process is driven by the following stochastic ordinary differential equation (ODE):

$$dx_c = v dt_c + \xi dt^{1/\alpha} \quad (5.42)$$

$$\xi \sim \text{AlphaStable}(\alpha, 0, 1, 0) \quad (5.43)$$

where dx_c denotes accumulated cognitive evidence in condition $c \in \{1, 2, 3, 4\}$, v_c denotes the average speed of information processing (drift) in condition c , and α controls the heaviness of the noise distribution's tails (i.e., smaller values increase the probability of outliers in the accumulation process).

Consider first a simple question of optimal experimental design. A behavioral researcher wants to conduct a response times (RT) experiment with four conditions and model performance via the Lévy flight model. How many trials are needed for accurate parameter recovery? To answer these questions, we can simulate multiple experiments with varying number of trials N per synthetic participant and then compute some practically relevant discrepancy between ground-truth parameters and their estimates. Afterwards, we can quantify computational faithfulness and model sensitivity with the particular number of trials N collected in the experiment. Note, that the mandatory prior predictive and posterior predictive checks are left out for conciseness of exposition.

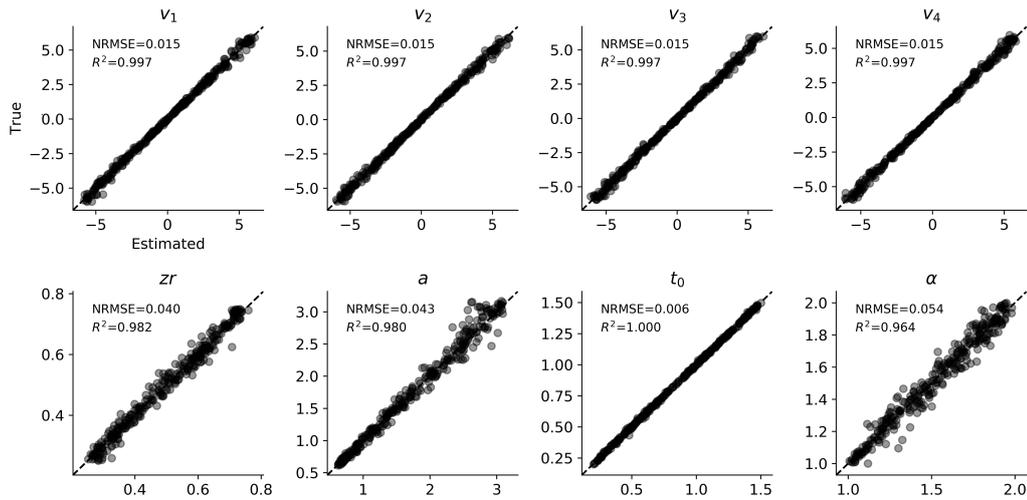


Figure 5.7: Parameter recovery (true vs. estimated) values for $N = 800$ simulated trials. Normalized root-mean-square error ($NRMSE$) and the coefficient of determination (R^2) are used to quantify discrepancy between posterior means and true parameter values.

Since the Lévy flight model is analytically intractable, such a simulation scenario is not feasible with non-amortized methods, which would need weeks on standard machines [169]. However, using a BayesFlow architecture, we can obtain an amortized neural sampler capable of working with variable number of trials (i.e., by using a permutation invariant summary network). The on-line training phase with Algorithm 3 took approximately one day on a standard laptop equipped with an NVIDIA[®] GTX1060 graphics card. Subsequent inference is then extremely efficient, as amortized Bayesian estimation on 500 simulated participants takes less than two seconds [137].

We visualize the results by plotting the average R^2 metric obtained from estimating the Lévy flight model on 300 simulated participants with N varying between 50 and 1000 (cf. Figure 5.6). Notably, recovery of the ground-truth parameters via posterior means is nearly perfect at higher trial numbers, and resembles a logarithmic function of N^3 . A similar plot can be created for posterior contraction as a function of N (see Ricker example in [133]).

Further, we can now apply the same network from the previous simulation example for executing fully Bayesian inference on real data. For this illustrative example, we estimate the Lévy flight model from eleven participants performing a long lexical decision task with $N = 800$ trials per condition [137]. Since the task had a 2×2 design, with a factor for *difficulty* (hard vs. easy), and a factor for *stimulus type* (word vs. non-word), we assume a different drift rate v_c for each design cell $c \in \{1, 2, 3, 4\}$.

Before performing inference on actually observed data, we compute SBC and evaluate parameter recovery using simulations with $N = 800$ trials per condition (aligned to the particular experimental design) in order to become a rough sense of computational faithfulness and model sensitivity. Importantly, these checks were performed within seconds with amortized inference and

³Strictly speaking, one should also ensure that inference is calibrated for each N , a step which is no more computationally expensive with BayesFlow and which we omit here for brevity.

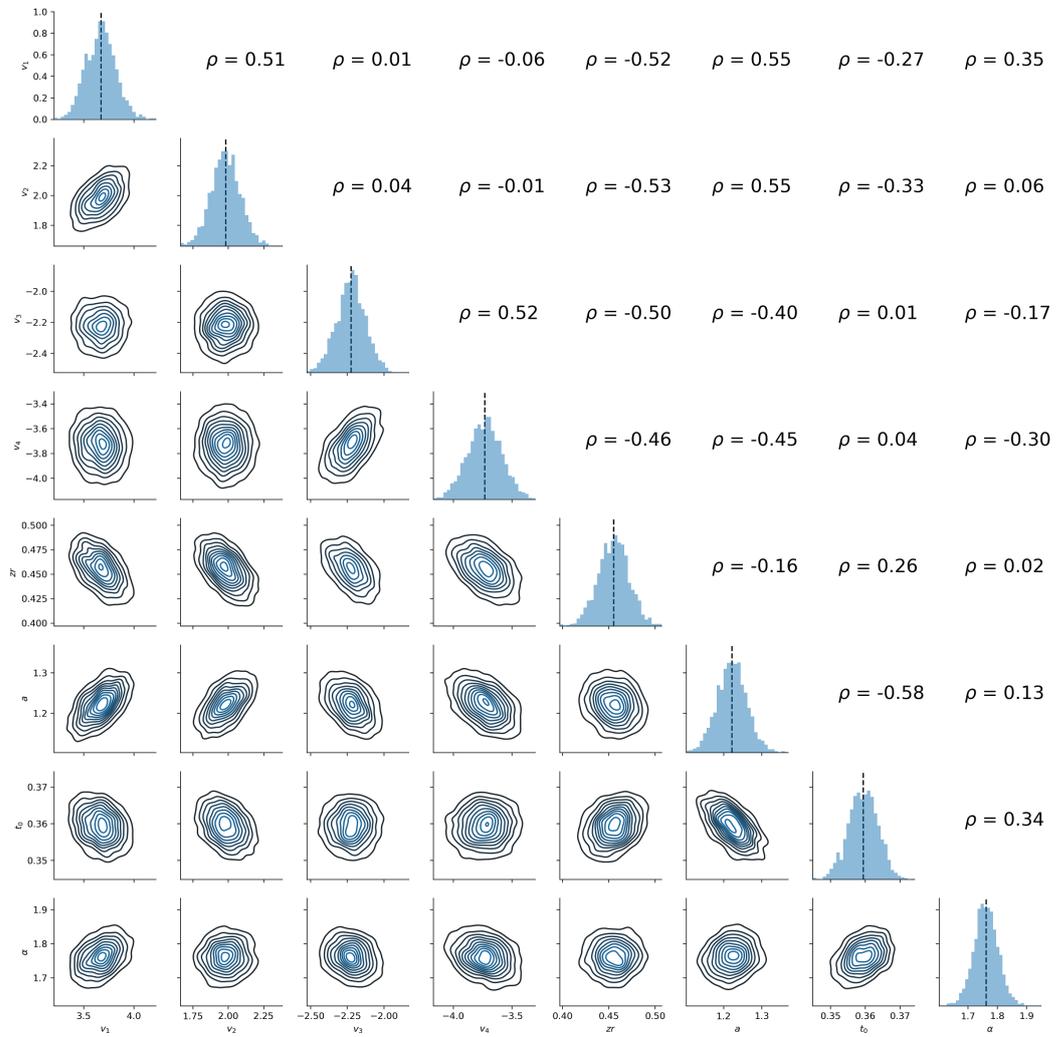


Figure 5.8: Individual bivariate posteriors obtained from data of one example participant in the lexical decision task.

would have been intractable with standard methods. Accordingly, marginal SBC and parameter recovery plots are depicted in Figures [Figure 5.5](#) and [Figure 5.7](#), respectively. The SBC histograms suggest no systematic biases across the approximate marginal posteriors for each parameters (e.g., no under- or overdispersion of the true posterior). Likewise, the recovery plots indicate excellent parameter recovery using posterior means as summaries of the full posteriors, a result which is also evident from the earlier [Figure 5.6](#).

Thus, we can interpret the results from these pre-data checks as hints of intact self-consistency and proceed to applying BayesFlow to real data. A typical output from applying BayesFlow to a single data set is depicted in [Figure 5.8](#), which presents marginal and bivariate posteriors. The latter allows us to visually inspect posterior correlations as indicators of disentanglement (linear independence) between the individual model parameters [137].

5.9 CONCLUDING REMARKS

This chapter introduced the building blocks of our BayesFlow framework and discussed its mathematical and algorithmic formulation at a relatively high level. More details regarding perfect convergence, training, and hyperparameter choice (e.g., learning rate, optimizer settings) can be found in our methodological work [133] as well as in the applied works [10, 31, 135, 147]. Details regarding implementation as well as templates for parameter estimation are also available at the corresponding code repository (<https://github.com/stefanradev93/BayesFlow>). Whereas a multitude of features and potential improvements remain to be explored in future research, our results from initial simulations and applications appear highly promising. Thus, we hope that our framework will accelerate model-based inference in a variety of scientific fields and prove its utility beyond the current applications.

6 AMORTIZED MODEL COMPARISON

When you have eliminated the impossible, whatever remains, however improbable, must be the truth.

— Arthur Conan Doyle

Researchers from various scientific fields face the task of selecting the most plausible theory for an empirical phenomenon among multiple competing theories. Theories in the cognitive and behavioral sciences are also not exempt from being subject to a relentless selection process. As already discussed, rigorous theories are often instantiated as formal models which describe how observable quantities arise from unobservable parameters in the language of mathematics. Focusing on the level of mathematical models, the problem of theory selection then becomes one of *model selection*.

For instance, neuroscientists might be interested in comparing different models describing the spiking patterns revealed by *in vivo* recordings of neural activity [76]. Epidemiologists, on the other hand, might consider different dynamic models for predicting the transmission rate or other characteristics of an unfolding infectious disease [167]. Crucially, the preference for one model over alternative models in these examples can have important consequences for research projects or social policies.

Accounting for complex natural phenomena often requires specifying complex models which entail some degree of randomness. Inherent stochasticity, incomplete description, or epistemic ignorance all call for some form of uncertainty awareness. As a further complication, empirical data on which models are fit are necessarily finite and can only be acquired with finite precision. Finally, the plausibility of many non-trivial models throughout various branches of science can be assessed only approximately, through rather costly simulation-based methods [26, 35, 76, 104, 139, 161].

Our evidential method aims to amortize Bayesian model comparison by combining latest ideas from simulation-based inference and uncertainty quantification for building efficient and uncertainty-aware neural classifiers. As such, it is intended to complement the toolbox of simulation-based methods for parameter estimation with crucial model comparison capabilities. Moreover, it incorporates a unique feature for estimation of higher-order uncertainty, which goes beyond the scope of standard ABC methods.

6.1 DESIDERATA

In the previous chapter, we introduced the nuts and bolts of our BayesFlow method for simulation-based Bayesian parameter estimation. This chapter will present our complementary framework for simulation-based Bayesian model comparison. The next chapter will discuss the potential and

challenges inherent in combining both frameworks into a single meta-framework. As with parameter estimation, we begin by stating our desiderata for building a useful model comparison method:

1. Estimated model probabilities should be, at least in theory, calibrated to the true model probabilities induced by an empirical problem
2. Estimated model probabilities should be accurate even for finite or small sample sizes
3. Preference for simpler models (i.e., the probabilistic Occam’s razor) should be encoded by the estimated model probabilities
4. The method should be applicable to complex models with implicit likelihoods within reasonable time limits
5. The method should enable full amortization over arbitrarily many models, data sets, and varying data set sizes
6. The method should automatically extract maximal information from the raw data and avoid information loss through insufficient summary statistics of the data

Evidently, the desiderata for model comparison are somewhat overlapping with those stated earlier for parameter estimation. Indeed, in this chapter, we will reuse many of the previous concepts for building algorithmically aligned *summary networks* in the BayesFlow framework.

6.2 BACKGROUND

The following section will briefly rehearse some of the core concepts related to Bayesian model comparison (see also Chapter 3), thereby setting the stage for the derivation of our evidential framework.

6.2.1 BAYESIAN MODEL COMPARISON

In Bayesian modeling, we typically start with a collection of J competing generative models, which we denoted as $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_J\}$. Each abstract model index \mathcal{M}_j is associated with a generative mechanism g_j , typically realized as a Monte Carlo simulation program, and a corresponding parameter space Θ_j equipped with a prior distribution $p(\theta_j | \mathcal{M}_j)$. Ideally, each g_j represents a theoretically plausible stochastic mechanism by which observable behavior \mathbf{x} arises from hidden time-invariant parameters θ_j and independent noise ξ :

$$\mathbf{x}_n = g_j(\theta_j, \xi_n) \text{ with } \theta_j \in \Theta_j \tag{6.1}$$

where Θ_j is the corresponding parameter space of model g_j and the subscript j explicates that each model might be specified over a different parameter space¹. We assume that the functional or algorithmic form of each g_j is known and that we have a sample (data set) $\{\mathbf{x}_i\}_{i=1}^N := \mathbf{x}_{1:N}$ of N

¹Also the noise distribution $p(\xi)$ and noise space Ξ might differ for each model, but we will keep this possibility implicit.

(multivariate) observations generated from an unknown process p^* . The task of Bayesian *model comparison* is to assign a plausibility score (e.g., a posterior probability) to each of the models in \mathcal{M} . The task of Bayesian *model selection* is then to choose the model in \mathcal{M} that best describes the observed data by balancing simplicity (sparsity) and predictive performance.

As already discussed in Chapter 3, Bayesian methods for model comparison can be categorized as either posterior predictive or prior predictive approaches [52], with our method falling into the latter category. Posterior predictive approaches are concerned with predicting upcoming observations using models extracted from the available data. In prior predictive approaches, models are conditioned only on prior information but not on the available data. Accordingly, all available data counts as new data for the purpose of prior predictive methods.

To recapitulate, the canonical measure of prior predictive performance is the *marginal likelihood*:

$$p(\mathbf{x}_{1:N} | \mathcal{M}_j) = \int_{\Theta_j} p(\mathbf{x}_{1:N} | \boldsymbol{\theta}_j, \mathcal{M}_j) p(\boldsymbol{\theta}_j | \mathcal{M}_j) d\boldsymbol{\theta}_j \quad (6.2)$$

which forms the basis for the computation of Bayes factors and posterior odds between pairs of competing models. If two models are equally likely *a priori*, the posterior odds equal the Bayes factor. Furthermore, if the Bayes factor, or, equivalently, the posterior odds equal one, the observed data provide no decisive evidence for one of the models over the other. However, a relative evidence of one does not distinguish whether the data are equally likely or equally unlikely under both models, as this is a question of absolute evidence. Needless to say, the distinction between relative and absolute evidence is of paramount importance for model comparison, so we address it in the next section on model comparison frameworks.

6.2.2 \mathcal{M} -FRAMEWORKS

Closely related to the distinction between relative and absolute evidence is the distinction between \mathcal{M} -closed and \mathcal{M} -complete frameworks [176]. Under an \mathcal{M} -closed framework, the true model is assumed to be in the predefined set of competing models \mathcal{M} , so relative evidence *is* identical to absolute evidence. Under an \mathcal{M} -complete framework, a true model is assumed to exist but is not necessarily assumed to be a member of \mathcal{M} . However, one still focuses on the models in \mathcal{M} due to computational or conceptual limitations².

Deciding on the particular \mathcal{M} -framework under which a model comparison problem is tackled is often a matter of prior theoretical considerations. However, since in most non-trivial research scenarios \mathcal{M} is a finite set and candidate models in \mathcal{M} are often simpler approximations to the true model, there will be *uncertainty* as to whether the observed data could have been generated by one of these models. In the following, we will refer to this uncertainty as *epistemic uncertainty*. Our method utilizes a data-driven way to calibrate its epistemic uncertainty in addition to the model probabilities through simulations under an \mathcal{M} -closed framework.

Consequently, given real observed data, a researcher can obtain a measure of uncertainty with regard to whether the generative model of the data is likely to be in \mathcal{M} or not. From this perspective, our method lies somewhere in the middle ground between \mathcal{M} -closed and an \mathcal{M} -complete framework as it provides information from both viewpoints.

²See also [176] for discussion of an \mathcal{M} -open framework, in which no true model is assumed to exist.

6.2.3 MODEL SELECTION AS CLASSIFICATION

In line with previous simulation-based approaches to model comparison (e.g., ABC), we will utilize the fact that we can simulate arbitrary amounts of data from each simulator g_j (to be described shortly). Following previous machine learning approaches to model selection [104, 132], we reinterpret the problem of model comparison as a probabilistic classification task. In other words, we seek to learn a mapping $f : \mathcal{X}^N \rightarrow \Delta^J$ from an arbitrary data space \mathcal{X}^N to a probability simplex Δ^J containing the multinomial posterior model probability $p(\mathcal{M} | \mathbf{x}_{1:N})$. Previously, different learning algorithms, such as random forests have been employed to tackle model comparison as classification [104]. Reusing the ideas from algorithmic alignment and probabilistic symmetry incorporated into the BayesFlow framework, our method parameterizes f_η via a specialized neural network with trainable parameters η which is aligned to the probabilistic structure of the generative models (i.e., a permutation invariant network for memoryless models or a recurrent network for stateful models).

In addition, our method differs from previous classification approaches to model comparison in the following aspects. First, it requires no hand-crafted summary statistics, since the most informative summary statistics are learned directly from data. Second, it can make use of online learning (i.e., on-the-fly simulations) which requires no storage of large reference tables or data grids. Third, the addition of new competing models does not require changing the architecture or re-training the network from scratch, since the underlying data domain remains the same. In line with the transfer learning literature, only the last layer of a pre-trained network needs to be changed and training can be resumed from where it had stopped. Last, our method is uncertainty-aware, as it returns a higher-order distribution over posterior model probabilities. From this distribution, one can extract both absolute and relative evidences, as well as quantify the model selection uncertainty implied by the observed data (more on this distinction later).

Intuitively, a converged evidential network encodes the probabilistic relationship between data and models through the network’s weights. Thus, once trained, the evidential network can be reused to perform instant model comparison on multiple real observations. As mentioned above, the addition of new models requires simply adjusting the pre-trained network, which requires much less time than re-training the network from scratch.

6.2.4 MULTI-MODEL FORWARD INFERENCE

Our evidential methods requires the ability to implement each candidate model as a simulator and efficiently generate synthetic observations from each model. This process amounts to performing forward inference in a multi-model context and is described in detail in Algorithm 7. Since we only need the simulated data sets and the corresponding model indices, we can run Algorithm 7 repeatedly to construct training batches of the form $\mathcal{D}_N^{(B)} := \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ with B simulated data sets of size N and B corresponding one-hot encoded model indices. We can then feed each batch to a specialized neural network which takes as input simulated data with variable sizes and returns a distribution over posterior model probabilities. Note, that similar considerations regarding computational efficiency and parallelism apply as previously discussed in the context of parameter estimation with BayesFlow.

Algorithm 7 Monte Carlo generation of synthetic data sets for model comparison

Require: $p(\mathcal{M})$ - prior over models, $\{p(\boldsymbol{\theta} | \mathcal{M}_j)\}$ - list of priors over model parameters, $\{g_j\}$ - list of stochastic simulators, $p(\boldsymbol{\xi})$ - noise distribution, $p(N)$ - distribution over data set sizes, B - number of data sets to generate per iteration (batch size).

- 1: Draw data set size: $N \sim p(N)$.
- 2: **for** $b = 1, \dots, B$ **do**
- 3: Draw model index from model prior: $\mathcal{M}_j^{(b)} \sim p(\mathcal{M})$.
- 4: Draw model parameters from prior: $\boldsymbol{\theta}_j^{(b)} \sim p(\boldsymbol{\theta}_j | \mathcal{M}_j^{(b)})$.
- 5: **for** $n = 1, \dots, N$ **do**
- 6: Sample noise instance: $\boldsymbol{\xi}_n \sim p(\boldsymbol{\xi})$.
- 7: Run simulator j to obtain n -th synthetic observation: $\mathbf{x}_n = g_j(\boldsymbol{\theta}_j^{(b)}, \boldsymbol{\xi}_n)$.
- 8: **end for**
- 9: Encode model index as a one-hot-encoded vector: $\mathbf{m}^{(b)} = \text{OneHotEncode}(\mathcal{M}_j^{(b)})$.
- 10: Store pair $(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})$ in data structure $\mathcal{D}_N^{(B)}$.
- 11: **end for**
- 12: **Return** mini-batch $\mathcal{D}_N^{(B)} := \{\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)}\}_{b=1}^B$.

6.3 TRAINING AMORTIZED EVIDENCE APPROXIMATORS

In our model comparison framework, evidential neural networks learn a higher-order uncertainty representation over the model posterior $p(\mathcal{M} | \mathbf{x}_{1:N})$ in a simulation-based manner. As we did in the BayesFlow framework, we split model comparison into two phases: an expensive training/simulation phase, in which neural network parameters are optimized via standard backpropagation; and a cheap inference phase, in which a pre-trained evidential network is applied to an arbitrary amount of real data sets $\mathbf{x}_{1:N}^{(obs)}$. The output of an evidential network is a higher-order distribution, from which we can obtain a vector of probabilities $p_\eta(\mathbf{m} | \mathbf{x}_{1:N}^{(obs)})$ which approximates the true model posterior $p(\mathcal{M} | \mathbf{x}_{1:N}^{(obs)})$ for any observable $\mathbf{x}_{1:N}^{(obs)}$. In addition, we can obtain local uncertainty information which serves as a proxy for absolute evidence.

6.3.1 EVIDENCE REPRESENTATION

How can we obtain a measure of absolute evidence by considering only a finite number of competing models? Indeed, such an undertaking has the appearance of an ill-posed problem from the very offset. Our approach will be to re-frame the problem as teaching a neural network to respond with *I don't know* when faced with data which could not have been generated by one of the models (i.e., has not been experienced during the simulation-based training phase). In general, however, a purely probabilistic approach is not well-suited for representing a lack of knowledge [74], since even the uniform distribution encodes the belief in *equally likely events*. In contrast, meta-probabilistic approaches propose to use second-order probabilities [80, 146] for representing the absence of any definite knowledge. In a Bayesian setting, we typically lack knowledge regarding the misspecification degree of the candidate models. Thus, our framework can also be viewed as an approach to quantifying model misspecification via higher-order uncertainty. In this way,

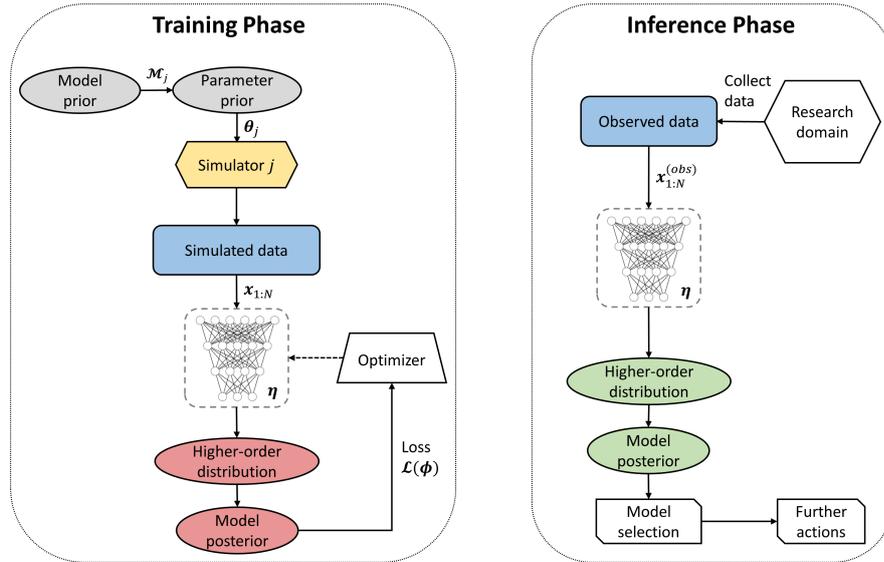


Figure 6.1: Both phases of our evidential model comparison framework. Left panel: During the training phase, an algorithmically aligned evidential network is trained jointly with random draws from the joint prior $p(\mathcal{M}, \theta)$ and synthetic data from the simulator; Right panel: During the inference phase, no training or optimization happens in this phase. The upfront training effort amortizes over an arbitrary number of models, observations and data sets from a research domain working on the same model class.

our approach differs from *likelihood-tempering* methods, which require an explicit evaluation of a *tilted* likelihood (raised to a power $0 < t < 1$) in order to prevent overconfident Bayesian updating [63].

In terms of the theory of subjective logic (SL, [80]), we can model second-order probabilities by placing a Dirichlet distribution over the estimated posterior model probabilities [145]. These second-order probabilities represent an uncertainty measure over quantities which are themselves probabilities. We use the second-order probabilities to capture epistemic uncertainty about whether the observed data has been generated by one of the candidate models considered during training.

The probability density function (PDF) of a Dirichlet distribution is given by:

$$\text{Dir}(\boldsymbol{\pi} \mid \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \quad (6.3)$$

where $\boldsymbol{\pi}$ belongs to the unit $J - 1$ simplex (i.e., $\boldsymbol{\pi} \in \Delta^J := \{\boldsymbol{\pi} \mid \sum_{j=1}^J \pi_j = 1\}$) and $B(\boldsymbol{\alpha})$ is the multivariate beta function [131]. The Dirichlet density is parameterized by a vector of *concentration parameters* $\boldsymbol{\alpha} \in \mathbb{R}_+^J$ which can be interpreted as evidences in the ST framework [80]. The sum of the individual evidence components $\alpha_0 = \sum_{j=1}^J \alpha_j$ is referred to as the Dirichlet strength, and it affects the precision of the higher-order distribution in terms of its variance. Intuitively, the Dirichlet strength governs the *peakedness* of the distribution, with larger values leading to more peaked densities (i.e., most of the density being concentrated in a smaller region of the

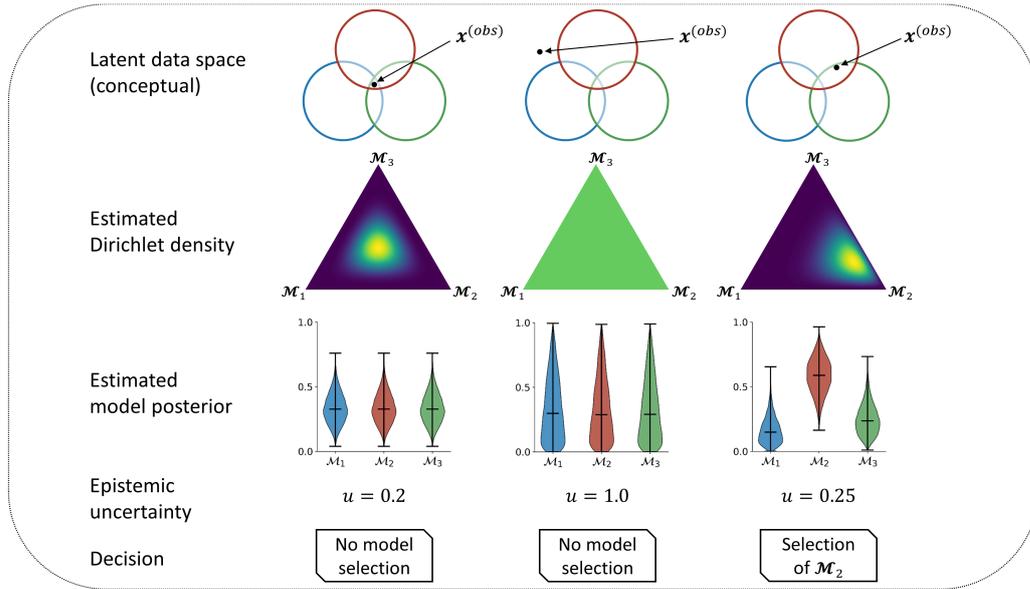


Figure 6.2: Three different hypothetical model comparison scenarios with different observations. The first column depicts observing a data set which is equally probable under all models. The second column depicts a data set which is beyond the generative scope of all models. The third column illustrates an observed data set which is most probable under model 2.

simplex). We can use the mean of the Dirichlet distribution, which is a vector of probabilities given by:

$$\mathbb{E}_{\boldsymbol{\pi} \sim \text{Dir}(\boldsymbol{\alpha})}[\boldsymbol{\pi}] = \frac{\boldsymbol{\alpha}}{\alpha_0} \quad (6.4)$$

to approximate the posterior model probabilities $p(\boldsymbol{m} \mid \boldsymbol{x}_{1:N})$, as will become clearer later in this section. A crucial advantage of such a Dirichlet representation is that it allows to look beyond model probabilities by inspecting the vector of computed evidences. For instance, imagine a scenario with three possible models. If $\boldsymbol{\alpha} = (5, 5, 5)$, the data provides equally strong evidence for all models (Figure 6.2, first column) – all models explain the data well. If, on the other hand, $\boldsymbol{\alpha} = (1, 1, 1)$, then the Dirichlet distribution reduces to a uniform on the simplex indicating no evidence for any of the models (Figure 6.2, second column) – no model explains the observations well. Note that in either case one cannot select a model on the basis of the data, because posterior model probabilities are equal, yet the interpretation of the two outcomes is very different: The second-order Dirichlet distribution allows one to distinguish between *equally likely* (first case) and *equally unlikely* (second case) models. The last column of Figure 6.2 illustrates a scenario with $\boldsymbol{\alpha} = (2, 7, 3)$ in which case one can distinguish between all models. Later, we will also demonstrate a scenario with data simulated from an actual model.

We can further quantify this distinction by computing an uncertainty score given by:

$$u = \frac{J}{\alpha_0} \quad (6.5)$$

where J is the number of candidate models. This uncertainty score ranges between 0 (total certainty) and 1 (total uncertainty) and has a straightforward interpretation. Accordingly, total uncertainty is given when $\alpha_0 = J$, which would mean that the data provide no evidence for any of the J candidate models. On the other hand, $u \ll 1$ implies a large Dirichlet strength $\alpha_0 \gg J$, which would read that the data provide plenty of evidence for one or more models in question. The uncertainty score corresponds to the concept of *vacuity* (i.e., epistemic uncertainty) in the terminology of SL [80]. We argue that epistemic uncertainty should be a crucial aspect in model selection, as it quantifies the strength of evidence, and, consequently, the strength of the theoretical conclusions we can draw given the observed data.

Consequently, model comparison in our framework consists in inferring the concentration parameters of a Dirichlet distribution given an observed or simulated data set. The problem of inferring posterior model probabilities can thus be reparameterized as:

$$p(\mathcal{M} | \mathbf{x}_{1:N}) \approx q_\eta(\mathbf{m} | \mathbf{x}_{1:N}) = \mathbb{E}_{\boldsymbol{\pi} \sim \text{Dir}(f_\eta(\mathbf{x}_{1:N}))}[\boldsymbol{\pi}] \quad (6.6)$$

where f_η is a neural network with positive outputs greater than one, that is, $f_\eta : \mathcal{X}^N \rightarrow [1, \infty]^J$. Additionally, we can also obtain a measure of absolute model evidence by considering the uncertainty encoded by the full Dirichlet distribution (Eq.6.5). Before elaborating on the latter point, we discuss the main concepts for learning relative evidence, since they form the backbone for further developments.

6.3.2 LEARNING EVIDENCE IN AN \mathcal{M} -CLOSED FRAMEWORK

How do we ensure that the outputs of the neural network match the true unknown model posterior probabilities? As per Algorithm 7, we have unlimited access to samples (simulations) from the joint model $p(\mathcal{M}, \mathbf{x}) = \int p(\mathcal{M}, \boldsymbol{\theta}, \mathbf{x}) d\boldsymbol{\theta}$. Consider, for ease of exposition, a data set with a single observation, that is $N = 1$ such that $\mathbf{x}_{1:N} = \mathbf{x}$. We use the mean of the Dirichlet distribution $q_\eta(\mathbf{m} | \mathbf{x})$ parameterized by an evidential neural network with parameters $\boldsymbol{\eta}$ to approximate $p(\mathcal{M} | \mathbf{x})$. To optimize the parameters of the neural network, we can minimize some loss \mathcal{L} in expectation over all possible data sets:

$$\boldsymbol{\eta}^* = \arg \min_{\boldsymbol{\eta}} \mathbb{E}_{p(\mathcal{M}, \mathbf{x})} [\mathcal{L}(q_\eta(\mathbf{m} | \mathbf{x}), \mathbf{m})] \quad (6.7)$$

$$= \arg \min_{\boldsymbol{\eta}} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{p(\mathcal{M} | \mathbf{x})} [\mathcal{L}(q_\eta(\mathbf{m} | \mathbf{x}), \mathbf{m})]] \quad (6.8)$$

where \mathbf{m} is a one-hot encoded vector of the true model index \mathcal{M}_j . We also require that \mathcal{L} be a *strictly proper loss* [59]. According to [59], a loss function in the context of simulation-based model comparison is strictly proper if and only if it attains its minimum when $q_\eta(\mathbf{m} | \mathbf{x}) = p(\mathcal{M} | \mathbf{x})$.

When we choose the Shannon entropy $\mathbb{H}(q_\eta(\mathbf{m} | \mathbf{x})) = -\sum_j q_\eta(\mathbf{m} | \mathbf{x})_j \log q_\eta(\mathbf{m} | \mathbf{x})_j$ for \mathcal{L} , we obtain the strictly proper logarithmic loss:

$$\mathcal{L}(q_\eta(\mathbf{m} | \mathbf{x}), \mathbf{m}) = -\sum_{j=1}^J m_j \log q_\eta(\mathbf{m} | \mathbf{x})_j \quad (6.9)$$

$$= -\sum_{j=1}^J m_j \log \left(\frac{f_\eta(\mathbf{x})_j}{\sum_{j'=1}^J f_\eta(\mathbf{x})_{j'}} \right) \quad (6.10)$$

where $m_j = 1$ when j is the true model index and 0 otherwise (i.e., standard one-hot encoding). Thus, in order to estimate ϕ , we can minimize the expected logarithmic loss over all simulated data sets where $f_\eta(\mathbf{x})_j$ denotes the j -th component of the Dirichlet density given by the evidential neural network. Since we use a strictly proper loss, the evidential network yields the true model posterior probabilities over all possible data sets when perfectly converged.

Intuitively, the logarithmic loss encourages high evidence for the true model and low evidences for the alternative models. Correspondingly, if a data set with certain characteristics can be generated by different models, evidence for these models will jointly increase. Additionally, the model which generates these characteristics most frequently will accumulate the most evidence and thus be preferred. However, we also require low evidence, or, equivalently, high epistemic uncertainty, for data sets which are implausible under all models. We address this problem in the next section.

6.3.3 LEARNING ABSOLUTE EVIDENCE THROUGH REGULARIZATION

We now propose a way to address the scenario in which no model explains the observed data well. In this case, we want the evidential network to estimate low evidence for all models in the candidate set. In order to attenuate evidence for data sets which are implausible under all models considered, we incorporate a Kullback-Leibler (KL) divergence into the criterion in Eq.6.9. We compute the KL divergence between the Dirichlet density generated by the neural network and a uniform Dirichlet density implying total uncertainty. Thus, the KL shrinks evidences which do not contribute to correct model assignments during training, so an implausible data set encountered in the inference phase will lead to low evidence under all models. This type of regularization has been used for capturing out-of-distribution (OOD) uncertainty in image classification tasks [145]. Curiously, the task of OOD detection closely resembles that of diagnosing model misspecification, so future developments in one of the areas would most likely benefit the other and *vice versa*.

Adding the KL regularization penalty, our modified optimization criterion becomes:

$$\boldsymbol{\eta}^* = \arg \min_{\boldsymbol{\eta}} \mathbb{E}_{p(\mathcal{M}, \mathbf{x})} [\mathcal{L}(q_\eta(\mathbf{m} | \mathbf{x}), \mathbf{m}) + \lambda \Omega(\tilde{\boldsymbol{\alpha}})] \quad (6.11)$$

with $\Omega(\tilde{\boldsymbol{\alpha}}) = \mathbb{KL}[\text{Dir}(\tilde{\boldsymbol{\alpha}}) || \text{Dir}(\mathbf{1})]$. The term $\tilde{\boldsymbol{\alpha}} = \mathbf{m} + (1 - \mathbf{m}) \odot \boldsymbol{\alpha}$ represents the estimated evidence vector after removing the evidence for the true model. This is possible, because we know the true model index sampled from the model prior $p(\mathcal{M})$ during the simulation-based training phase. During the inference phase, knowing the ground truth is not required anymore, since $\hat{\boldsymbol{\eta}}$ has already been obtained at this point.

The KL regularizer penalizes evidences for the false models and drives these evidences towards unity. Equivalently, it acts as a ground-truth preserving prior on the higher-order Dirichlet distribution which preserves evidence for the true model and attenuates misleading evidences for the false models. The hyperparameter λ controls the regularization weight and encodes the tolerance of the algorithm to accept implausible (out-of-distribution) data sets during inference. With large values of λ , it becomes possible to detect cases where all models are deficient (i.e., misspecified); with $\lambda = 0$, only relative evidence can be generated. Note, that in the latter case, we recover our original proper criterion without penalization. The KL weight λ should be selected through prior empirical considerations on how well the simulations cover the plausible set of real-world data sets.

Importantly, the introduction of the KL regularizer renders the loss no longer *strictly proper*. Therefore, a large regularization weight λ would lead to poorer calibration of the approximate model posteriors, as the regularized loss is no longer minimized by the true model posterior. However, since the KL prior is ground-truth preserving, the accuracy of recovering the true model should not be affected. Indeed, we observe this behavior across a number of simulated experiments. More analytical research is needed on the rate of miscalibration induced by a particular choice of λ .

To make optimization of Equation 6.11 tractable in practice, we utilize the fact that we can easily simulate batches of the form $\mathcal{D}_N^{(B)} = \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ via Algorithm 7 and approximate Eq.6.11 via standard backpropagation by minimizing the following loss:

$$\mathcal{L}(\boldsymbol{\eta}) = \frac{1}{B} \sum_{b=1}^B \left[- \sum_{j=1}^J m_j^{(b)} \log \left(\frac{f_{\boldsymbol{\eta}}(\mathbf{x}_{1:N}^{(b)})_j}{\sum_{j'=1}^J f_{\boldsymbol{\eta}}(\mathbf{x}_{1:N}^{(b)})_{j'}} \right) + \lambda \Omega(\tilde{\boldsymbol{\alpha}}^{(b)}) \right] \quad (6.12)$$

over multiple batches to converge at a Monte Carlo estimator $\hat{\boldsymbol{\eta}}$ of the optimal neural network parameters $\boldsymbol{\eta}^*$. In practice, convergence can be determined as the point at which the loss stops decreasing, a criterion similar to *early stopping*. Alternatively, the network can be trained for a pre-defined number of epochs. Note, that, at least in principle, we can train the network arbitrarily long, since we assume that we can access the full joint Bayesian distribution $p(\mathcal{M}, \mathbf{x}, N)$ through simulation (cf. Figure 6.1, left panel). In practice, early stopping seems to work reasonably well, since it requires no prior considerations on the (most likely unknown) optimal number of simulations or interventions during training.

6.3.4 IMPLICIT PREFERENCE FOR SIMPLER MODELS

Perfect convergence of the evidential network for a given model comparison problem implies $q_{\boldsymbol{\eta}}(\mathbf{m} | \mathbf{x}_{1:N}) \propto p(\mathbf{x}_{1:N} | \mathcal{M})p(\mathcal{M})$. Thus, a perfectly converged evidential network automatically encodes a preference for simpler models (Bayesian Occam’s razor). This is due to the fact that we are approximating an expectation over all possible data sets, parameters, and models (i.e., the full Bayesian distribution). Accordingly, a simple model has a narrow generative scope, so data sets generated by a simpler model will tend to be more similar compared to those from a more complex competitor. Therefore, during training, certain data sets which are plausible under multiple models will be generated most often by the simplest model. Thus, a perfectly converged evidential

network will capture this behavior by assigning higher posterior probability to the simplest model (assuming equal prior probabilities). Therefore, at least in theory, our method captures complexity differences arising purely from the generative behavior of the models and does not presuppose an *ad hoc* measure of complexity (e.g., number of parameters).

6.3.5 TRAINING AND INFERENCE

Algorithm 8 (Online) Training phase and inference phase for amortized Bayesian model comparison with regularization-based uncertainty estimation.

Require: f_η - evidential neural network, $\{\mathbf{x}_{1:N_i}^{(obs)}\}_{i=1}^I$ - list of observed data sets for inference, λ - regularization weight, B - number of simulations at each iteration (batch size).

- 1: *Simulation-based training phase:*
 - 2: **repeat**
 - 3: Generate a training batch $\mathcal{D}_N^{(B)} = \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ via Algorithm 7.
 - 4: Compute evidences for each simulated data set in $\mathcal{D}_N^{(B)}: \alpha^{(b)} = f_\eta(\mathbf{x}_{1:N}^{(b)})$.
 - 5: Compute loss according to Equation 6.12.
 - 6: Update neural network parameters η via backpropagation.
 - 7: **until** convergence to $\hat{\eta}$
 - 8: *Amortized inference phase:*
 - 9: **for** $i = 1, \dots, I$ **do**
 - 10: Compute model evidences $\alpha_i^{(obs)} = f_{\hat{\eta}}(\mathbf{x}_{1:N_i}^{(obs)})$.
 - 11: Compute uncertainty $u_i = J / \sum_{j=1}^J \alpha_{i,j}^{(obs)}$.
 - 12: Approximate true model posterior probabilities $p(\mathcal{M} | \mathbf{x}_{1:N_i}^{(obs)})$ via $q_\eta(\mathbf{m} | \mathbf{x}_{1:N_i}) = \alpha_i^{(obs)} / \sum_{j=1}^J \alpha_{i,j}^{(obs)}$.
 - 13: **end for**
 - 14: Choose further actions.
-

The training phase in our evidential framework can be carried out using the same ideas and considerations explored in the context of BayesFlow. Accordingly, one can choose between online, offline, or a hybrid learning regime, depending on the computational resources available, the complexity of the candidate models and the simulation budget allocated for performing model comparison. Thus, in order to avoid repetition, Algorithm 8 summarizes both the training and inference phase with our evidential method using online learning during training. Note, that steps 2-7 and 9-13 can be executed in parallel and with GPU support in order to dramatically accelerate convergence and inference. Importantly, if the priors over model parameters change or additional models need to be considered, the parameters η of a pre-trained network can be augmented to η' by adding additional output nodes for the new models. Training can then be resumed from where it had previously stopped without optimizing η' from scratch.

6.3.6 SOURCES OF ERROR

Since our evidential framework embodies some of the principles implemented in the BayesFlow framework (simulation-based training, amortized inference), it also inherits some of the error sources described in Section 5.6, namely, simulation gap, amortization gap, and approximation error. Most notably, our regularization approach to model misspecification is precisely designed to detect the presence of simulation gaps. Since our evidential networks distill global knowledge about the models' generative scopes, it is supposed to assign low evidences to models which encounter a simulation gap during inference. However, whenever an evidential network is trained to minimize Equation 6.12 with $\lambda = 0$ (i.e., no regularization is applied), simulation gaps remain an undetectable issue, at least until posterior predictive checks are performed. In any case, errors due to an amortization gap (i.e., learning a global neural estimator) and Monte Carlo approximation (i.e., estimating an expectation) remain something to be aware of when performing amortized neural model comparison.

Another source of error is an underexpressive evidential network which is unable to properly encode the probabilistic relationship between data and models. In this case, the evidential network will be poorly calibrated, that is, its outputs would not represent the true posterior distribution $p(\mathcal{M} | \mathbf{x})$. Fortunately, due to amortized inference, we can easily estimate and visualize the expected calibration error (ECE, [64]) of an evidential network over multiple simulations from $p(\mathcal{M}, \boldsymbol{\theta}, \mathbf{x})$. Accordingly, ECE values close to 0 indicate proper calibration of the network. Some model comparison scenarios may prioritize different metrics, such as accuracy or precision/recall ratios, common to classification tasks in machine learning applications.

6.4 A SIMULATED EXPERIMENT

As a brief illustrative example (described in detail in [134]), we applied our evidential method to distinguish between complex nested spiking neuron models describing the properties of biological cells in the nervous system. The purpose of this experiment was twofold. On the one hand, we wanted to assess the ability of our method to classify models deploying a variety of neural patterns accounting for different cortical and sub-cortical neuronal activity. On the other hand, we wanted to investigate the network's ability to detect biologically implausible data patterns, as indexed by our measure of epistemic uncertainty. To this aim, we rely on a renowned computational model of biological neural dynamics.

6.4.1 MODEL COMPARISON SETTING

In computational neuroscience, mathematical models of neuronal electrical dynamics serve as a basis to explain the functional organization of the brain from both single neuron and large-scale neuronal networks processing perspectives [1, 16, 70, 75]. A multitude of different neuron models have been proposed during the last decades, ranging from completely abstract to biologically plausible models. The former offer a simplified mathematical representation which takes the main functional properties of spiking neurons into account. The latter provide a detailed analogy between models' state variables and ion channels in biological neurons [129]. Importantly, these

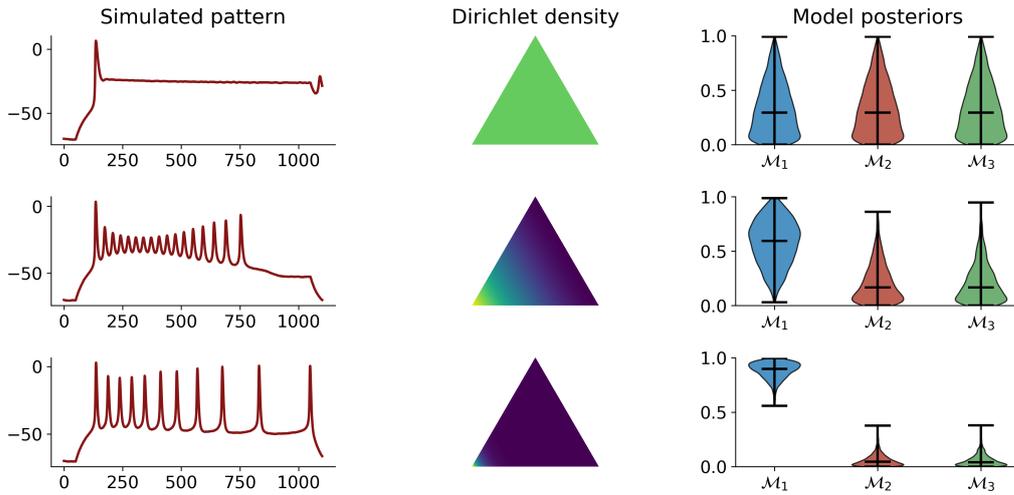


Figure 6.3: Three simulated firing patterns, corresponding estimated Dirichlet densities and model posteriors [134]. Each row represents a different value of the parameter \bar{g}_K , $\bar{g}_K = 0.1$, $\bar{g}_K = 0.5$, and $\bar{g}_K = 0.75$, respectively. An increase in the parameter \bar{g}_K is accompanied by a decrease in epistemic uncertainty (as measured via Eq.6.5). An implausible value of \bar{g}_k (first row) results in a flat Dirichlet density as an index of total epistemic uncertainty (uniform green areas). As the parameter value surpasses the plausibility boundary (second and third rows), the Dirichlet simplex becomes peaked towards the lower left edge encoding \mathcal{M}_1 .

computational models differ in their capability to reproduce firing patterns observed in real cortical neurons [76].

For this simulated experiment, we consider a Hodgkin-Huxley stateful model of cortical and thalamic neurons [70, 130]. The forward model is formulated as a set of five ordinary differential equations (ODEs) describing how the neuron membrane potential $V(t)$ changes over time as a function of the injected current $I_{inj}(t)$ and of various ion channels properties (see [134] for more details regarding the forward model).

To set up the model comparison problem, we treat different types of conductance, g_L , \bar{g}_{Na} , \bar{g}_K and \bar{g}_M , as free parameters, and define different neural models based on different parameter specifications. In particular, we formulate three models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3\}$ defined by the parameter sets $\theta_1 = (\bar{g}_{Na}, \bar{g}_K)$, $\theta_2 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M)$, and $\theta_3 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M, g_L)$, respectively.

In order to evaluate performance, we train an unregularized recurrent evidential network for 60 epochs resulting in 60000 backpropagation updates. At each iteration, we draw a random input current duration $T \sim \mathcal{U}_D(100, 400)$ (in units of milliseconds), keeping a constant input current, I_{inj} . T reflects the physical time window in which biological spiking patterns can occur. Since the sampling rate of membrane potential is fixed ($dt = 0.2$), T affects both the span of observable spiking behavior and the number of simulated data points.

6.4.2 VALIDATION RESULTS

The entire training phase using online learning (Algorithm 8) took approximately 2.5 hours of wall-clock time. On the other hand, model comparison on 5000 neural time-series simulated for validation took approximately 0.7 seconds, which is a remarkable efficiency gain.

Regarding model selection performance, the network exhibited accuracies above 0.92 across all T s, with no gains in accuracy for increasing T . This result highlights the fact that even short input currents are sufficient for reliably distinguishing between these complex models. Further, we observed good calibration for all three models, with all ECEs less than 0.1. Notably, we observed no overconfidence for all three models.

In order to assess how well we can detect biologically implausible patterns, we train an identical recurrent evidential network with a gradually increasing regularization weight up to $\lambda = 1.0$. We then fix the parameter $\bar{g}_{Na} = 4.0$ of model \mathcal{M}_1 and gradually increase its second parameter \bar{g}_K from 0.1 to 2.0. Since spiking patterns observed with low values of \bar{g}_K are quite implausible and have not been encountered during training, we expect uncertainty to gradually decrease. Indeed, Figure 6.3 shows this pattern. On the other hand, changing the sign of the output membrane potential, which also results in biologically implausible patterns, leads to a trivial selection of \mathcal{M}_3 . This is contrary to expectations, and shows that absolute evidence is also relative to the model knowledge the evidential network has learned during training. Future research should therefore focus on making the latent space of the evidential network interpretable, in order to make the conceptual visualization from Figure 6.2 tractable.

6.5 CONCLUDING REMARKS

This chapter introduced the building blocks of our evidential framework for Bayesian model comparison and discussed its mathematical and algorithmic formulation. In contrast to BayesFlow, applications of our amortization approach to real-world model comparison/selection problems are still underway. One reason for this is that recent developments in the field of simulation-based statistical inference have focused predominantly on parameter estimation and model comparison has often played a secondary role (or has been too costly to perform). Thus, we hope that our framework (or underlying ideas) can enhance and enrich model-based analysis and inference in many fields dealing with competing computational models of complex natural processes. We leave it to future research to investigate whether there are more elegant ways to quantify absolute evidence or detect model misspecification from a simulation-based perspective. Details regarding training, hyperparameter choice, and validation metrics can be found in our methodological paper [134]. Further details regarding implementation as well as templates for model comparison are also available at the code repository (<https://github.com/stefanradev93/BayesFlow>).

7 META-AMORTIZED INFERENCE

What is now proved was once only imagined.

— William Blake

In this chapter, we explore the idea of a universal neural Bayesian inference architecture for performing simultaneous parameter estimation and model comparison in a purely simulation-based way. We build on our previous methods by using ideas from *multi-task* learning [20] and meta-amortized variational inference [25]. In this way, we propose to enable and amortize all steps of a Bayesian workflow within a unified framework involving a single training/optimization phase. Instead of presenting a ready-made solution, this chapter merely intends to point out towards a speculative future development aimed at scaling up an entire Bayesian workflow to complex models.

7.1 THE BAYESIAN HARDSHIPS

A Bayesian analysis consists of more than just parameter estimation and model comparison. The big picture of Bayesian inference involves a rather significant allocation of creative, computational, financial, and decision making resources (cf. Figure 7.1 for an illustrative overview). Most recently, attempts have been made to systematize Bayesian analysis into a principled, step-by-step workflow reminiscent of a cooking recipe [49, 56, 144]. Naturally, it is beyond the scope of this chapter to review these comprehensive works. Thus, we will attempt to extract the most basic elements of a Bayesian workflow which can directly benefit from the notion of amortized inference.

The starting point of our Bayesian workflow of interest is a collection of observed data sets $\mathcal{D}^{(obs)} = \{\mathbf{x}_{1:N_i}^{(obs)}\}_{i=1}^I$, with $I = 1$ in the case of a single data set and $N = 1$ in the case of a single observation, and a collection of J competing models $\mathcal{M} = \{\mathcal{M}_j\}_{j=1}^J$. Most Bayesian pipelines would go through multiple steps involving, among other things, a considerable computational burden. These steps are represented by the red-shaded boxes in Figure 7.1 and summarized as follows:

1. Parameter estimation
2. Evaluation of computational faithfulness
3. Evaluation of model adequacy/sensitivity
4. Posterior predictive checks
5. Model comparison/model aggregation

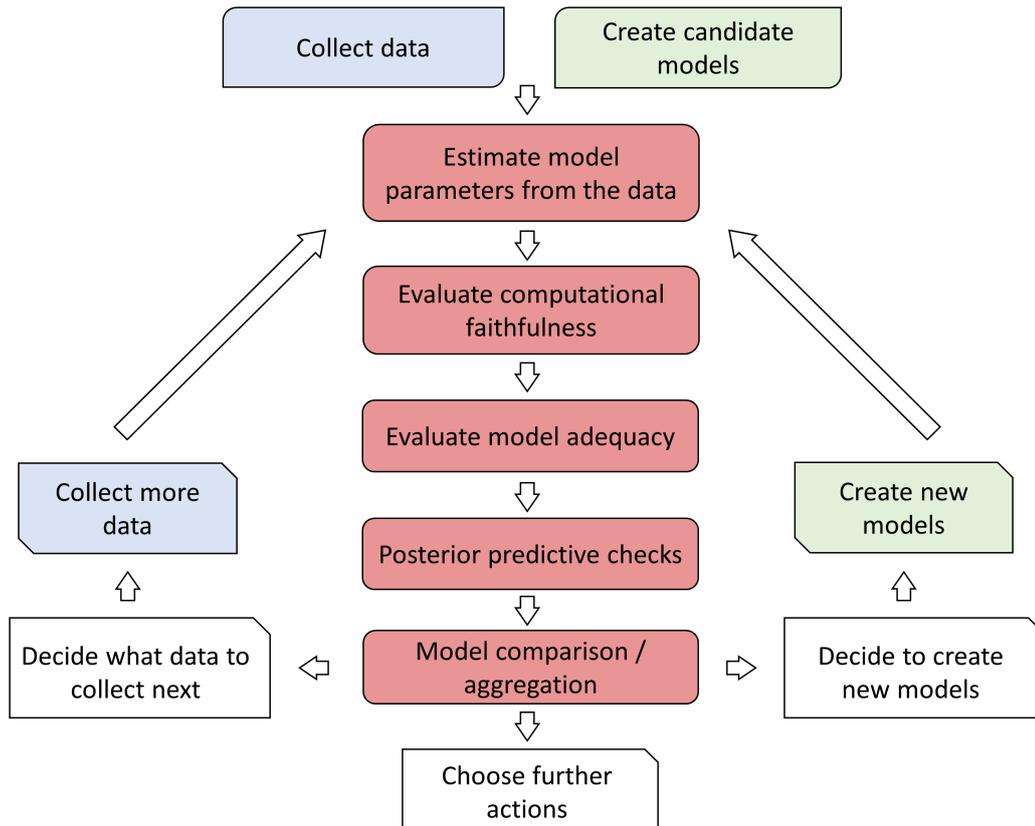


Figure 7.1: The basic conceptual steps of a Bayesian analysis pipeline (workflow) associated with allocation of different types of resources. White indicates allocation of decisional resources. Blue indicates allocation of financial resources. Green indicates allocation of creative resources. Red indicates allocation of computational resources. Especially the latter can profit to a great extent from amortization.

As we saw in the previous two chapters, we can tackle parameter estimation as well as checks of computational faithfulness, model sensitivity, and posterior predictions with our BayesFlow framework. Furthermore, we can circumvent fitting all models explicitly to each data set in $\mathcal{D}^{(obs)}$ and perform efficient model comparison with our evidential framework. However, as it currently stands, these frameworks seem rather disconnected and conceptualized to work on their own.

For instance, in order to perform parameter estimation for all models in \mathcal{M} , one has to train and store J neural density estimators. When using BayesFlow for parameter estimation, one might reuse the same (pre-trained) summary network over all models, thus effectively pooling some of the resources. However, separate inference networks would still be needed for each of the models. Needless to say, such an approach does not scale well when J is large and needs the resources of a computational cluster to be feasible in practice.

In addition, in order to perform (prior predictive) model comparison, a disjoint training phase for an evidential network needs to be introduced. Researchers then need to ensure that simulations are shared between the parameter estimation steps and the model comparison step, other-

wise a considerable portion of the simulation budget would be wasted by discarding simulations. Thus, the need for a framework which amortizes all of the above steps in a Bayesian workflow becomes immediately obvious. For such a framework to be useful in practice, it needs to enable, at least in theory, amortization over an arbitrary number of models, data sets, and observations (data set sizes).

7.2 AMORTIZED INFERENCE REVISITED

In the following, we continue our discussion on different *levels* of amortization. In Chapter 5, we introduced three different types of amortized Bayesian inference bootstrapped by neural density estimation: case-wise, model-wise, and meta-amortized. We further assume, for the sake of our discussion, that the neural networks employed in these approaches are all capable of fully Bayesian inference (i.e., return a full posterior distribution) and use the raw simulated or observed data directly (i.e., do not rely on manual selection of summary statistics). Note, that the different types of amortization have not been explicitly distinguished in the literature on simulation-based inference, so our nomenclature is rather non-standard.

Case-wise amortized methods require a separate optimization loop for each observed data set and model. When case-wise methods incorporate a training phase (e.g., APT in a sequential regime [60]), it must be repeated for each new data set and model, since the observed data is part of the optimization criterion. The general form of the case-wise optimization criterion for obtaining optimal neural network parameters is given by:

$$\varphi_{i,j}^* = \arg \min_{\varphi} \mathbb{E}_{p(\boldsymbol{\theta}_j | \mathbf{x}_{1:N_i}^{(obs)}, \mathcal{M}_j)} [-\log q_{\varphi}(\boldsymbol{\theta}_j | \mathbf{x}_{1:N}, \mathcal{M}_j)] \quad (7.1)$$

where we have a separate set of neural network parameters $\varphi_{i,j}$ for each data set i and model j . The case-wise approach to amortizing Bayesian inference is illustrated in Figure 7.2.

Model-wise amortized methods require a global upfront training phase before any real data are collected via simulations from each joint Bayesian model $p(\boldsymbol{\theta}_j, \mathbf{x}, N | \mathcal{M}_j)$. During inference, model-wise methods operate entirely in a feed-forward manner, that is, they involve no training or optimization in this phase. Thus, the upfront training effort amortizes over all observed data sets from the generative scope of model \mathcal{M}_j defined by its corresponding prior $p(\boldsymbol{\theta}_j)$ and simulator $g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi})$. However, in the frequent case of multiple candidate models, one still needs to perform separate optimization loops for each of the models in \mathcal{M} . The general form of the model-wise optimization criterion for obtaining optimal neural network parameters is given by:

$$\varphi_j^* = \arg \min_{\varphi} \mathbb{E}_{p(\boldsymbol{\theta}_j, \mathbf{x}, N | \mathcal{M}_j)} [-\log q_{\varphi}(\boldsymbol{\theta}_j | \mathbf{x}_{1:N}, \mathcal{M}_j)] \quad (7.2)$$

where we only have a separate set of neural network parameters for each model j . This is due to the fact that the expectation runs over the model-implied joint distribution $p(\boldsymbol{\theta}_j, \mathbf{x}, N | \mathcal{M}_j)$ and the observed data does not enter the optimization phase. The model-wise approach to amortizing Bayesian inference is illustrated in Figure 7.3.

Finally, a meta-amortized approach requires an even costlier upfront training phase involving simulations from multiple models (i.e., multi-model forward inference). The resulting benefit

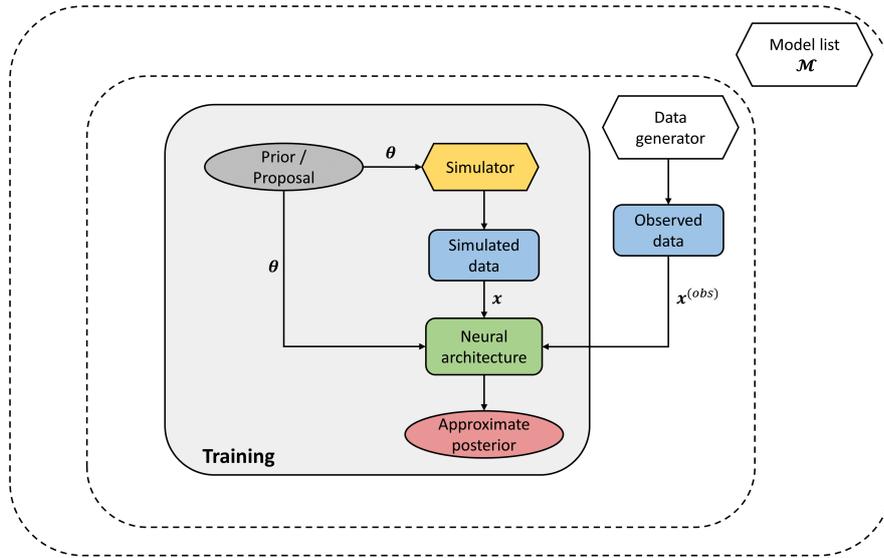


Figure 7.2: Case-wise amortized Bayesian inference with a generative neural architecture capable of fully Bayesian inference. The gray-shaded plane indicates the scope of amortization.

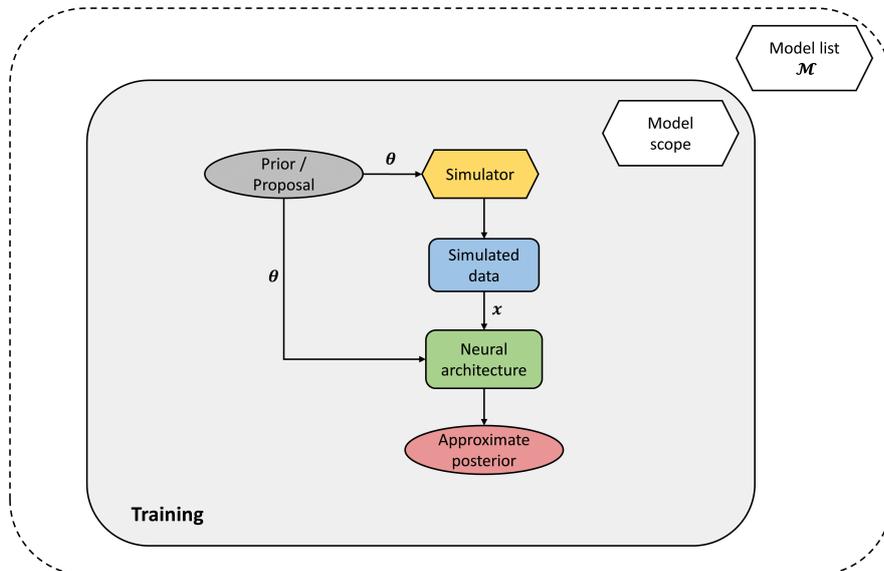


Figure 7.3: Model-wise amortized Bayesian inference with a generative neural architecture capable of fully Bayesian inference. The gray-shaded plane indicates the scope of amortization and, in contrast to case-wise approaches, includes the entire generative scope of the model.

manifests itself in the outcome that a single (composite) network is able to account for all models and all possible data sets arising from the models. Accordingly, the general form of the meta-amortized optimization criterion is given by:

$$\varphi^* = \arg \min_{\varphi} \mathbb{E}_{p(\mathcal{M}, \boldsymbol{\theta}, \mathbf{x}, N)} [-\log q_{\varphi}(\boldsymbol{\theta}, \mathcal{M} | \mathbf{x}_{1:N})] \quad (7.3)$$

where the expectation runs over the full joint model $p(\boldsymbol{\theta}, \mathbf{x}, N, \mathcal{M})$. The meta-amortized criterion not only opens new possibilities but also brings new challenges to which we turn next.

7.3 LEARNING A MULTI-MODEL POSTERIOR

In order to approximate the multi-model posterior $p(\boldsymbol{\theta}, \mathcal{M} | \mathbf{x}_{1:N})$, we seek neural network parameters φ which minimize our meta-amortized criterion from Equation 7.3. We can expand the latter as follows:

$$\varphi^* = \arg \min_{\varphi} \mathbb{E}_{p(\mathcal{M})} [\mathbb{E}_{p(\mathbf{x}, N | \mathcal{M})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x}_{1:N}, \mathcal{M})} [-\log q_{\varphi}(\boldsymbol{\theta}, \mathcal{M} | \mathbf{x}_{1:N})]]] \quad (7.4)$$

$$= \arg \min_{\varphi} - \sum_{j=1}^J \int_{\mathcal{X}} \int_{\Theta_j} p(\mathcal{M}_j, \boldsymbol{\theta}_j, \mathbf{x}_{1:N}) \log q_{\varphi}(\boldsymbol{\theta}_j, \mathcal{M}_j | \mathbf{x}_{1:N}) d\boldsymbol{\theta}_j d\mathbf{x} \quad (7.5)$$

which we can approximate via Monte Carlo simulations from $p(\mathcal{M}, \boldsymbol{\theta}, \mathbf{x}_{1:N})$ (i.e., multi-model forward inference):

$$\hat{\varphi} = \arg \min_{\varphi} - \frac{1}{B} \sum_{b=1}^B \log q_{\varphi}(\boldsymbol{\theta}_j^{(b)}, \mathcal{M}_j^{(b)} | \mathbf{x}_{1:N}^{(b)}) \quad (7.6)$$

Correspondingly, we can treat Equation 7.6 as a loss function and minimize it with any stochastic gradient descent method. To derive a tractable criterion, we can further expand Equation 7.6 into:

$$\hat{\varphi} = \arg \min_{\varphi} \frac{1}{B} \sum_{b=1}^B -\log q_{\varphi}(\boldsymbol{\theta}_j^{(b)} | \mathbf{x}_{1:N}^{(b)}, \mathcal{M}_j^{(b)}) - \log q_{\varphi}(\mathcal{M}_j^{(b)} | \mathbf{x}_{1:N}^{(b)}) \quad (7.7)$$

The above formulation has two important components: the (amortized) approximate parameter posterior $q_{\varphi}(\boldsymbol{\theta} | \mathbf{x}_{1:N}, \mathcal{M})$ and the (amortized) approximate model posterior $q_{\varphi}(\mathcal{M} | \mathbf{x}_{1:N})$. We will explore an architecture for meta-amortized inference consisting of three neural network components: an inference network (parameterized by ϕ), a summary network (parameterized by ψ), and an evidence network (parameterized by η). Thus, φ represents the collection of all neural network parameters, $\varphi = (\phi, \psi, \eta)$. The inference network is responsible for approximating each parameter posterior $p(\boldsymbol{\theta}_j | \mathbf{x}_{1:N}, \mathcal{M}_j)$. The summary network is responsible for extracting maximally informative summary vectors from raw data. Last, the evidence network is responsible for approximating the model posterior $p(\mathcal{M}, \mathbf{x}_{1:N})$. We now describe how to render optimization of Equation 7.7 tractable.

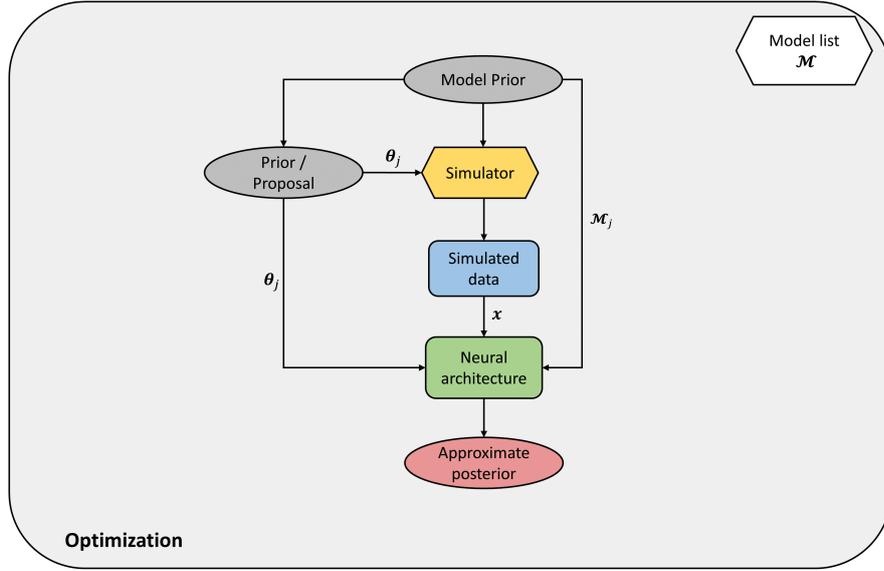


Figure 7.4: Meta-amortized Bayesian inference with a generative neural architecture capable of fully Bayesian inference. The gray-shaded plane indicates the scope of amortization. In contrast to case-wise and model-wise approaches, the latter includes the generative scopes of all models in the model list \mathcal{M} .

First, we can represent the (multi-model) parameter posterior via a doubly conditional INN implementing a normalizing flow between θ_j and z_j for all j

$$q_\phi(\theta_j | \mathbf{x}_{1:N}, \mathcal{M}_j) = p(z_j = f_\phi(\theta_j; h_\psi(\mathbf{x}_{1:N}), \mathbf{m})) \left| \det \left(\frac{\partial f_\phi(\theta_j; h_\psi(\mathbf{x}_{1:N}), \mathbf{m})}{\partial \theta_j} \right) \right| \quad (7.8)$$

where h_ψ is any (algebraically aligned) summary network with trainable parameters ψ and \mathbf{m} is a one-hot encoded vector representation of the abstract model index \mathcal{M}_j . Writing \mathbf{J}_{f_ϕ} as a shorthand for the Jacobian of the learnable transformation, we can derive the following loss function for a batch of B simulated model indices, parameters, and data sets:

$$\mathcal{L}_{KL}(\phi, \psi) = \frac{1}{B} \sum_{b=1}^B \left(\frac{\|f_\phi(\theta_j^{(b)}; h_\psi(\mathbf{x}_{1:N}^{(b)}), \mathbf{m}^{(b)})\|_2^2}{2} - \log |\det \mathbf{J}_{f_\phi}^{(b)}| \right), \quad (7.9)$$

which is a modified version of the BayesFlow criterion (Equation 5.40) with a model index included as a further conditioning input for the cINN.

Second, we can represent the model posterior using our evidential formulation

$$q_\eta(\mathcal{M} | \mathbf{x}_{1:N}) = \mathbb{E}_{\text{Dir}(f_\eta(\mathbf{x}_{1:N}))}[\boldsymbol{\pi}], \quad (7.10)$$

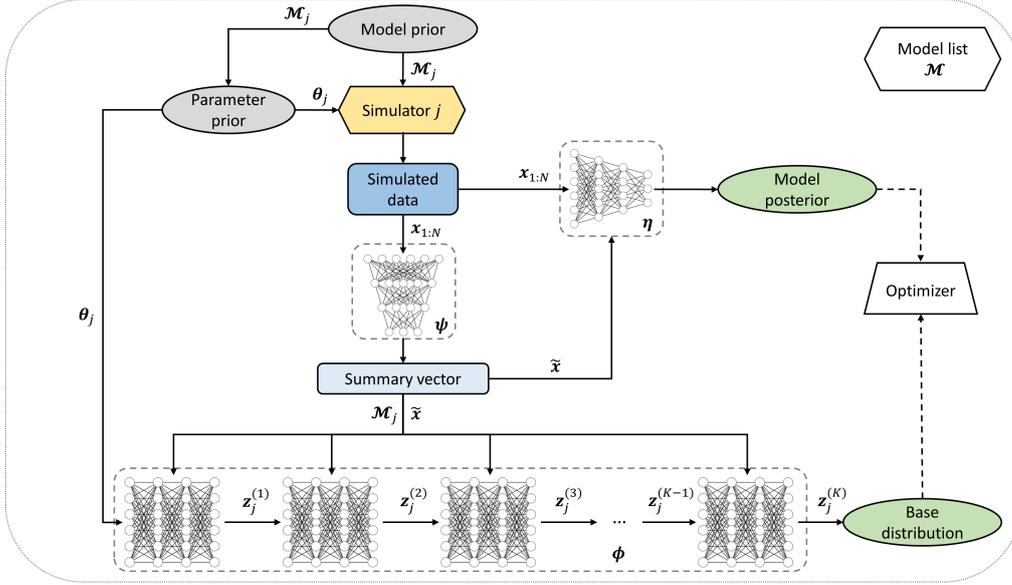


Figure 7.5: A possible framework for meta-amortized Bayesian inference connecting together an inference network (ϕ), a summary network (ψ), and an evidence network (η). All three networks are optimized together and trained with an arbitrary number of simulations from the joint Bayesian model $p(\mathcal{M})p(\boldsymbol{\theta} | \mathcal{M})p(\mathbf{x} | \boldsymbol{\theta}, \mathcal{M})$.

where $\text{Dir}(f_{\eta}(\mathbf{x}_{1:N}))$ denotes a Dirichlet density with concentration parameters provided by an evidential network $f_{\eta}(\mathbf{x}_{1:N})$ with parameters $\boldsymbol{\eta}$. Note also, that the evidential network might directly use the representation provided by the summary network h as a single input, $f_{\eta}(h_{\psi}(\mathbf{x}_{1:N}))$, or as an additional input concatenated with the raw data, $f_{\eta}(\mathbf{x}_{1:N}, h_{\psi}(\mathbf{x}_{1:N}))$. Moreover, our evidential formulation allows us to re-use the concepts for learning absolute evidence introduced in Chapter 6. However, if model misspecification is not considered an issue, one could use any probabilistically calibrated classifier for $q_{\eta}(\mathcal{M} | \mathbf{x}_{1:N})$. Accordingly, we can minimize the following (unregularized cross-entropy) loss function:

$$\mathcal{L}_{CE}(\boldsymbol{\eta}, \boldsymbol{\psi}) = \frac{1}{B} \sum_{b=1}^B \left[- \sum_{j=1}^J m_j^{(b)} \log \left(\frac{f_{\eta}(\mathbf{x}_{1:N}^{(b)}, h_{\psi}(\mathbf{x}_{1:N}^{(b)}))_j}{\sum_{j'=1}^J f_{\eta}(\mathbf{x}_{1:N}^{(b)}, h_{\psi}(\mathbf{x}_{1:N}^{(b)}))_{j'}} \right) \right], \quad (7.11)$$

which assumes that the evidential network processes the output of the summary network, in addition to the raw simulated data.

Finally, putting the two together, our composite loss for meta-amortized inference becomes

$$\mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\psi}, \boldsymbol{\eta}) = \mathcal{L}_{KL}(\boldsymbol{\phi}, \boldsymbol{\psi}) + \mathcal{L}_{CE}(\boldsymbol{\eta}, \boldsymbol{\psi}) \quad (7.12)$$

In this way, multi-model Bayesian inference is amortized through a single set of network parameters $\hat{\boldsymbol{\varphi}} = (\hat{\boldsymbol{\phi}}, \hat{\boldsymbol{\psi}}, \hat{\boldsymbol{\eta}})$ obtained via backpropagation through the entire composite architecture. An example implementation of such an architecture is illustrated in Figure 7.5. The practical utility

and advantages of such a framework are yet to be demonstrated in simulated experiments and on observed data in a research context.

7.4 USE CASES AND CHALLENGES FOR META-AMORTIZED INFERENCE

A realization of the imaginable architecture depicted in Figure 7.5 would ensure that the computational burden of many Bayesian tasks (see red-shaded boxes in Figure 7.1) is attenuated via amortization. During inference, posterior draws from each parameter posterior $p(\boldsymbol{\theta}_j | \mathbf{x}_{1:N}^{(obs)}, \mathcal{M}_j)$ can be efficiently obtained by feeding the observed data and the desired model index through the summary and inference networks. Amortized evaluation of computational faithfulness and model adequacy as well as posterior predictive checks are all consequences of amortizing multi-model posterior inference. The model posterior $p(\mathcal{M} | \mathbf{x}_{1:N})$ can be estimated by feeding the output of the summary network from the previous step together with the observed data to the evidence network. Thus, model comparison via Bayes factors or model aggregation via Bayesian model averaging can also be performed efficiently upon convergence.

However, the actual implementation of the neural architecture depicted in Figure 7.5 is not as straightforward as its conceptualization. When it comes to the architecture of the network responsible for parameter inference, it appears necessary to distinguish between three main scenarios encountered in multi-model inference.

The first corresponds to *prior sensitivity analysis* in which the consequences of different prior configurations for subsequent Bayesian updating are systematically investigated. In this case, the simulator for each \mathcal{M}_j is the same and only the corresponding $p(\boldsymbol{\theta} | \mathcal{M}_j)$ differ in their distributional form. Prior sensitivity analysis is perhaps also the easiest case to tackle, since it requires no essential structural changes to the cINN, which is augmented to accept the one-hot encoded model index as an additional conditioning variable.

The second corresponds to a setting in which, once again, the simulator remains conceptually and functionally the same, but some components of $\boldsymbol{\theta}$ are treated as fixed and some as varying. In this case, some parts of the latent space \mathbf{z} are shared among all models, whereas others are missing whenever a subset of the model parameters is treated as fixed. In other words, a complete parameter vector is reduced from $\boldsymbol{\theta} \in \mathbb{R}^D$ to $\boldsymbol{\theta}_j \in \mathbb{R}^{D_j}$ with $D_j \leq D$ for each model j . Two potential approaches for performing meta-amortized Bayesian inference in such a setting appear viable. In the first, each reduced parameter space Θ_j is augmented to Θ'_j , with additional dummy parameters following a simple distribution (e.g., Gaussian), such that each augmented $\boldsymbol{\theta}'_j$ has the same dimensionality. A disadvantage of this approach is, that the inference network needs to learn an identity transformation between the dummy parameters and the corresponding latent variables, which seems inelegant. In the second approach, missing parameters are encoded with zeros and the loss function is masked, such that each \mathbf{z}_j is optimized only from the remaining parameters. A disadvantage of this approach is that it is less straightforward to implement and might lead to training instabilities.

The third scenario corresponds to the task of comparing essentially different generative mechanisms assumed to account for the same data. In this case, each \mathcal{M}_j is implemented as a different simulator g_j and the number of parameters might or might not differ between the simulators.

Such a scenario would require learning disjoint latent spaces z_j for each simulator. It might even necessitate separate inference networks (with a shared summary network) for each model in \mathcal{M} or a completely different neural architecture altogether.

As could be gathered from this short exposition, many problems lurking on the path towards a universal framework for meta-amortized simulation-based inference remain unsolved and many unsuspected challenges are yet to reveal themselves. However, since meta-amortized inference appears to be a desirable goal for many scientific domains, future research should explore various promising avenues for actually attaining it.

8 APPLICATIONS

This chapter briefly reviews four concrete applications of our ideas for amortized Bayesian inference to real-world modeling problems. The details of these modeling scenarios have already been described in the corresponding papers [31, 94, 135, 172], so the purpose of this chapter is merely to convey the gist of each particular application and describe how our developed frameworks contributed to solving the problem of inference.

8.1 A BAYESIAN BRAIN MODEL OF ADAPTIVE BEHAVIOR

In this work¹, we proposed and validated a new computational Bayesian model accounting for individual performance in the Wisconsin Card Sorting Test (WCST), an established clinical tool for measuring set-shifting and deficient inhibitory processes on the basis of environmental feedback [6, 67].

Performance in WCST is usually measured via a rough summary metric such as the number of correct/incorrect responses or pre-defined psychological scoring criteria (see for instance [67]). These metrics form the basis for inferring the underlying cognitive processes recruited by the task. However, a major shortcoming of this approach is that it merely assumes the cognitive processes to be inferred without specifying an explicit process model. Moreover, summary measures do not utilize the full information present in the data, such as trial-by-trial fluctuations or relevant agent-environment interactions. For this reason, crude scoring measures are often insufficient to disentangle the dynamics of the relevant cognitive components.

To address this shortcoming, we formalized the interaction between the task's structure, the received feedback, and the participant's behavior by building a model of the underlying information processing mechanisms used to infer the hidden rules of the task environment. Furthermore, we embedded the new model within the mathematical framework of the Bayesian Brain Theory (BBT), according to which beliefs about hidden environmental states are dynamically updated following the logic of Bayesian inference [48, 91].

The simple controlled setting (environment) realized by the WCST consists of a target and a set of stimulus cards with geometric figures which vary according to three perceptual features: color (red, green, blue, yellow), shape (triangle, star, cross, circle), and number of objects (1, 2, 3, 4). Participants in the test have to infer the correct classification principle by trial and error using the examiner's or computer's feedback. The feedback carries a positive or a negative signal informing the participant whether her choice of action was appropriate or not. Moreover, modeling adaptive behavior in the WCST from a Bayesian perspective is straightforward, since observable actions

¹M. D'Alessandro, S. T. Radev, A. Voss, and L. Lombardi. "A Bayesian brain model of adaptive behavior: an application to the Wisconsin Card Sorting Task". *PeerJ* 8, 2020, e10316

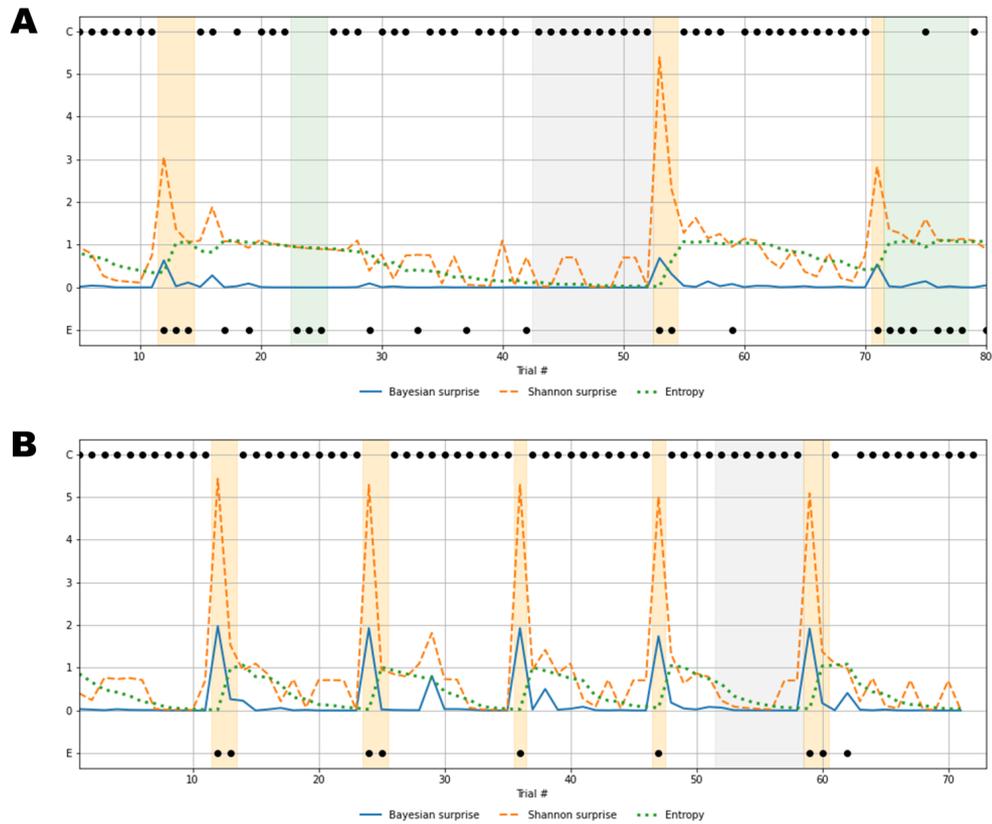


Figure 8.1: Estimated information processing dynamics of two exemplary individuals [31]. **(A)** Trial-by-trial information-theoretic measures of an individual with SD characterized by very low flexibility and very high information loss; **(B)** Trial-by-trial information-theoretic measures of a healthy control individual characterized by relatively high flexibility and low information loss. Labels C and E on the y -axis indicate correct and error responses.

emerge from the interaction between the internal probabilistic model of the agent and a set of discrete environmental states.

The main contributions of our modeling work were thus threefold. First, we developed a two-level model of adaptive agent-environment interaction, consisting of a cognitive and an information-theoretic component. The cognitive component decomposes performance in the WCST into two interpretable parameters: *flexibility* and *information loss*. The information-theoretic component transforms the parameters into dynamic measures of belief updating, surprise, and internal model uncertainty. Second, we performed extensive simulation studies for ensuring reasonable computational faithfulness and model sensitivity. Third, we applied the model to a sample of individuals with substance dependence (SD) and a sample of healthy controls to account for (mal)adaptive task performance in a principled way. For the latter two tasks, we had to resort to amortized inference with BayesFlow, since the likelihood function of our custom dynamic model is unknown. Overcoming this intractability with BayesFlow (using a recurrent summary network), we could perform both simulation-based calibration, parameter recovery, and inference on real data in a matter of seconds, once training had converged. The entire training phase took approximate 12 hours wall-clock time on a laptop with a graphics card.

Our initial application showed promising results in explaining adaptive behavior in the WCST. Figure 8.1 depicts the model-derived information processing dynamics of an individual with SD (upper panel) and a healthy control (lower panel). Indeed, patterns of belief updating (*Bayesian surprise*), surprisal (*Shannon surprise*), and model uncertainty (*Entropy*) are very different for the two individuals, highlighting the ability of the model to discriminate between sub-optimal and nearly-optimal performance via multiple sources of information (see [31] for a detailed interpretation).

8.2 JUMPING TO CONCLUSION? A LÉVY-FLIGHT MODEL OF DECISION MAKING

In this work², we formally tested whether a Lévy flight model, assuming an α -stable noise distribution with a free parameter α can provide a more accurate description of performance in simple binary decision making tasks than a classical diffusion model, assuming a Gaussian noise distribution.

Distributions with *fat tails*, such as the Cauchy distribution or the Lévy distribution, are characterized by an increased probability for extreme events, compared to a Gaussian distribution [95]. Moreover, fat-tailed models incorporating so-called Lévy flights have been applied in a variety of research contexts. For instance, such models have proven useful to account for animal foraging behavior. The Lévy flight foraging hypothesis states that in certain natural environments, (truncated) Lévy flights optimize random searches. Accordingly, the hypothesis implies that biological organisms have evolved to exploit occasional large divergences in their wandering movements during foraging, which are best accounted for by Lévy flights [168].

In our study on human decision making [172], we compared the relative fit of four evidence accumulation models applied to a color discrimination and a lexical decision task. In the color

²E. M. Wieschen, A. Voss, and S. Radev. “Jumping to conclusion? a lévy flight model of decision making”. *TQMP* 16:2, 2020, pp. 120–132

discrimination task, participants were tasked to indicate whether there were more orange or more blue pixels in a set of specifically designed stimuli. In the lexical decision task, participants were tasked to indicate whether a presented string of letters was an existing German word or a meaningless sequence. The candidate models included: \mathcal{M}_1 - a parsimonious version of the diffusion model with Gaussian noise; \mathcal{M}_2 - a model with an α -stable distribution for the noise of evidence accumulation; \mathcal{M}_3 a full version of the diffusion model with inter-trial variability for drift, starting point and non-decision time; and \mathcal{M}_4 a model with alpha as a free parameter and all previous inter-trial variability parameters. Note, that the all-time favorite Gaussian distribution is a special case of the stable family with a stability value of $\alpha = 2.0$. Lower values of α imply a higher probability of extreme events and thus occasional large *jumps* in the evidence accumulation process.

The inclusion of stable noise in the accumulation process renders numerical evaluation of the likelihood intractable. Thus, we train a separate BayesFlow network for each of the candidate models, in order to ensure that all models are comparably estimated within the same Bayesian framework. The training phase for each model took less than 12 hours on a laptop with GPU acceleration. In contrast, subsequent parameter recovery, calibration checks, and inference on the experimental data took a couple of seconds.

Our initial results suggest, in accordance with previous results [169], that the simple Lévy model (\mathcal{M}_1) yields a superior fit than the simple diffusion model (\mathcal{M}_2) for both experimental tasks. In addition, the complex Lévy model (\mathcal{M}_4) had a superior fit than the complex diffusion model (\mathcal{M}_3). Finally, each of the complex models (including inter-trial variability parameters) exhibited a superior fit than each of the simpler models. We speculate, that such a result might be explained by a particular property of the experiments: The longer duration possibly induces larger fluctuations in performance which is best captured by the inter-trial variability parameters (see [172] for a more in-depth interpretation).

8.3 INSIGHTS FROM BAYESIAN MODELING IN A ONE MILLION SAMPLE

In this work, we applied BayesFlow to elucidate cross-sectional age differences in cognitive parameters as indexed by the main diffusion model parameters. A large number of studies from the last decades have reported that processing speed, typically measured as mean response time (RT) in simple cognitive tasks, significantly slows down in old age and starts to decline in young and middle adulthood [78, 143]. We challenged this notion by carrying out a comprehensive model-based analysis on a massive, publicly available data set ($M > 1\,000\,000$) collected during the course of Project Implicit [173]. Notably, this sample was multiple orders of magnitude larger than the data sets used in all previous diffusion model studies combined. Accordingly, our approach was able to provide unique and robust findings on age-related patterns regarding processing speed, decision caution, and non-decision components of RTs.

Furthermore, applying Bayesian diffusion modeling to a sample of such magnitude appears to be a task insurmountable by standard Bayesian (or non-Bayesian, for that matter) methods. Thus, we resorted to BayesFlow with a deep invariant summary network for efficient amortized inference. In this way, fully Bayesian inference with BayesFlow (training and inference phase) on the entire sample took less than two days on a standard computer. We also estimated, that

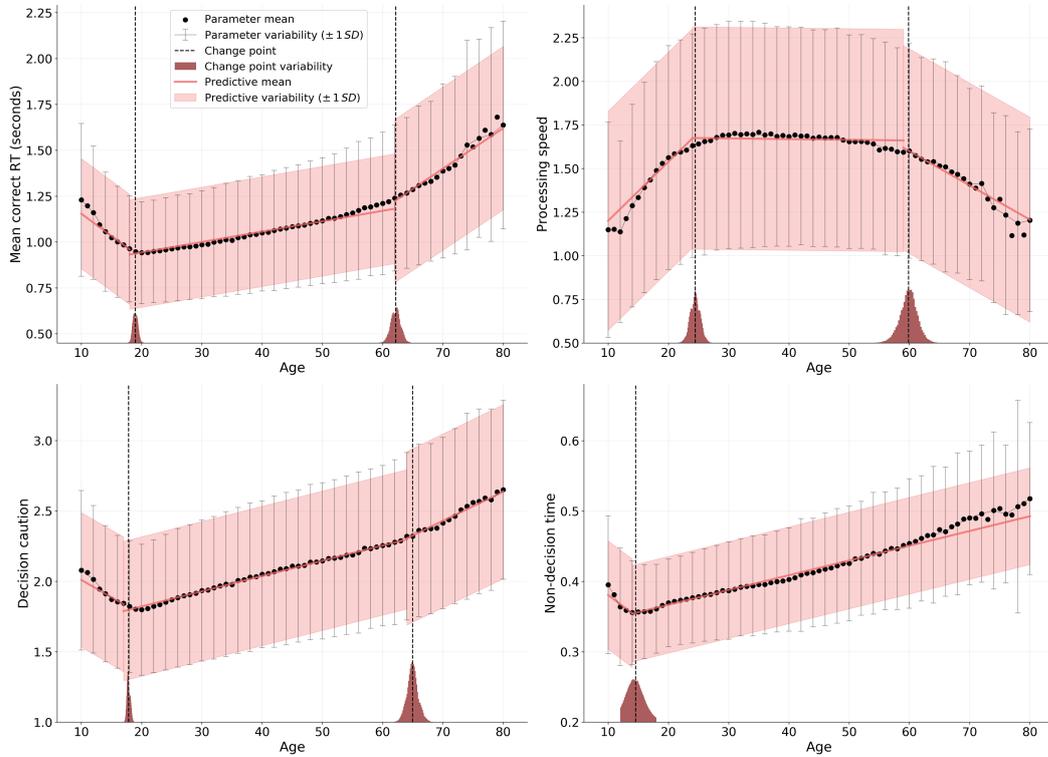


Figure 8.2: Mean correct response times (RTs) and main diffusion model parameters as a function of age. Black points indicate parameter means computed separately for each age (in years). Bars indicate standard deviations (only shown for every second year). Red lines denote the Bayesian piecewise ridge regression model's mean predictions, which describe the observed means fairly well. The shaded red region denotes the uncertainty (standard deviation) of the piece-wise model's predictions. The dashed lines indicate the mean change points estimated from the per-age-group averaged data, with the full posterior distributions (scaled for readability) of the change points shown at the bottom of each plot. Both the data- and model-implied standard deviations highlight the great variability within each year of age. Nevertheless, the year-specific means suggest a clear and consistent pattern for mean correct RT and each parameter. The figure depicts drift rates and boundary separations for the incongruent condition and non-decision times obtained from correct responses.

applying MCMC for the same number of posterior samples per data set (participant) would have taken more years than are currently available to any human scientist.

Our analysis pipeline followed the steps advocated by a principled Bayesian workflow [144] which ensure a transparent presentation of computational faithfulness and model adequacy. Following parameter estimation, we applied a Bayesian change point regression of each cognitive estimate (and also mean RT) on age. The results from one of the two experimental conditions are depicted in Figure 8.2 (parameter estimates in the other condition do not exhibit qualitatively different age-related patterns).

Importantly, our results suggest a clear non-linear association between drift rate (as an index of processing speed) and age, which was strikingly different than the one implied by mean RTs and far more informative than the age differences found in previous diffusion model studies (cf. Figure 8.2). Thus, our model-based analysis suggests a picture of age differences in cognitive parameters yielding a radically different implication than the one based on model-free analysis of raw RT data.

8.4 OUTBREAKFLOW: MODEL-BASED BAYESIAN INFERENCE OF DISEASE OUTBREAK DYNAMICS

In this work³, we applied a version of BayesFlow for dynamic (stateful) models to infer important disease characteristics and transmission dynamics of the Covid-19 pandemic in Germany. Inference of hidden disease-related parameters is of utmost importance in the case of new outbreaks in order to forecast their progression and guide effective public health measurements. Accordingly, mathematical models that provide a reliable representation of the processes driving the dynamics of an epidemic are an essential tool for this task (see for example [81]).

In order to account for the specific nature of the initial Covid-19 outbreak in Germany, we specified a custom compartmental model consisting of three sub-models: a disease model, an observation model, and an intervention model.

The *disease model* is represented by a system of non-linear ordinary differential equations (ODEs) comprising six compartments: susceptible (S), exposed (E - infected individuals who do not show symptoms and are not yet infectious), infected (I - symptomatic cases that are infectious), carrier (C - infectious individuals who recover without being detected), recovered (R), and dead (D).

The *intervention model* represents the time-varying transmission rate $\lambda(t)$. Following [34], we defined three change points encoding an assumed transmission rate reduction in response to public health measures (e.g., lockdown, social distancing) imposed by the German authorities. Each change point is represented by a piece-wise linear function with three degrees of freedom: the strength of interventions and the time interval (start and end point) for the effect to fully manifest itself.

The *observation model* represents the deviations between officially reported case counts and their true values. It comprises three sources of systematic and unsystematic errors: the reporting

³S. T. Radev, F. Graw, S. Chen, N. Mutters, V. Eichel, T. Bärnighausen, and U. Köthe. “Model-based Bayesian inference of disease outbreak with invertible neural networks”. *arXiv preprint arXiv:2010.00300*, 2020

8.4 OutbreakFlow: Model-Based Bayesian Inference of Disease Outbreak Dynamics

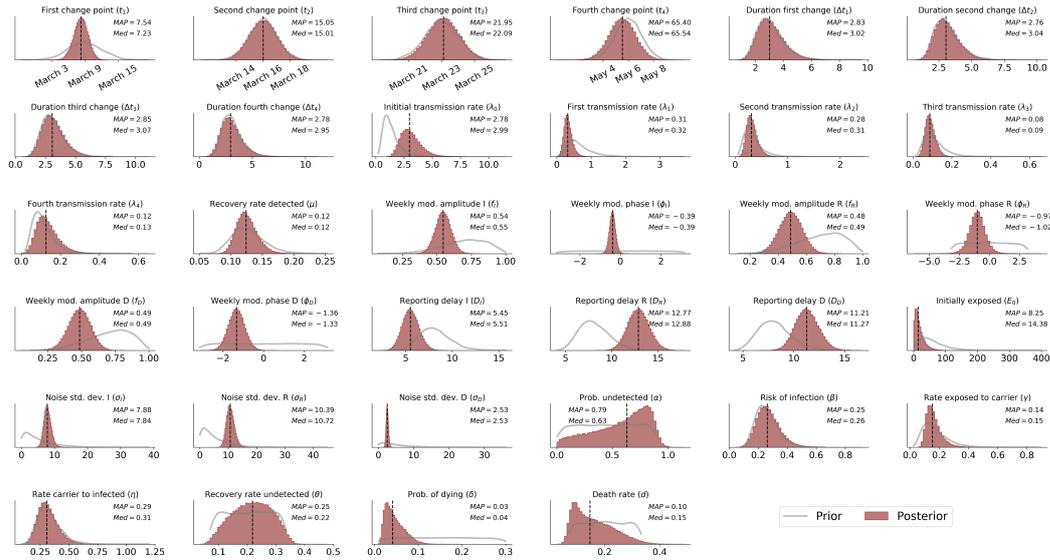


Figure 8.3: Marginal posteriors of all 34 model parameters inferred from data from entire Germany alongside median and MAP summary statistics. Gray lines depict prior distributions for comparison with the posteriors. Vertical dashed lines indicate posterior medians.

delay, the weekly modulation (since testing and reporting activities are considerably reduced on weekends), and a symmetric t -distributed noise term describing random fluctuations.

Due to the complexity of this composite model and the need to apply the same model repeatedly to different federal states, we resort to efficient simulation-based inference with BayesFlow (see [135] for details regarding network architectures). Furthermore, amortized Bayesian inference appears especially advantageous in epidemiological contexts, where the same model is estimated in multiple populations (countries, cultures) or at different scales (states, regions). Indeed, in the current application, we were able to demonstrate efficient amortized inference and excellent predictive performance with a single architecture applied simultaneously to epidemiological data from Germany as a whole and all sixteen German federal states.

Marginal parameter posteriors for Germany as a whole are depicted in Figure 8.3. Posterior predictions and forecasts for new infections, recoveries, and deaths are further depicted in Figure 8.4. We observe that median predictions of our model follow very closely the reported cumulative number of cases across all federal states. Furthermore, the officially reported cases are very well represented by the uncertainty bounds derived from the parameter posteriors, with prediction uncertainty growing as we move towards the future (cf. predictions after the dotted vertical lines in Figure 8.4). When interpreting these results, the reader should be aware that mechanistic models like ours only describe the average behavior of entire compartments, in contrast to agent-based models. Accordingly, the given CIs quantify our uncertainty about the inferred parameter averages and *cannot* be interpreted as a measure of the variability between individual cases.

Our estimates suggest that a considerable number of individuals (a fraction of 60-80 % of cases) might have gone undetected through the course of the Covid-19 outbreak in Germany, confirming results from previous studies in other countries [32, 115, 128]. However, our posteriors also suggest

8 Applications

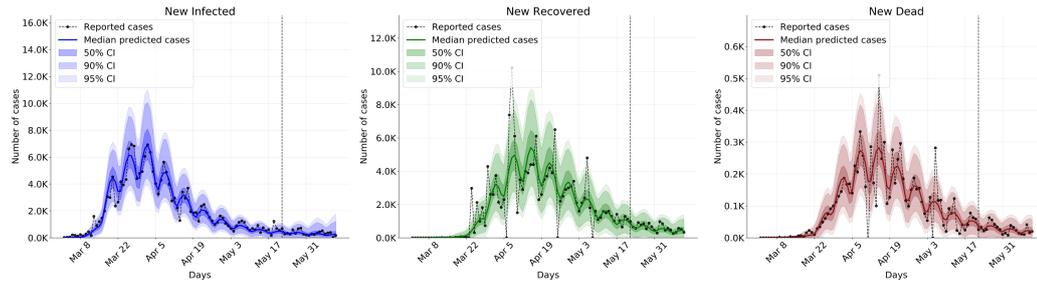


Figure 8.4: Model-based predictions and forecasts of new cases obtained by inferring model parameters from epidemiological data available for reported infected, assumed recovered and deaths by Covid-19 from entire Germany. Cases to the left of the vertical dashed line were used for posterior checking (model training) and cases to the right for posterior forecasts (predictions) on upcoming data.

that there is non-negligible uncertainty surrounding this estimate when derived in a purely model-based manner. Moreover, different summary statistics (e.g., means, medians, MAPs) derived from non-symmetric posteriors offer slightly different conclusions. The latter observation highlights the need to consider the full posteriors and corresponding credibility intervals when aiming to draw substantive conclusions and possible forecasts for the progression of the epidemic or the effect of specific public health interventions.

9 OUTLOOK

The current thesis presented frameworks and ideas for scaling up many steps of a complete Bayesian workflow with a focus on cognitive modeling. A cornerstone notion of this thesis was to employ generative neural networks for amortized Bayesian parameter inference, model comparison and validation when working with intractable simulators whose behavior as a whole is too complex to be described analytically. We presented various frameworks for tackling different types of models (e.g., stateless vs. stateful) and Bayesian tasks (e.g., parameter estimation, model comparison, model calibration). A common theme was splitting Bayesian analysis into two conceptual phases: i) a training phase, in which the networks gradually become domain experts in solving the Bayesian tasks they are optimized for, and ii) a downstream inference phase, in which the networks are efficiently applied to extract information from real-world observations about quantities of interest (e.g., model parameters or model plausibility). Further, we explored potential developments towards meta-amortized Bayesian inference and discussed related challenges standing in the way of such a generalized framework. Finally, we presented some applications of BayesFlow to a number of complex estimation problems. In the following, we briefly go through some further topics left for future research beyond those mentioned in previous chapters.

Dynamic parameters In this thesis, we have focused exclusively on models defined by a fixed number of parameters θ . However, some dynamic models might incorporate parameters which are a function of time (or another variable), implying that a different set of parameters $\theta(t)$ is available at each time point t . For instance, realistic spatio-temporal models of disease outbreaks strive to capture relevant disease characteristics $\theta(t, s)$ as a function of time t and space s [24], introducing yet another dimension to the parameter space. While this setting poses no inherent problems for our evidential framework for model comparison, it presents a challenge for parameter estimation with BayesFlow.

One approach to amortized inference with such models would be to re-parameterize the problem so that we can estimate a fixed-size set of parameters ω of which the time-varying parameter vector $\theta(\omega, t)$ is a deterministic function. This approach is then easily amenable to amortized inference with BayesFlow, and is the one we followed in [31] for recovering trial-by-trial information processing dynamics or in [135] for recovering the time-varying transmission rate of Covid-19 in Germany. However, not all models naturally admit such a re-parameterization, so estimating the original $\theta(t)$ might be inescapable in these cases. Thus, another viable approach is to modify the summary and inference networks in order to capture the dynamic structure of the problem. For instance, the inference network can be easily implemented as a generative recurrent network which outputs a latent embedding $z(t)$ for each $\theta(t)$. In this way, it can interact with a recurrent summary network in the form of a probabilistic sequence-to-sequence architecture [124, 151] Transformer networks utilizing neural attention mechanisms [83, 164] appear as another promising option for tackling dynamic models without recurrent networks.

Hyperparameter optimization A common aspect of our frameworks for amortized Bayesian inference is the potentially large number of hyperparameter settings that might require fine-tuning by the user for optimal performance on a given Bayesian task. For instance, the most important hyperparameters for BayesFlow are: optimizer settings (e.g., learning rate, adaptive weights, decay); summary network design (e.g., type of modular architecture, number of layers and neurons in each module); inference network design (e.g., type of flow, number of layers, structure of each internal network); training schedule (e.g., online vs. offline learning). This makes general experience with neural networks highly advantageous when working with our frameworks, to say the least.

So far, in our simulated experiments and applications, we could empirically ascertain that some hyperparameters are more important than others. For instance, optimizer settings appear to be vital for stable training. When working with the Adam optimizer [84], smaller learning rates (i.e., $\alpha < 0.001$) and the inclusion of learning rate decay generally lead to more stable convergence than larger learning rates and no decay. On the other hand, using larger networks consisting of 3 to 10 coupling layers does not seem to hurt performance or destabilize training, even if the simulator to be inverted is relatively simple [133]. Based on our results, we expect that a single architecture should be able to perform well on similar simulators from a given domain (e.g., one architecture for all EAMs [172] or one architecture for all compartmental models [135]). Future research should investigate the impact of modern hyperparameter optimization methods, such as Bayesian optimization [41]. Moreover, Bayesian optimization, or other black-box optimization methods or search algorithms can easily be integrated into our frameworks (e.g., in a pre-training phase with utilizing a small number of simulations).

Optimal experimental design Optimal experimental design (OED) in a Bayesian context is a mathematical framework for making efficient use of limited experimental resources when performing Bayesian modeling [21, 46, 110]. A majority of OED approaches revolve around the notion of expected information gain (EIG), which quantifies the expected reduction in entropy (uncertainty) when replacing the prior with the posterior under the marginal distribution over experimental observations. For instance, in static design optimization (DO), a researcher sets up a simulation involving different models in different experimental contexts and picks the configuration which yields the highest EIG. No further optimization happens during the actual experiment. Differently, in adaptive design optimization (ADO), the EIG is computed (estimated) on each trial and subsequent trials are chosen in order to maximize the discriminability between candidate models or the sharpness of the posterior in a single-model setting.

Unfortunately, obtaining accurate approximations of the EIG even for simpler models is computationally demanding and nearly infeasible with non-amortized methods for Bayesian updating. Variational OED offers a promising approach for amortizing different aspects of OED [46]. Moreover, utilizing neural networks for efficiently maximizing a lower bound on the EIG (i.e., variational autoencoders, VAEs, [87]), variational methods are able to yield considerable efficiency gains. However, vanilla VAEs maximizing a lower bound on the actual (intractable) criterion, the so called ELBO criterion, suffer from some rather consequential problems, as aptly demonstrated by [178]. Further, in contrast to normalizing flows realized via invertible neural networks, vanilla variational methods offer no theoretical guarantee for learning the correct target posterior when employed for the task of Bayesian updating.

Our proposed frameworks for amortized Bayesian parameter estimation and model comparison appear suitable for amortizing OED. For instance, BayesFlow can be adapted for estimating the EIG in either static DO or ADO. Further, our evidential framework can be employed in multi-model DO or ADO contexts. Such an integration is possible, since our networks can be augmented to process arbitrary contextual information. Moreover, we can emulate Bayesian updating by training the networks with a variable number of observations N (using algorithmically aligned networks), such that amortization over increasing N is enabled during inference. Thus, future research should investigate the utility of our frameworks for amortizing OED and compare them to variational approaches.

Model-aware learning Finally, our frameworks currently operate in a model-agnostic manner, that is, the neural estimators treat the simulator purely as a black-box data generator. For researchers, on the other hand, the neural networks (in addition to reality itself) are uninterpretable black-boxes while the simulator serves, at least in theory, as a human-interpretable, white-box computational model. Thus, it is possible that the networks can profit from some prior “knowledge” of the model’s structure (e.g., in the form of gradients or other information) or generative scope in order to learn even more efficiently the resulting probabilistic mappings. The networks themselves could guide the model to produce more realistic artificial observations, for instance, by restricting the generative scope of the simulator.

Sequential neural posterior estimation (SNPE, [60]) methods offer a neat way to gradually transform the prior $p(\boldsymbol{\theta})$ into a sharper proposal $\hat{p}(\boldsymbol{\theta})$ (eventually becoming the target posterior) through iterative refinement. In this way, the generative scope of the joint Bayesian model $p(\boldsymbol{\theta}, \boldsymbol{x})$ also becomes narrower, concentrating around the actually observed data $\boldsymbol{x}^{(obs)}$. Thus, SNPE methods implement one promising form of network-simulator interaction. One consequence of this approach, however, is that such an interaction necessarily reduces amortization to a case-wise level, since the neural density estimator needs to be re-trained for each observed data set. Needless to say, such repetitions become increasingly computationally taxing in data-rich applications. Future research should therefore explore other forms of *model-aware* learning, which enable model-wise or even meta-amortized Bayesian inference.

10 CONCLUSION

Reality is noisy and messy, and there is no grand simulator of things in sight. What is more, our models can only restrict their scope to increasingly smaller empirical nooks, solving, at best, tiny fractions of an infinite jigsaw puzzle. As we, researchers in need of cognition, go through the process of building new models and discarding old ones, we require the right tools to foster our epistemic achievements. At the time of writing, deep learning continues to enjoy a vibrant hype, refurbishing the methodological equipment of many quantitative sciences. The behavioral sciences, despite being more resistant to change than fellow disciplines, are also enjoying their fair share of the rapidly expanding trend. When it comes to model-based inference, deep learning innovations are currently transforming the way models are fit to data and employed for drawing substantive conclusions or deriving reliable forecasts. Moreover, uncertainty (an ancient concept) and its quantification are becoming more and more important in deep learning theory and practice. Bayesian methods, deeply rooted in probability theory, are currently viewed by many researchers as the gold-standard for uncertainty-aware inference, but other approaches or generalizations might push through in the not-so-distant future. In a way, this thesis ventured into a discourse between deep learning and scientific modeling with a focus on cognitive science and mathematical psychology. It brought together ideas for dramatically accelerating Bayesian inference by using non-Bayesian neural networks designed to deal with the data types encountered by researchers working in various areas of knowledge. The general idea of using black-box estimators to learn white-box scientific models from computer simulations is certainly not new, but is still largely underutilized in the behavioral and cognitive sciences. Most importantly, future research should further foster the discourse between deep learning (or *artificial intelligence*, in marketing jargon) and the behavioral sciences due to the potential upside of such a creative interaction. Indeed, human researchers and decision-makers can definitely learn something from deep learning agents surpassing human performance in various real-world tasks. On the other hand, neural networks can also learn something from studying the structure of human behavior and cognition. Luckily, the global connectedness of the modern world makes such mutual learning a rather effortless endeavor. Finally, models of human behavior and cognition need to come to terms with the buzzword of the century: *complexity*. We have begun to realize, that simple models can only do so much in aiding our understanding of emergent phenomena. On the other hand, we have grown suspicious of opaque, overparameterized neural networks capable of solving overly specific tasks. It is thus very well possible, that future models of cognition and behavior become more and more uninterpretable (i.e., more black-boxy) in the pursuit of complexity. On the other hand, future neural network might become more and more interpretable (i.e., less black-boxy) due to our need for cognition. As always, predictions are hard, especially about the future.

10 Conclusion

BIBLIOGRAPHY

1. L. Abbott and T. B. Kepler. “Model neurons: from hodgkin-huxley to hopfield”. In: *Statistical mechanics of neural networks*. Springer, 1990, pp. 5–18.
2. J. R. Anderson. *How can the human mind occur in the physical universe?* Vol. 3. Oxford University Press, 2009.
3. L. Ardizzone, J. Kruse, S. Wirkert, D. Rahner, E. W. Pellegrini, R. S. Klessen, L. Maier-Hein, C. Rother, and U. Köthe. “Analyzing inverse problems with invertible neural networks”. In: *Intl. Conf. on Learning Representations*. 2019.
4. L. Ardizzone, C. Lüth, J. Kruse, C. Rother, and U. Köthe. “Guided image generation with conditional invertible neural networks”. *arXiv preprint arXiv:1907.02392*, 2019.
5. M. Belkin, D. Hsu, S. Ma, and S. Mandal. “Reconciling modern machine-learning practice and the classical bias–variance trade-off”. *Proceedings of the National Academy of Sciences* 116:32, 2019, pp. 15849–15854.
6. E. A. Berg. “A simple objective technique for measuring flexibility in thinking”. *The Journal of general psychology* 39:1, 1948, pp. 15–22.
7. J. O. Berger. “Robust Bayesian analysis: sensitivity to the prior”. *Journal of statistical planning and inference* 25:3, 1990, pp. 303–328.
8. M. Betancourt. “A conceptual introduction to Hamiltonian Monte Carlo”. *arXiv preprint arXiv:1701.02434*, 2017.
9. A. Biedermann and F. Taroni. “Bayesian networks for evaluating forensic DNA profiling evidence: a review and guide to literature”. *Forensic Science International: Genetics* 6:2, 2012, pp. 147–157.
10. S. Bieringer, A. Butter, T. Heimel, S. Höche, U. Köthe, T. Plehn, and S. T. Radev. “Measuring QCD splittings with invertible networks”. *arXiv preprint arXiv:2012.09873*, 2020.
11. D. Bigoni, O. Zahm, A. Spantini, and Y. Marzouk. “Greedy inference with layers of lazy maps”. *arXiv:1906.00031*, 2019.
12. D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. “Variational inference: A review for statisticians”. *Journal of the American statistical Association* 112:518, 2017, pp. 859–877.
13. B. Bloem-Reddy and Y. W. Teh. “Probabilistic symmetries and invariant neural networks”. *Journal of Machine Learning Research* 21:90, 2020, pp. 1–61.
14. M. G. Blum and O. François. “Non-linear regression models for Approximate Bayesian Computation”. *Statistics and computing* 20:1, 2010, pp. 63–73.

15. O. Breidbach. “Weltordnungen und Körperwelten. Das Tableau des Gewussten und seine Repräsentation bei Robert Fludd”. *Schramm, H./Schwarte, L./Lazardzig, J.(Hg.): Instrumente in Kunst und Wissenschaft. Zur Architektonik kultureller Grenzen im 17*, 2006, pp. 41–65.
16. A. N. Burkitt. “A review of the integrate-and-fire neuron model: I. Homogeneous synaptic input”. *Biological cybernetics* 95:1, 2006, pp. 1–19.
17. P.-C. Bürkner, J. Gabry, and A. Vehtari. “Approximate leave-future-out cross-validation for Bayesian time series models”. *Journal of Statistical Computation and Simulation* 90:14, 2020, pp. 2499–2523.
18. G. Buzsaki. *Rhythms of the Brain*. Oxford University Press, 2006.
19. B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell. “Stan: A probabilistic programming language”. *Journal of statistical software* 76:1, 2017.
20. R. Caruana. “Multitask learning”. *Machine learning* 28:1, 1997, pp. 41–75.
21. D. R. Cavagnaro, J. I. Myung, M. A. Pitt, and J. V. Kujala. “Adaptive design optimization: A mutual information-based approach to model discrimination in cognitive science”. *Neural computation* 22:4, 2010, pp. 887–905.
22. P. C. Cheeseman. “In defense of probability.” In: *IJCAI*. Vol. 2. 1985, pp. 1002–1009.
23. L. Chen, G. Zhang, and E. Zhou. “Fast greedy map inference for determinantal point process to improve recommendation diversity”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 5622–5633.
24. Y. Chen, Q. Li, H. Karimian, X. Chen, and X. Li. “Spatio-temporal distribution characteristics and influencing factors of COVID-19 in China”. *Scientific Reports* 11:1, 2021, pp. 1–12.
25. K. Choi, M. Wu, N. Goodman, and S. Ermon. “Meta-amortized variational inference and learning”. In: *International Conference on Learning Representations*. 2019.
26. J. M. J. da Costa, H. R. B. Orlande, and W. B. da Silva. “Model selection and parameter estimation in tumor growth models using approximate Bayesian computation-ABC”. *Computational and Applied Mathematics* 37:3, 2018, pp. 2795–2815.
27. K. Cranmer, J. Brehmer, and G. Louppe. “The frontier of simulation-based inference”. *Proceedings of the National Academy of Sciences*, 2020.
28. F. Crick and J. Clark. “The astonishing hypothesis”. *Journal of Consciousness Studies* 1:1, 1994, pp. 10–16.
29. K. Csilléry, M. G. Blum, O. E. Gaggiotti, and O. François. “Approximate Bayesian computation (ABC) in practice”. *Trends in Ecology & Evolution* 25:7, 2010, pp. 410–418.
30. M. D’Alessandro, G. Gallitto, A. Greco, and L. Lombardi. “A Joint modelling approach to analyze risky decisions by means of diffusion tensor imaging and behavioural data”. *Brain sciences* 10:3, 2020, p. 138.

31. M. D'Alessandro, S. T. Radev, A. Voss, and L. Lombardi. "A Bayesian brain model of adaptive behavior: an application to the Wisconsin Card Sorting Task". *PeerJ* 8, 2020, e10316.
32. M. Day. *Covid-19: four fifths of cases are asymptomatic, China figures indicate*. 2020.
33. B. De Finetti. *Fondamenti logici del ragionamento probabilistico*. Azzoguidi, 1930.
34. J. Dehning, J. Zierenberg, F. P. Spitzner, M. Wibral, J. P. Neto, M. Wilczek, and V. Priesemann. "Inferring change points in the spread of COVID-19 reveals the effectiveness of interventions". *Science*, 2020.
35. T. S. Deisboeck and G. S. Stamatakos. *Multiscale cancer modeling*. CRC press, 2010.
36. G. Detommaso, J. Kruse, L. Ardizzone, C. Rother, U. Köthe, and R. Scheichl. "HINT: Hierarchical invertible neural transport for general and sequential Bayesian inference". *arXiv:1905.10687*, 2019.
37. L. Dinh, D. Krueger, and Y. Bengio. "Nice: Non-linear independent components estimation". *arXiv preprint arXiv:1410.8516*, 2014.
38. L. Dinh, J. Sohl-Dickstein, and S. Bengio. "Density estimation using real nvp". *arXiv preprint arXiv:1605.08803*, 2016.
39. M. P. DOUGHERTY, C. GETTYS, and E. OGDEN. "MINERVA-DM: A memory processes model for judgments of likelihood". *Psychological review* 106:1, 1999, pp. 180–209.
40. C. Durkan, I. Murray, and G. Papamakarios. "On contrastive learning for likelihood-free inference". In: *International Conference on Machine Learning*. PMLR. 2020, pp. 2771–2781.
41. K. Eggensperger, M. Feurer, F. Hutter, J. Bergstra, J. Snoek, H. Hoos, and K. Leyton-Brown. "Towards an empirical foundation for assessing bayesian optimization of hyperparameters". In: *NIPS workshop on Bayesian Optimization in Theory and Practice*. Vol. 10. 2013, p. 3.
42. S. Engblom, R. Eriksson, and S. Widgren. "Bayesian epidemiological modeling over high-resolution network data". *Epidemics* 32, 2020, p. 100399.
43. A. Etz, Q. F. Gronau, F. Dablander, P. A. Edelsbrunner, and B. Baribault. "How to become a Bayesian in eight easy steps: An annotated reading list". *Psychonomic Bulletin & Review* 25:1, 2018, pp. 219–234.
44. N. J. Evans and M. Servant. "A comparison of conflict diffusion models in the flanker task through pseudolikelihood Bayes factors." *Psychological review* 127:1, 2020, p. 114.
45. S. Farrell and S. Lewandowsky. *Computational modeling of cognition and behavior*. Cambridge University Press, 2018.
46. A. Foster, M. Jankowiak, E. Bingham, P. Horsfall, Y. W. Teh, T. Rainforth, and N. Goodman. "Variational bayesian optimal experimental design". *arXiv preprint arXiv:1903.05480*, 2019.
47. J. Friedman, T. Hastie, R. Tibshirani, et al. *The elements of statistical learning*. Vol. 1. 10. Springer series in statistics New York, 2001.

Bibliography

48. K. Friston. “A theory of cortical responses”. *Philosophical transactions of the Royal Society B: Biological sciences* 360:1456, 2005, pp. 815–836.
49. J. Gabry, D. Simpson, A. Vehtari, M. Betancourt, and A. Gelman. “Visualization in Bayesian workflow”. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 182:2, 2019, pp. 389–402.
50. F. Galton. *Vox populi*. 1907.
51. A. E. Gelfand. “Gibbs sampling”. *Journal of the American statistical Association* 95:452, 2000, pp. 1300–1304.
52. A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian data analysis*. CRC press, 2013.
53. A. Gelman, J. Hwang, and A. Vehtari. “Understanding predictive information criteria for Bayesian models”. *Statistics and computing* 24:6, 2014, pp. 997–1016.
54. A. Gelman and X.-L. Meng. “Simulating normalizing constants: From importance sampling to bridge sampling to path sampling”. *Statistical science*, 1998, pp. 163–185.
55. A. Gelman, D. B. Rubin, et al. “Inference from iterative simulation using multiple sequences”. *Statistical science* 7:4, 1992, pp. 457–472.
56. A. Gelman, A. Vehtari, D. Simpson, C. C. Margossian, B. Carpenter, Y. Yao, L. Kennedy, J. Gabry, P.-C. Bürkner, and M. Modrák. “Bayesian workflow”. *arXiv preprint arXiv:2011.01808*, 2020.
57. F. A. Gers, J. A. Schmidhuber, and F. A. Cummins. “Learning to forget: continual prediction with LSTM”. *Neural Computation* 12:10, 2000, pp. 2451–2471.
58. T. Gneiting, F. Balabdaoui, and A. E. Raftery. “Probabilistic forecasts, calibration and sharpness”. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69:2, 2007, pp. 243–268.
59. T. Gneiting and A. E. Raftery. “Strictly proper scoring rules, prediction, and estimation”. *Journal of the American Statistical Association* 102:477, 2007, pp. 359–378.
60. D. Greenberg, M. Nonnenmacher, and J. Macke. “Automatic posterior transformation for likelihood-free inference”. In: *International Conference on Machine Learning*. 2019, pp. 2404–2414.
61. A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. “A kernel two-sample test”. *The Journal of Machine Learning Research* 13:1, 2012, pp. 723–773.
62. Q. F. Gronau, A. Sarafoglou, D. Matzke, A. Ly, U. Boehm, M. Marsman, D. S. Leslie, J. J. Forster, E.-J. Wagenmakers, and H. Steingroever. “A tutorial on bridge sampling”. *Journal of mathematical psychology* 81, 2017, pp. 80–97.
63. P. Grünwald, T. Van Ommen, et al. “Inconsistency of Bayesian inference for misspecified linear models, and a proposal for repairing it”. *Bayesian Analysis* 12:4, 2017, pp. 1069–1103.
64. C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. “On calibration of modern neural networks”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 1321–1330.

65. F. Hartig, J. M. Calabrese, B. Reineking, T. Wiegand, and A. Huth. “Statistical inference for stochastic simulation models—theory and application”. *Ecology letters* 14:8, 2011, pp. 816–827.
66. W. K. Hastings. “Monte Carlo sampling methods using Markov chains and their applications”, 1970.
67. R. Heaton. “Wisconsin card sorting test manual; revised and expanded”. *Psychological Assessment Resources*, 1981, pp. 5–57.
68. J. Hermans, V. Begy, and G. Louppe. “Likelihood-free MCMC with approximate likelihood ratios”. *arXiv preprint*, 2019.
69. J. M. Hernández-Lobato and R. Adams. “Probabilistic backpropagation for scalable learning of bayesian neural networks”. In: *International Conference on Machine Learning*. PMLR, 2015, pp. 1861–1869.
70. A. L. Hodgkin and A. F. Huxley. “A quantitative description of membrane current and its application to conduction and excitation in nerve”. *The Journal of physiology* 117:4, 1952, pp. 500–544.
71. M. D. Hoffman and A. Gelman. “The No-U-Turn sampler: adaptively setting path lengths in Hamiltonian Monte Carlo.” *J. Mach. Learn. Res.* 15:1, 2014, pp. 1593–1623.
72. T. Hu, Z. Chen, H. Sun, J. Bai, M. Ye, and G. Cheng. “Stein neural sampler”. *arXiv preprint arXiv:1810.03545*, 2018.
73. F. Huber, T. Krisztin, and P. Piribauer. “Forecasting global equity indices using large Bayesian VARs”. *Bulletin of Economic Research* 69:3, 2017, pp. 288–308.
74. E. Hüllermeier and W. Waegeman. “Aleatoric and epistemic uncertainty in machine learning: A tutorial introduction”. *arXiv preprint arXiv:1910.09457*, 2019.
75. E. M. Izhikevich. “Simple model of spiking neurons”. *IEEE Transactions on neural networks* 14:6, 2003, pp. 1569–1572.
76. E. M. Izhikevich. “Which model to use for cortical spiking neurons?” *IEEE transactions on neural networks* 15:5, 2004, pp. 1063–1070.
77. E. T. Jaynes. *Probability theory: The logic of science*. Cambridge university press, 2003.
78. A. R. Jensen. *Clocking the mind: Mental chronometry and individual differences*. Elsevier, 2006.
79. B. Jiang, T.-y. Wu, C. Zheng, and W. H. Wong. “Learning summary statistic for approximate Bayesian computation via deep neural network”. *Statistica Sinica*, 2017, pp. 1595–1618.
80. A. Jsang. *Subjective Logic: A formalism for reasoning under uncertainty*. Springer Publishing Company, Incorporated, 2018.
81. M. J. Keeling and P. Rohani. *Modeling infectious diseases in humans and animals*. Princeton University Press, 2011.

82. A. Kendall and Y. Gal. “What uncertainties do we need in Bayesian deep learning for computer vision?” In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 2017, pp. 5580–5590.
83. T. H. Kim, M. S. Sajjadi, M. Hirsch, and B. Scholkopf. “Spatio-temporal transformer network for video restoration”. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018, pp. 106–122.
84. D. P. Kingma and J. Ba. “Adam: A method for stochastic optimization”. *arXiv preprint arXiv:1412.6980*, 2014.
85. D. P. Kingma, T. Salimans, and M. Welling. “Variational dropout and the local reparameterization trick”. In: *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*. 2015, pp. 2575–2583.
86. D. P. Kingma, M. Welling, et al. “An Introduction to Variational Autoencoders”. *Foundations and Trends® in Machine Learning* 12:4, 2019, pp. 307–392.
87. D. P. Kingma and M. Welling. “Auto-encoding variational bayes”. *arXiv preprint arXiv:1312.6114*, 2013.
88. D. P. Kingma and P. Dhariwal. “Glow: Generative flow with invertible 1x1 convolutions”, 2018. arXiv: [1807.03039](https://arxiv.org/abs/1807.03039) [stat.ML].
89. D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling. “Improved variational inference with inverse autoregressive flow”. In: *Advances in Neural Information Processing Systems*. 2016, pp. 4743–4751.
90. E. Klinger, D. Rickert, and J. Hasenauer. “pyABC: distributed, likelihood-free inference”. *Bioinformatics* 34:20, 2018, pp. 3591–3593.
91. D. C. Knill and A. Pouget. “The Bayesian brain: the role of uncertainty in neural coding and computation”. *TRENDS in Neurosciences* 27:12, 2004, pp. 712–719.
92. I. Kobyzev, S. Prince, and M. Brubaker. “Normalizing flows: An introduction and review of current methods”. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020, pp. 1–1. ISSN: 1939-3539. DOI: [10.1109/tpami.2020.2992934](https://doi.org/10.1109/tpami.2020.2992934). URL: <http://dx.doi.org/10.1109/TPAMI.2020.2992934>.
93. L. Konicar, S. Radev, K. Prillinger, M. Klöbl, R. Diehm, N. Birbaumer, R. Lanzenberger, P. Plener, and L. Poustka. “Volitional modification of brain activity in adolescents with Autism Spectrum Disorder: A Bayesian analysis of Slow Cortical Potential neurofeedback”. *NeuroImage: Clinical*, 2021, p. 102557.
94. M. von Krause, S. T. Radev, and A. Voss. “Processing speed is high until age 60: insights from Bayesian modeling in a one million sample (with a little help of deep learning)”, Manuscript submitted for publication.
95. S. Levy. “Wealthy people and fat tails: An explanation for the Lévy distribution of stock returns”, 1998.
96. Y. Li and Y. Gal. “Dropout inference in Bayesian neural networks with alpha-divergences”. In: *International conference on machine learning*. PMLR. 2017, pp. 2052–2061.

97. H. Lin and S. Jegelka. “Resnet with one-neuron hidden layers is a universal approximator”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 6169–6178.
98. R. Liu and J. Zou. “The effects of memory replay in reinforcement learning”. In: *2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE. 2018, pp. 478–485.
99. C. Louizos and M. Welling. “Multiplicative normalizing flows for variational bayesian neural networks”. In: *International Conference on Machine Learning*. PMLR. 2017, pp. 2218–2227.
100. J.-M. Lueckmann, G. Bassetto, T. Karaletsos, and J. H. Macke. “Likelihood-free inference with emulator networks”. In: *Symposium on Advances in Approximate Bayesian Inference*. PMLR. 2019, pp. 32–53.
101. J.-M. Lueckmann, P.J. Gonçalves, G. Bassetto, K. Öcal, M. Nonnenmacher, and J. H. Macke. “Flexible statistical inference for mechanistic models of neural dynamics”. In: *Proceedings of the 31st International Conference on Neural Information Processing Systems*. 2017, pp. 1289–1299.
102. D. J. MacKay. “Bayesian neural networks and density networks”. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 354:1, 1995, pp. 73–80.
103. D. J. MacKay and D. J. Mac Kay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
104. J.-M. Marin, P. Pudlo, A. Estoup, and C. Robert. *Likelihood-free model choice*. Chapman and Hall/CRC Press Boca Raton, FL, 2018.
105. D. W. Marquardt. “An algorithm for least-squares estimation of nonlinear parameters”. *Journal of the society for Industrial and Applied Mathematics* 11:2, 1963, pp. 431–441.
106. D. Marr. *Vision: A computational investigation into the human representation and processing of visual information*. MIT press, 2010.
107. U. K. Mertens, A. Voss, and S. Radev. “ABrox—A user-friendly Python module for approximate Bayesian computation with a focus on model comparison”. *PLoS one* 13:3, 2018, e0193981.
108. M. Mestdagh, S. Verdonck, K. Meers, T. Loossens, and F. Tuerlinckx. “Prepaid parameter estimation without likelihoods”. *PLoS computational biology* 15:9, 2019, e1007181.
109. N. Metropolis, A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. “Equation of state calculations by fast computing machines”. *The journal of chemical physics* 21:6, 1953, pp. 1087–1092.
110. J. I. Myung, D. R. Cavagnaro, and M. A. Pitt. “A tutorial on adaptive design optimization”. *Journal of mathematical psychology* 57:3-4, 2013, pp. 53–67.
111. B. Neal, S. Mittal, A. Baratin, V. Tantia, M. Scicluna, S. Lacoste-Julien, and I. Mitliagkas. “A modern take on the bias-variance tradeoff in neural networks”. *arXiv preprint arXiv:1810.08591*, 2018.

112. R. M. Neal et al. "MCMC using Hamiltonian dynamics". *Handbook of markov chain monte carlo* 2:11, 2011, p. 2.
113. K. A. Nicoli, S. Nakajima, N. Strodthoff, W. Samek, K.-R. Müller, and P. Kessel. "Asymptotically unbiased estimation of physical observables with neural samplers". *Physical Review E* 101:2, 2020, p. 023304.
114. S. Nowozin, B. Cseke, and R. Tomioka. "f-gan: Training generative neural samplers using variational divergence minimization". In: *Advances in neural information processing systems*. 2016, pp. 271–279.
115. D. P. Oran and E. J. Topol. "Prevalence of Asymptomatic SARS-CoV-2 Infection: A Narrative Review". *Annals of Internal Medicine*, 2020.
116. P. Orbanz and D. M. Roy. "Bayesian models of graphs, arrays and other exchangeable random structures". *IEEE transactions on pattern analysis and machine intelligence* 37:2, 2014, pp. 437–461.
117. B. Paige and F. Wood. "Inference networks for sequential Monte Carlo in graphical models". *International Conference on Machine Learning* 48, 2016, pp. 3040–3049.
118. J. J. Palestro, G. Bahg, P. B. Sederberg, Z.-L. Lu, M. Steyvers, and B. M. Turner. "A tutorial on joint models of neural and behavioral measures of cognition". *Journal of Mathematical Psychology* 84, 2018, pp. 20–48.
119. J. J. Palestro, P. B. Sederberg, A. F. Osth, T. Van Zandt, and B. M. Turner. *Likelihood-free methods for cognitive science*. Springer, 2018.
120. S. Palminteri, V. Wyart, and E. Koechlin. "The importance of falsification in computational cognitive modeling". *Trends in cognitive sciences* 21:6, 2017, pp. 425–433.
121. G. Papamakarios and I. Murray. "Fast ϵ -free inference of simulation models with Bayesian conditional density estimation". In: *Proceedings of the 30th International Conference on Neural Information Processing Systems*. 2016, pp. 1036–1044.
122. G. Papamakarios, D. Sterratt, and I. Murray. "Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows". In: *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR. 2019, pp. 837–848.
123. M. Park, W. Jitkrittum, and D. Sejdinovic. "K2-ABC: Approximate Bayesian computation with kernel embeddings". In: *Artificial Intelligence and Statistics*. PMLR. 2016, pp. 398–407.
124. S. H. Park, B. Kim, C. M. Kang, C. C. Chung, and J. W. Choi. "Sequence-to-sequence prediction of vehicle trajectory via LSTM encoder-decoder architecture". In: *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE. 2018, pp. 1672–1678.
125. M. D. Parno and Y. M. Marzouk. "Transport Map Accelerated Markov Chain Monte Carlo". *SIAM/ASA Journal on Uncertainty Quantification* 6:2, 2018, pp. 645–682.
126. J. Piironen and A. Vehtari. "Comparison of Bayesian predictive methods for model selection". *Statistics and Computing* 27:3, 2017, pp. 711–735.

127. M. Plummer et al. “JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling”. In: *Proceedings of the 3rd international workshop on distributed statistical computing*. Vol. 124. 125.10. Vienna, Austria. 2003, pp. 1–10.
128. P. Poletti, M. Tirani, D. Cereda, F. Trentini, G. Guzzetta, G. Sabatino, V. Marziano, A. Castrofino, F. Grosso, G. Del Castillo, et al. “Probability of symptoms and critical disease after SARS-CoV-2 infection”. *arXiv preprint arXiv:2006.08471*, 2020.
129. M. Pospischil, Z. Piwkowska, T. Bal, and A. Destexhe. “Comparison of different neuron models to conductance-based post-stimulus time histograms obtained in cortical pyramidal cells using dynamic-clamp in vitro”. *Biological cybernetics* 105:2, 2011, p. 167.
130. M. Pospischil, M. Toledo-Rodriguez, C. Monier, Z. Piwkowska, T. Bal, Y. Frégnac, H. Markram, and A. Destexhe. “Minimal Hodgkin–Huxley type models for different classes of cortical and thalamic neurons”. *Biological cybernetics* 99:4-5, 2008, pp. 427–441.
131. W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. *Numerical recipes 3rd edition: The art of scientific computing*. Cambridge university press, 2007.
132. P. Pudlo, J.-M. Marin, A. Estoup, J.-M. Cornuet, M. Gautier, and C. P. Robert. “Reliable ABC model choice via random forests”. *Bioinformatics* 32:6, 2015, pp. 859–866.
133. S. T. Radev, U. K. Mertens, A. Voss, L. Ardizzone, and U. Köthe. “BayesFlow: Learning complex stochastic models with invertible neural networks”. *IEEE Transactions on Neural Networks and Learning Systems*, 2020, pp. 1–15. DOI: [10.1109/TNNLS.2020.3042395](https://doi.org/10.1109/TNNLS.2020.3042395).
134. S. T. Radev, M. D’Alessandro, P.-C. Bürkner, U. K. Mertens, A. Voss, and U. Köthe. “Amortized Bayesian model comparison with evidential deep learning”, Manuscript submitted for publication in *IEEE Transactions on Neural Networks and Learning Systems*.
135. S. T. Radev, F. Graw, S. Chen, N. Mutters, V. Eichel, T. Bärnighausen, and U. Köthe. “Model-based Bayesian inference of disease outbreak with invertible neural networks”. *arXiv preprint arXiv:2010.00300*, 2020.
136. S. T. Radev, U. K. Mertens, A. Voss, and U. Köthe. “Towards end-to-end likelihood-free inference with convolutional neural networks”. *British Journal of Mathematical and Statistical Psychology* 73:1, 2020, pp. 23–43.
137. S. T. Radev, A. Voss, E. M. Wieschen, and P.-C. Bürkner. “Amortized Bayesian inference for models of cognition”. *International Conference on Cognitive Modelling (ICCM) Conference Proceedings*, 2020.
138. R. Ratcliff, A. Thapar, P. Gomez, and G. McKoon. “A diffusion model analysis of the effects of aging in the lexical-decision task.” *Psychology and aging* 19:2, 2004, p. 278.
139. O. Ratmann, C. Andrieu, C. Wiuf, and S. Richardson. “Model criticism based on likelihood-free inference, with an application to protein network evolution”. *Proceedings of the National Academy of Sciences* 106:26, 2009, pp. 10576–10581.
140. L. Raynal, J.-M. Marin, P. Pudlo, M. Ribatet, C. P. Robert, and A. Estoup. “ABC random forests for Bayesian parameter inference”. *Bioinformatics* 35:10, 2018, pp. 1720–1728.
141. C. Robert and G. Casella. “A short history of Markov chain Monte Carlo: Subjective recollections from incomplete data”. *Statistical Science*, 2011, pp. 102–115.

Bibliography

142. N. Said, M. Engelhart, C. Kirches, S. Körkel, and D. V. Holt. “Applying mathematical optimization methods to an ACT-R instance-based learning model”. *PloS one* 11:7, 2016, e0158832.
143. T. A. Salthouse. “Selective review of cognitive aging”. *Journal of the International Neuropsychological Society: JINS* 16:5, 2010, p. 754.
144. D. J. Schad, M. Betancourt, and S. Vasishth. “Toward a principled Bayesian workflow in cognitive science.” *Psychological Methods*, 2020.
145. M. Sensoy, L. Kaplan, and M. Kandemir. “Evidential deep learning to quantify classification uncertainty”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 3179–3189.
146. G. Shafer. *A mathematical theory of evidence*. Vol. 42. Princeton university press, 1976.
147. T. Shiono. “Estimation of agent-based models using Bayesian deep learning approach of BayesFlow”. Available at SSRN 3640351, 2020.
148. D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. “Mastering the game of Go with deep neural networks and tree search”. *nature* 529:7587, 2016, pp. 484–489.
149. S. Steegen, F. Tuerlinckx, A. Gelman, and W. Vanpaemel. “Increasing transparency through a multiverse analysis”. *Perspectives on Psychological Science* 11:5, 2016, pp. 702–712.
150. M. Sunnåker, A. G. Busetto, E. Numminen, J. Corander, M. Foll, and C. Dessimoz. “Approximate bayesian computation”. *PLoS computational biology* 9:1, 2013, e1002803.
151. I. Sutskever, O. Vinyals, and Q. V. Le. “Sequence to sequence learning with neural networks”. *arXiv preprint arXiv:1409.3215*, 2014.
152. C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. “Going deeper with convolutions”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
153. N. A. Taatgen, C. Lebiere, and J. R. Anderson. “Modeling paradigms in ACT-R”. *Cognition and multi-agent interaction: From cognitive modeling to social simulation*, 2006, pp. 29–52.
154. N. N. Taleb. *The black swan: The impact of the highly improbable*. Vol. 2. Random house, 2007.
155. S. Talts, M. Betancourt, D. Simpson, A. Vehtari, and A. Gelman. “Validating Bayesian inference algorithms with simulation-based calibration”. *arXiv preprint*, 2018.
156. V. Tchuiev and V. Indelman. “Inference over distribution of posterior class probabilities for reliable bayesian classification and object-level perception”. *IEEE Robotics and Automation Letters* 3:4, 2018, pp. 4329–4336.
157. T. Toni. “ABC SMC for parameter estimation and model selection with applications in systems biology”. *Nature Precedings*, 2011, pp. 1–1.

158. T. Toni, D. Welch, N. Strelkowa, A. Ipsen, and M. P. Stumpf. "Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems". *Journal of the Royal Society Interface* 6:31, 2009, pp. 187–202.
159. B. M. Turner, B. U. Forstmann, M. Steyvers, et al. *Joint models of neural and behavioral data*. Springer, 2019.
160. B. M. Turner and P. B. Sederberg. "A generalized, likelihood-free method for posterior estimation". *Psychonomic bulletin & review* 21:2, 2014, pp. 227–250.
161. B. M. Turner, P. B. Sederberg, and J. L. McClelland. "Bayesian analysis of simulation-based models". *Journal of Mathematical Psychology* 72, 2016, pp. 191–199.
162. M. Usher and J. L. McClelland. "The time course of perceptual choice: the leaky, competing accumulator model." *Psychological review* 108:3, 2001, p. 550.
163. S. Van Erp, J. Mulder, and D. L. Oberski. "Prior sensitivity analysis in default Bayesian structural equation modeling." *Psychological Methods* 23:2, 2018, p. 363.
164. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. "Attention is all you need". *arXiv preprint arXiv:1706.03762*, 2017.
165. A. Vehtari, A. Gelman, and J. Gabry. "Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC". *Statistics and computing* 27:5, 2017, pp. 1413–1432.
166. A. Vehtari, J. Ojanen, et al. "A survey of Bayesian predictive methods for model assessment, selection and comparison". *Statistics Surveys* 6, 2012, pp. 142–228.
167. R. Verity, L. C. Okell, I. Dorigatti, P. Winskill, C. Whittaker, N. Imai, G. Cuomo-Dannenburg, H. Thompson, P. G. Walker, H. Fu, et al. "Estimates of the severity of coronavirus disease 2019: a model-based analysis". *The Lancet Infectious Diseases*, 2020.
168. G. Viswanathan, E. Raposo, and M. Da Luz. "Lévy flights and superdiffusion in the context of biological encounters and random searches". *Physics of Life Reviews* 5:3, 2008, pp. 133–150.
169. A. Voss, V. Lerche, U. Mertens, and J. Voss. "Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models". *Psychonomic bulletin & review* 26:3, 2019, pp. 813–832.
170. J. Walker, C. Doersch, A. Gupta, and M. Hebert. "An uncertain future: Forecasting from static images using variational autoencoders". In: *European Conference on Computer Vision*. Springer, 2016, pp. 835–851.
171. S. Watanabe. "WAIC and WBIC are information criteria for singular statistical model evaluation". In: *Proceedings of the Workshop on Information Theoretic Methods in Science and Engineering*. 2013, pp. 90–94.
172. E. M. Wieschen, A. Voss, and S. Radev. "Jumping to conclusion? a lévy flight model of decision making". *TQMP* 16:2, 2020, pp. 120–132.
173. K. Xu, B. Nosek, and A. Greenwald. "Psychology data from the race implicit association test on the project implicit demo website". *Journal of Open Psychology Data* 2:1, 2014.

Bibliography

174. K. Xu, J. Li, M. Zhang, S. S. Du, K.-i. Kawarabayashi, and S. Jegelka. “What can neural networks reason about?” *arXiv preprint arXiv:1905.13211*, 2019.
175. Z. Yang, Z. Hu, R. Salakhutdinov, and T. Berg-Kirkpatrick. “Improved variational autoencoders for text modeling using dilated convolutions”. In: *International conference on machine learning*. PMLR, 2017, pp. 3881–3890.
176. Y. Yao, A. Vehtari, D. Simpson, A. Gelman, et al. “Using stacking to average Bayesian predictive distributions (with discussion)”. *Bayesian Analysis* 13:3, 2018, pp. 917–1007.
177. T. Yarkoni and J. Westfall. “Choosing prediction over explanation in psychology: Lessons from machine learning”. *Perspectives on Psychological Science* 12:6, 2017, pp. 1100–1122.
178. S. Zhao, J. Song, and S. Ermon. “Infovae: Information maximizing variational autoencoders”. *arXiv preprint arXiv:1706.02262*, 2017.
179. V. M. Zolotarev. *One-dimensional stable distributions*. Vol. 65. American Mathematical Soc., 1986.

APPENDIX A1 - MANUSCRIPT 1

Manuscript 1: BayesFlow: Learning complex stochastic models with invertible neural network.

BayesFlow: Learning Complex Stochastic Models with Invertible Neural Networks

Stefan T. Radev, Ulf K. Mertens, Andreas Voss, Lynton Ardizzone, and Ullrich Köthe

Abstract—Estimating the parameters of mathematical models is a common problem in almost all branches of science. However, this problem can prove notably difficult when processes and model descriptions become increasingly complex and an explicit likelihood function is not available. With this work, we propose a novel method for globally amortized Bayesian inference based on invertible neural networks which we call BayesFlow. The method uses simulation to learn a global estimator for the probabilistic mapping from observed data to underlying model parameters. A neural network pre-trained in this way can then, without additional training or optimization, infer full posteriors on arbitrary many real datasets involving the same model family. In addition, our method incorporates a summary network trained to embed the observed data into maximally informative summary statistics. Learning summary statistics from data makes the method applicable to modeling scenarios where standard inference techniques with hand-crafted summary statistics fail. We demonstrate the utility of BayesFlow on challenging intractable models from population dynamics, epidemiology, cognitive science and ecology. We argue that BayesFlow provides a general framework for building amortized Bayesian parameter estimation machines for any forward model from which data can be simulated.

I. INTRODUCTION

THE goal of Bayesian analysis is to infer the underlying characteristics of some natural process of interest given observable manifestations \mathbf{x} . In a Bayesian setting, we assume that we already possess sufficient understanding of the forward problem, that is, a suitable model of the mechanism that generates observations from a given configuration of the hidden parameters θ . This forward model can be provided in two forms: In likelihood-based approaches, the likelihood function $p(\mathbf{x}|\theta)$ is *explicitly known* and can be *evaluated* analytically or numerically for any pair (\mathbf{x}, θ) . In contrast, likelihood-free approaches only require the ability to *sample* from the likelihood. The latter approaches are typically realized by simulation programs, which generate synthetic observations by means of a deterministic function g of parameters θ and independent noise (i.e., random numbers) ξ :

$$\mathbf{x}_i \sim p(\mathbf{x}|\theta) \iff \mathbf{x}_i = g(\theta, \xi_i) \text{ with } \xi_i \sim p(\xi) \quad (1)$$

In this case, the likelihood $p(\mathbf{x}|\theta)$ is only defined *implicitly* via the action of the simulation program g , but calculation of its actual numerical value for a simulated observation \mathbf{x}_i is impossible. This, in turn, prohibits standard statistical inference.

Likelihood-free problems arise, for example, when $p(\mathbf{x}|\theta)$ is not available in closed-form, or when the forward model is defined by a stochastic differential equation, a Monte-Carlo simulation, or a complicated algorithm [27], [49], [47], [51]. In this paper, we propose a new Bayesian solution to the likelihood-free setting in terms of *invertible neural networks*.

Bayesian modeling leverages the available knowledge about the forward model to get the best possible estimate of the posterior distribution of the inverse model:

$$p(\theta|\mathbf{x}_{1:N}) = \frac{p(\mathbf{x}_{1:N}|\theta)p(\theta)}{\int p(\mathbf{x}_{1:N}|\theta)p(\theta)d\theta}$$

In Bayesian inference, the posterior encodes all information about θ obtainable from a set of observations $\mathbf{x}_{1:N} = \{\mathbf{x}_i\}_{i=1}^N$. The observations are assumed to arise from N runs of the forward model with fixed, but unknown, true parameters θ^* . Bayesian inverse modeling is challenging for three reasons:

- 1) The right-hand side of Bayes' formula above is always intractable in the likelihood-free case and must be approximated.
- 2) The forward model is usually non-deterministic, so that there is intrinsic uncertainty about the true value of θ .
- 3) The forward model is typically not information-preserving, so that there is ambiguity among possible values of θ .

The standard solution to these problems is offered by *approximate Bayesian computation* (ABC) methods [45], [10], [39], [47]. ABC methods approximate the posterior by repeatedly sampling parameters from a proposal (prior) distribution $\theta^{(l)} \sim p(\theta)$ and then simulating multiple datasets by running the forward model $\mathbf{x}_i \sim p(\mathbf{x}|\theta^{(l)})$ for $i = 1 \dots N$. If the resulting dataset is sufficiently similar to the actually observed dataset $\mathbf{x}_{1:N}^o$, the corresponding $\theta^{(l)}$ is retained as a sample from the desired posterior, otherwise rejected. Stricter similarity criteria lead to more accurate approximations of the desired posterior at the price of higher and oftentimes prohibitive rejection rates.

More efficient methods for approximate inference, such as sequential Monte Carlo (ABC-SMC), Markov-Chain Monte Carlo variants [44], or the recent neural density estimation methods [16], [38], [30], optimize sampling from a proposal distribution in order to balance the speed-accuracy trade-off of vanilla ABC methods. More details can be found in the section **Related Work** and in the excellent review by [9].

All sampling methods described above operate on the level of individual datasets, that is, for each observation sequence $\mathbf{x}_{1:N}$, the entire estimation procedure for the posterior must be run again from scratch. Therefore, we refer to this approach as *case-based inference*. Running estimation for each individual dataset separately stands in contrast to *amortized inference*, where estimation is split into a potentially expensive *upfront* training phase, followed by a much cheaper inference phase. The goal of the upfront training phase is to learn an approximate posterior $\hat{p}(\theta|\mathbf{x}_{1:N})$ that works well for *any* observation sequence

$\mathbf{x}_{1:N}$. Evaluating this model for specific observations $\mathbf{x}_{1:N}^o$ is then very fast, so that the training effort amortizes over repeated evaluations (see Figure 1 for a graphical illustration). The break-even between case-based and amortized inference depends on the application and model types, and we will report comparisons in the experimental section. Our main aim in this paper, however, is the introduction of a general approach to amortized Bayesian inference and the demonstration of its excellent accuracy in posterior estimation for a variety of popular forward models.

To make amortized inference feasible in practice, it must work well for arbitrary dataset sizes N . Depending on data acquisition circumstances, the number of available observations for a fixed model parameter setting may vary considerably, ranging from $N = 1$ to several hundreds and beyond. This has not only consequences for the required architecture of our density approximators, but also for their behavior: They must exhibit correct *posterior contraction*. Accordingly, the estimated posterior $\hat{p}(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ should get sharper (i.e., more peaked) as the number N of available observations increases. In the simplest case, the posterior variance should decrease at rate $1/N$, but more complex behavior can occur for difficult (e.g., multi-modal) true posteriors $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$.

We incorporate these considerations into our method by integrating two separate deep neural networks modules (detailed in the **Methods** section; see also Figure 1), which are trained jointly on simulated data from the forward model: a *summary network* and an *inference network*.

The *summary network* is responsible for reducing a set of observations $\mathbf{x}_{1:N}$ of variable size to a fixed-size vector of *learned* summary statistics. In traditional likelihood-free approaches, the method designer is responsible for selecting suitable statistics for each application *a priori* [33], [32], [43], [45]. In contrast, our summary networks learn the most informative statistics directly from data, and we will show experimentally (see **Experiment 3.7**) that these statistics are superior to manually constructed ones. Summary networks differ from standard feed-forward networks because they should be independent of the input size N and respect the inherent functional and probabilistic symmetries of the data. For example, permutation invariant networks are required for *i.i.d.* observations [6], and recurrent networks [15] or convolutional networks [29] for data with temporal or spatial dependencies.

The *inference network* is responsible for learning the true posterior of model parameters given the summary statistics of the observed data. Since it sees the data only through the lens of the summary network, all symmetries captured by the latter are automatically inherited by the posterior. We implement the inference network as an *invertible neural network*. Invertible neural networks are based on the recent theory and applications of *normalizing flows* [3], [25], [18], [13], [26]. Flow-based methods can perform exact inference under perfect convergence and scale favourably from simple low-dimensional problems to high-dimensional distributions with complex dependencies, for instance, the pixels of an image [25]. For each application/model of interest, we train an invertible network jointly with a corresponding summary network using simulated data from the respective known

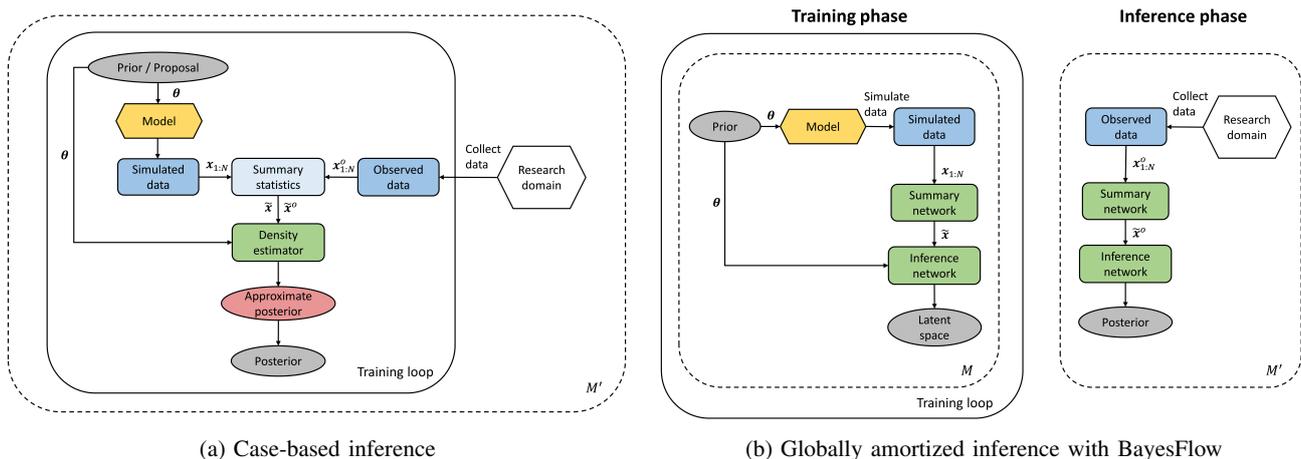
forward model with reasonable priors. After convergence of this forward training, the network’s invertibility ensures that a model for the inverse model is obtained for free, simply by running inference through the model backwards. Thus, our networks can perform fast amortized Bayesian inference on arbitrary many datasets from a given application domain without expensive case-based optimization. We call our method *BayesFlow*, as it combines ideas from Bayesian inference and flow-based deep learning.

BayesFlow draws on major advances in modern deep probabilistic modeling, also referred to as deep generative modeling [6], [25], [2], [24]. A hallmark idea in deep probabilistic modeling is to represent a complicated target distribution as a non-linear bijective transformation of some simpler latent distribution (e.g., Gaussian or uniform), a so called *pushforward*. Density estimation of the target distribution, a very complex problem, is thus reduced to learning a non-linear transformation, a task that is ideally suited for gradient-based neural network training via standard backpropagation. During the *inference phase*, samples from the target distribution are obtained by sampling from the simpler latent distribution and applying the inverse transformation learned during the training phase (see Figure 1b for a high-level overview). Using this approach, recent applications of deep probabilistic models have achieved unprecedented performance on hitherto intractable high-dimensional problems [6], [25], [18].

In the context of Bayesian inference, the target distribution is the posterior $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ of model parameters given observed data. We leverage the fact that we can simulate arbitrarily large amounts of training data from the forward model in order to ensure that the summary and invertible networks approximate the true posterior as well as possible. During the inference phase, our model can either numerically evaluate the posterior probability of any candidate parameter $\boldsymbol{\theta}$, or can generate a posterior sample $\boldsymbol{\theta}^{(1)}, \boldsymbol{\theta}^{(2)}, \dots, \boldsymbol{\theta}^{(L)}$ of likely parameters for the observed data $\mathbf{x}_{1:N}^o$. In the **Methods** section, we show that our networks indeed sample from the correct posterior under perfect convergence. In summary, the contributions of our BayesFlow method are the following:

- Globally amortized approximate Bayesian inference with invertible neural networks;
- Learning maximally informative summary statistics from raw datasets with variable number of observations instead of relying on restrictive hand-crafted summary statistics;
- Theoretical guarantee for sampling from the true posterior distribution with arbitrary priors and posteriors;
- Parallel computations applicable to both forward simulations and neural network optimization;

To illustrate the utility of BayesFlow, we first apply it to two toy models with analytically tractable posteriors. The first is a multivariate Gaussian with a full covariance matrix and a unimodal posterior. The second is a Gaussian mixture model with a multimodal posterior. Then, we present applications to challenging models with intractable likelihoods from population dynamics, cognitive science, epidemiology, and ecology and demonstrate the utility of BayesFlow in terms of speed, accuracy of recovery, and probabilistic calibration. Alongside,



(a) Case-based inference

(b) Globally amortized inference with BayesFlow

Fig. 1: Graphical illustration of the main differences between case-based (neural) density estimation methods and BayesFlow. (a) Case-based methods require a separate optimization loop for each observed dataset from a given research domain. When case-based methods incorporate a training phase (e.g., APT), it must be repeated for each new dataset. Summary statistics are manually selected and may thus be sub-optimal; (b) BayesFlow incorporates a global upfront training (before any real data are collected) via simulations from the forward model (left panel). Summary and inference network are trained jointly, resulting in higher accuracy than hand-crafted summary statistics. In the inference phase (right panel), BayesFlow works entirely in a feed-forward manner, that is, no training or optimization happens in this phase. The upfront training effort is therefore amortized over arbitrary many observed datasets from a research domain working on the same model family. Note that the solid and dashed plates are swapped between case-based Bayesian inference and the training phase of BayesFlow.

we introduce several performance validation tools.

A. Related Work

BayesFlow incorporates ideas from previous machine learning and deep learning approaches to likelihood-free inference [31], [41], [33], [43], [22]. The most common approach has been to cast the problem of parameter estimation as a supervised learning task. In this setting, a large dataset of the form $\mathbf{D} = \{(h(\mathbf{x}_{1:N}^{(m)}), \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ is created by repeatedly sampling from $p(\boldsymbol{\theta})$ and simulating an artificial datasets $\mathbf{x}_{1:N}$ by running $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ with the sampled parameters. Usually, the dimensionality of the simulated data is reduced by computing summary statistics with a fixed summary function $h(\mathbf{x}_{1:N})$. Then, a supervised learning algorithm (e.g., random forest [43], or a neural network [41]) is trained on the summary statistics of the simulated data to output an estimate of the true data generating parameters. Thus, an attempt is made to approximate the intractable inverse model $\boldsymbol{\theta} = g^{-1}(\mathbf{x}, \boldsymbol{\xi})$. A main shortcoming of supervised approaches is that they provide only limited information about the posterior (e.g., point-estimates, quantiles or variance estimates) or impose overly restrictive distributional assumptions on the shape of the posterior (e.g., Gaussian).

Our ideas are also closely related to the concept of *optimal transport maps* and its application in Bayesian inference [12], [40], [8], [5]. A transport map defines a transformation between (probability) measures which can be constructed in a way to *warp* a simple probability distribution into a more complex one. In the context of Bayesian inference, transport maps have been applied to accelerate MCMC sampling [40], to perform sequential inference [12], and to solve inference

problems via direct optimization [5]. In fact, BayesFlow can be viewed as a parameterization of invertible transport maps via invertible neural networks. An important distinction is that BayesFlow does not require an explicit likelihood function for approximating the target posteriors and is capable of amortized inference.

Similar ideas for likelihood-free inference are incorporated in the recent automatic posterior transformation (APT) [16], and the sequential neural likelihood (SNL) [38] methods. APT iteratively refines a proposal distribution via masked autoregressive flow (MAF) networks to generate parameter samples which closely match a particular observed dataset. SNL, in turn, trains a masked autoencoder density estimator (MADE) neural network within an MCMC loop to speed-up convergence to the true posterior. Even though these methods also entail a relatively expensive learning phase and a cheap inference phase, posterior inference is amortized only for a single dataset. Thus, the learning phase needs to be run through for each individual dataset (see Figure 1a). In contrast, we propose to learn the posterior globally over the entire range of plausible parameters and datasets by employing a conditional invertible neural network (cINN) estimator (see Figure 1b). Previously, INNs have been successfully employed to model data from astrophysics and medicine [2]. We adapt the model to suit the task of parameter estimation in the context of mathematical modeling and develop a probabilistic architecture for performing fully Bayesian and globally amortized inference with complex mathematical models.

II. METHODS

A. Notation

In the following, the number of parameters of a mathematical model will be denoted as D , and the number of observations in a dataset as N . We denote data simulated from the mathematical model of interest as $\mathbf{x}_{1:N} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$, where each individual \mathbf{x}_i can represent a scalar or a vector. Observed or test data will be marked with a superscript o (i.e., $\mathbf{x}_{1:N}^o$). The parameters of a mathematical model are represented as a vector $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_D)$, and all trainable parameters of the invertible and summary neural networks as ϕ and ψ , respectively. When a dataset consists of observations over a period of time, the number of observations will be denoted as T .

B. Learning the Posterior

Assume that we have an invertible function $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D$, parameterized by a vector of parameters ϕ , for which the inverse $f_\phi^{-1} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ exists. For now, consider the case when raw simulated data $\mathbf{x}_{1:N}$ of size $N = 1$ is entered directly into the invertible network without using a summary network. Our goal is to train an invertible neural network which approximates the true posterior as accurately as possible:

$$p_\phi(\boldsymbol{\theta} | \mathbf{x}) \approx p(\boldsymbol{\theta} | \mathbf{x}) \quad (2)$$

for all possible $\boldsymbol{\theta}$ and \mathbf{x} . We reparameterize the approximate posterior p_ϕ in terms of a conditional invertible neural network (cINN) f_ϕ which implements a normalizing flow between $\boldsymbol{\theta}$ and a Gaussian latent variable \mathbf{z} :

$$\boldsymbol{\theta} \sim p_\phi(\boldsymbol{\theta} | \mathbf{x}) \iff \boldsymbol{\theta} = f_\phi^{-1}(\mathbf{z}; \mathbf{x}) \text{ with } \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \quad (3)$$

Accordingly, we need to ensure that the outputs of $f_\phi^{-1}(\mathbf{z}; \mathbf{x})$ follow the target posterior $p(\boldsymbol{\theta} | \mathbf{x})$. Thus, we seek neural network parameters $\hat{\phi}$ which minimize the Kullback-Leibler (KL) divergence between the true and the model-induced posterior for all possible datasets \mathbf{x} . Therefore, our objective becomes:

$$\hat{\phi} = \operatorname{argmin}_{\phi} \mathbb{E}_{p(\mathbf{x})} [\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}) || p_\phi(\boldsymbol{\theta} | \mathbf{x}))] \quad (4)$$

$$= \operatorname{argmin}_{\phi} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [\log p(\boldsymbol{\theta} | \mathbf{x}) - \log p_\phi(\boldsymbol{\theta} | \mathbf{x})]] \quad (5)$$

$$= \operatorname{argmax}_{\phi} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{p(\boldsymbol{\theta} | \mathbf{x})} [\log p_\phi(\boldsymbol{\theta} | \mathbf{x})]] \quad (6)$$

$$= \operatorname{argmax}_{\phi} \iint p(\mathbf{x}, \boldsymbol{\theta}) \log p_\phi(\boldsymbol{\theta} | \mathbf{x}) d\mathbf{x} d\boldsymbol{\theta} \quad (7)$$

Note, that the log posterior density $p(\boldsymbol{\theta} | \mathbf{x})$ can be dropped from the optimization objective in Eq.6, as it does not depend on the neural network parameters ϕ . In other words, we seek neural network parameters $\hat{\phi}$ which maximize the posterior probability of data-generating parameters $\boldsymbol{\theta}$ given observed data \mathbf{x} for all $\boldsymbol{\theta}$ and \mathbf{x} . Since $f_\phi(\boldsymbol{\theta}; \mathbf{x}) = \mathbf{z}$ by design, the change of variable rule of probability yields:

$$p_\phi(\boldsymbol{\theta} | \mathbf{x}) = p(\mathbf{z} = f_\phi(\boldsymbol{\theta}; \mathbf{x})) \left| \det \left(\frac{\partial f_\phi(\boldsymbol{\theta}; \mathbf{x})}{\partial \boldsymbol{\theta}} \right) \right| \quad (8)$$

Thus, we can re-write our objective as:

$$\hat{\phi} = \operatorname{argmax}_{\phi} \iint p(\mathbf{x}, \boldsymbol{\theta}) \log p_\phi(\boldsymbol{\theta} | \mathbf{x}) d\mathbf{x} d\boldsymbol{\theta} \quad (9)$$

$$= \operatorname{argmax}_{\phi} \iint p(\mathbf{x}, \boldsymbol{\theta}) (\log p(f_\phi(\boldsymbol{\theta}; \mathbf{x})) + \log |\det \mathbf{J}_{f_\phi}|) d\mathbf{x} d\boldsymbol{\theta} \quad (10)$$

where we have abbreviated $\partial f_\phi(\boldsymbol{\theta}; \mathbf{x}) / \partial \boldsymbol{\theta}$ (the Jacobian of f_ϕ evaluated at $\boldsymbol{\theta}$ and \mathbf{x}) as \mathbf{J}_{f_ϕ} . Due to the architecture of our cINN, the $\log |\det \mathbf{J}_{f_\phi}|$ is easy to compute (see next section for details).

Utilizing simulations from the forward model (Eq.1), we can approximate the expectations by minimizing the Monte-Carlo estimate of the negative of Eq.10. Accordingly, for a batch of M simulated datasets and data-generating parameters $\{(\mathbf{x}^{(m)}, \boldsymbol{\theta}^{(m)})\}_{m=1}^M$ we have:

$$\hat{\phi} = \operatorname{argmin}_{\phi} \frac{1}{M} \sum_{m=1}^M -\log p_\phi(\boldsymbol{\theta}^{(m)} | \mathbf{x}^{(m)}) \quad (11)$$

$$= \operatorname{argmin}_{\phi} \frac{1}{M} \sum_{m=1}^M \left(-\log p(f_\phi(\boldsymbol{\theta}^{(m)}; \mathbf{x}^{(m)})) - \log |\det \mathbf{J}_{f_\phi}^{(m)}| \right) \quad (12)$$

$$= \operatorname{argmin}_{\phi} \frac{1}{M} \sum_{m=1}^M \left(\frac{\|f_\phi(\boldsymbol{\theta}^{(m)}; \mathbf{x}^{(m)})\|_2^2}{2} - \log |\det \mathbf{J}_{f_\phi}^{(m)}| \right) \quad (13)$$

We treat Eq.13 as a loss function $\mathcal{L}(\phi)$ which can be minimized with any stochastic gradient descent method. The first term follows from Eq.12 due to the fact that we have prescribed a unit Gaussian distribution to \mathbf{z} . It represents the negative log of $\mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \propto \exp(-\frac{1}{2}\|\mathbf{z}\|_2^2)$. The second term controls the rate of volume change induced by the learned non-linear transformation from $\boldsymbol{\theta}$ to \mathbf{z} achieved by f_ϕ . Thus, minimizing Eq.13 ensures that \mathbf{z} follows the prescribed unit Gaussian.

The correctness of the learned posterior can be guaranteed in the following way, assuming the network is able to reach the global minimum of the loss (i.e. under perfect convergence).

Proposition 1. Assume that the cINN architecture and domain of ϕ are chosen such that $\hat{\phi}$ is the global minimum of the objective in Eq.10. Then, the latent output distribution will be statistically independent of the conditioning data, $p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) \perp p(\mathbf{x})$. As a result, the samples transformed backwards from $p(\mathbf{z})$ will follow the true posterior, that is:

$$f_{\hat{\phi}}^{-1}(\mathbf{z}; \mathbf{x}) \sim p(\boldsymbol{\theta} | \mathbf{x}) \text{ with } \mathbf{z} \sim \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I}) \quad (14)$$

Proof. For short, we denote $p(\mathbf{z}) := \mathcal{N}_D(\mathbf{z} | \mathbf{0}, \mathbb{I})$, and the distribution of network outputs as $p(f_{\hat{\phi}}(\boldsymbol{\theta}; \mathbf{x})) := p_{\hat{\phi}}(\mathbf{z} | \mathbf{x})$. Due to $\mathbb{KL}(\cdot || \cdot) \geq 0$ (Gibbs' inequality), the global minimum of the objective is achieved exactly when the argument in Eq.4 becomes 0. To relate this to the sampling process, we note the invariance of \mathbb{KL} under diffeomorphic transformations, from which it follows that

$$\mathbb{KL}(p_{\hat{\phi}}(\mathbf{z} | \mathbf{x}) || p(\mathbf{z})) = 0. \quad (15)$$

Considering $p(z) \perp p(\mathbf{x})$ and $p_{\hat{\phi}}(z | \mathbf{x}) = p(z)$ (from Eq.15), this also implies $p_{\hat{\phi}}(z | \mathbf{x}) \perp p(\mathbf{x})$, which means the latent output distribution is the same for any fixed \mathbf{x} we choose. This motivates the validity of taking samples from $p(z)$ and transforming them back using the condition, to generate samples from the posterior. By definition of the model, the generated samples $f_{\hat{\phi}}^{-1}(z, \mathbf{x})$ with $z \sim p(z)$ follow $p_{\hat{\phi}}(\boldsymbol{\theta} | \mathbf{x})$. The proposition therefore holds when the argument in Eq.4 is zero. \square

We now generalize our formulation to datasets with arbitrary numbers of observations. If we let the number of observations N vary and train a summary network $\tilde{\mathbf{x}} = h_{\psi}(\mathbf{x}_{1:N})$ together with the cINN, our main objective changes to:

$$\hat{\phi}, \hat{\psi} = \operatorname{argmax}_{\phi, \psi} \mathbb{E}_{p(\mathbf{x}, \boldsymbol{\theta}, N)} [\log p_{\phi}(\boldsymbol{\theta} | h_{\psi}(\mathbf{x}_{1:N}))] \quad (16)$$

and its Monte Carlo estimate to:

$$\hat{\phi}, \hat{\psi} = \operatorname{argmin}_{\phi, \psi} \frac{1}{M} \sum_{m=1}^M \left(\frac{\|f_{\phi}(\boldsymbol{\theta}^{(m)}; h_{\psi}(\mathbf{x}_{1:N}^{(m)}))\|_2^2}{2} - \log \left| \det \left(\mathbf{J}_{f_{\phi}}^{(m)} \right) \right| \right) \quad (17)$$

In order to make the estimation of $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$ tractable, we assume that there exists a vector $\boldsymbol{\eta}$ of sufficient statistics that captures all information about $\boldsymbol{\theta}$ contained in $\mathbf{x}_{1:N}$ in a fixed-size (vector) representation. For $h_{\psi}(\mathbf{x}_{1:N})$ to be a useful estimator for $\boldsymbol{\eta}$, both should convey the same information about $\boldsymbol{\theta}$, as measured by the mutual information:

$$MI(\boldsymbol{\theta}, h_{\psi}(\mathbf{x}_{1:N})) \approx MI(\boldsymbol{\theta}, \boldsymbol{\eta}) \quad (18)$$

Since we do not know $\boldsymbol{\eta}$, we can enforce this requirement only indirectly by minimizing the Monte Carlo estimate of Eq.16. The following proposition states that, under perfect convergence, samples from a cINN still follow the true posterior given the outputs of a summary networks.

Proposition 2. *Assume that we have a perfectly converged cINN f_{ϕ} and a perfectly converged summary network h_{ψ} . Assume also, that there exists a vector $\boldsymbol{\eta}$ of sufficient summary statistics for $\mathbf{x}_{1:N}$. Then, independently sampling $z \sim p(z)$ and applying $f_{\phi}^{-1}(z; h_{\psi}(\mathbf{x}_{1:N}))$ to each z yields independent samples from $p(\boldsymbol{\theta} | \mathbf{x}_{1:N})$.*

Proof. Perfect convergence of the networks under Eq.16 implies $\mathbb{KL}(p(\boldsymbol{\theta} | \mathbf{x}_{1:N}) || p_{\phi}(\boldsymbol{\theta} | h_{\psi}(\mathbf{x}_{1:N}))) = 0$. This, in turn, implies that $MI(\boldsymbol{\theta}, h_{\psi}(\mathbf{x}_{1:N})) = MI(\boldsymbol{\theta}, \boldsymbol{\eta})$, because a perfect match of the densities would be impossible if $h_{\psi}(\mathbf{x}_{1:N})$ contained less information about $\boldsymbol{\theta}$ than $\boldsymbol{\eta}$. Therefore, the proof reduces to that of **Proposition 1**. Note, that whenever the KL divergence is driven to a minimum, $h_{\psi}(\mathbf{x}_{1:N})$ is a *maximally informative statistic* [11]. \square

In summary, the approximate posteriors obtained by the BayesFlow method are correct if the summary and invertible networks are perfectly converged. In practice, however, perfect convergence is unrealistic and there are three sources of error which can lead to incorrect posteriors. The first is the Monte Carlo error introduced by using simulations from $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ to

approximate the expectation in Eq.16. The second is due to a summary network which may not fully capture the relevant information in the data or when sufficient summary statistics do not exist. The third is due to an invertible network which does not accurately transform the true posterior into the prescribed Gaussian latent space. Even though we can mitigate the Monte Carlo error by running the simulator $g(\boldsymbol{\theta}, \boldsymbol{\xi})$ more often, the latter two can be harder to detect and alleviate in a principled way. Nevertheless, recent work on *probabilistic symmetry* [6] and *algorithmic alignment* [52] can provide some guidelines on how to choose the right summary network for a particular problem. Additionally, the depth as well as the building blocks (to be explained shortly) of the invertible chain can be tuned to increase the expressiveness of the learned transformation from $\boldsymbol{\theta}$ -space to z -space. The benefits of neural network depth have been confirmed both in theory and practice [28], [4], so we expect better performance in complex settings with increasing network depth.

C. Composing Invertible Networks

The basic building block of our cINN is the affine coupling block (ACB) [13]. Each ACB consists of four separate fully connected neural networks denoted as $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$. An ACB performs an invertible non-linear transformation, which means that in addition to a parametric mapping $f_{\phi} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ it also learns the inverse mapping $f_{\phi}^{-1} : \mathbb{R}^D \rightarrow \mathbb{R}^D$ for free. Denoting the input vector of f_{ϕ} as \mathbf{u} and the output vector as \mathbf{v} , it follows that $f_{\phi}(\mathbf{u}) = \mathbf{v}$ and $f_{\phi}^{-1}(\mathbf{v}) = \mathbf{u}$. Invertibility is achieved by splitting the input vector into two parts $\mathbf{u} = (\mathbf{u}_1, \mathbf{u}_2)$ with $\mathbf{u}_1 = u_{1:D/2}$ and $\mathbf{u}_2 = u_{D/2+1:D}$ (where $D/2$ is understood as a floor division) and performing the following operations on the split input:

$$\mathbf{v}_1 = \mathbf{u}_1 \odot \exp(s_1(\mathbf{u}_2)) + t_1(\mathbf{u}_2) \quad (19)$$

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_2(\mathbf{v}_1)) + t_2(\mathbf{v}_1) \quad (20)$$

where \odot denotes element-wise multiplication. The outputs $\mathbf{v} = (\mathbf{v}_1, \mathbf{v}_2)$ are then concatenated again and passed to the next ACB. The inverse operation is given by:

$$\mathbf{u}_2 = (\mathbf{v}_2 - t_2(\mathbf{v}_1)) \odot \exp(-s_2(\mathbf{v}_1)) \quad (21)$$

$$\mathbf{u}_1 = (\mathbf{v}_1 - t_1(\mathbf{u}_2)) \odot \exp(-s_1(\mathbf{u}_2)) \quad (22)$$

This formulation ensures that the Jacobian of the affine transformation is a strictly upper or a lower triangular matrix and therefore its determinant is very cheap to compute. Furthermore, the internal functions $s_1(\cdot), s_2(\cdot), t_1(\cdot), t_2(\cdot)$ can be represented by arbitrarily complex neural networks, which themselves need not be invertible, since they are only ever evaluated in the forward direction during both the forward and the inverse pass through the ACBs. In our applications, we parameterize the internal functions as fully connected neural networks with exponential linear units (ELU).

In order to ensure that the neural network architecture is expressive enough to represent complex distributions, we chain multiple ACBs, so that the output of each ACB becomes the input to the next one. In this way, the whole chain remains invertible from the first input to the last output and can

be viewed as a single function parameterized by trainable parameters ϕ .

In our applications, the input to the first ACB is the parameter vector θ , and the output of the final ACB is a d -dimensional vector z representing the non-linear transformation of the parameters. As described in the previous section, we ensure that z follows a unit Gaussian distribution via optimization, that is, $p(z) = \mathcal{N}_D(z | \mathbf{0}, \mathbb{I})$. Fixed permutation matrices are used before each ACB to ensure that each axis of the transformed parameter space z encodes information from all components of θ .

In order to account for the observed data, we feed the learned summary vectors into all internal networks of each ACB (explained shortly). Intuitively, in this way we realize the following process: the forward pass maps data-generating parameters θ to z -space using conditional information from the data $\mathbf{x}_{1:N}$, while the inverse pass maps data points from z -space to the data-generating parameters of interest using the same conditional information.

D. Summary Network

Since the number of observations usually varies in practical scenarios (e.g., different number of measurements or time points) and since datasets might exhibit various redundancies, the cINN can profit from some form of dimensionality reduction. As previously mentioned, we want to avoid information loss through restrictive hand-crafted summary statistics and, instead, learn the most informative summary statistics directly from data. Therefore, instead of feeding the raw simulated or observed data to each ACB, we pass the data through an additional summary network to obtain a fixed-sized vector of learned summary statistics $\tilde{\mathbf{x}} = h_{\psi}(\mathbf{x}_{1:N})$.

The architecture of the summary network should be aligned with the probabilistic symmetry of the observed data. An obvious choice for time series-data is an LSTM-network [15], since recurrent networks can naturally deal with long sequences of variable size. Another choice might be a 1D fully convolutional network [29], which has already been applied in the context of likelihood-free inference [41]. A different architecture is needed when dealing with *i.i.d.* samples of variable size. Such data are often referred to as *exchangeable*, or *permutation invariant*, since changing the order of individual elements does not change the associated likelihood or posterior. In other words, if $\mathbb{S}_N(\cdot)$ is an arbitrary permutation of N elements, the following should hold for the posterior:

$$p(\theta | \mathbf{x}_{1:N}) = p(\theta | \mathbb{S}_N(\mathbf{x}_{1:N})) \quad (23)$$

Following [6], we encode probabilistic permutation invariance by implementing a permutation invariant function through an equivariant non-linear transformation followed by a pooling operator (e.g., sum or mean) and another non-linear transformation:

$$\tilde{\mathbf{x}} = h_{\psi_1} \left(\sum_{i=1}^N h_{\psi_2}(\mathbf{x}_i) \right) \quad (24)$$

where h_{ψ_1} and h_{ψ_2} are two different fully connected neural networks. In practice, we stack multiple equivariant and

invariant functions into an invariant network in order to achieve higher expressiveness [6].

We optimize the parameters ψ of the summary network jointly with those of the cINN chain via backpropagation. Thus, training remains completely end-to-end, and BayesFlow learns to generalize to datasets of different sizes by suitably varying N during training of a permutation invariant summary network or varying sequence length during training of a recurrent/convolutional network.

To incorporate the observed or simulated data $\mathbf{x}_{1:N}$, each of the internal networks of each ACB is augmented to take the learned summary vector $\tilde{\mathbf{x}}$ of the data as an additional input. The output of each ACB then becomes:

$$\mathbf{v}_1 = \mathbf{u}_1 \odot \exp(s_1(\mathbf{u}_2, \tilde{\mathbf{x}})) + t_1(\mathbf{u}_2, \tilde{\mathbf{x}}) \quad (25)$$

$$\mathbf{v}_2 = \mathbf{u}_2 \odot \exp(s_2(\mathbf{v}_1, \tilde{\mathbf{x}})) + t_2(\mathbf{v}_1, \tilde{\mathbf{x}}) \quad (26)$$

Thus, a complete pass through the entire conditional invertible chain can be expressed as $f_{\phi}(\theta; \tilde{\mathbf{x}}) = z$ together with the inverse operation $f_{\phi}^{-1}(z; \tilde{\mathbf{x}}) = \theta$. The inverse transformation during inference is depicted in Figure 2.

E. Putting It All Together

Algorithm 1 describes the essential steps of the BayesFlow method using an arbitrary summary network and employing an online learning approach.

The backpropagation algorithm works by computing the gradients of the loss function with respect to the parameters of the neural networks and then adjusting the parameters, so as to drive the loss function to a minimum. We experienced no instability or convergence issues during training with the loss function given by Eq.16. Note, that steps 2-13 and 16-20 of **Algorithm 1** can be executed in parallel with GPU support in order to dramatically accelerate convergence and inference. Moreover, steps 16-20 can be applied in parallel to an arbitrary number of observed datasets after convergence of the networks (see Figure 2 for a full graphical illustration).

In what follows, we apply BayesFlow to two toy models with a unimodal and multimodal posteriors, respectively, and then use it to perform Bayesian inference on challenging models from population dynamics, cognitive science, epidemiology, and ecology.¹ We deem these models suitable for an initial validation, since they differ widely in the generative mechanisms they implement and the observed data they aim to explain. Therefore, good performance on these disparate examples underlines the broad empirical utility of the BayesFlow method. Details for models' setup can be found in **Appendix C**.

III. EXPERIMENTS

A. Training the Networks

We train all invertible and summary networks described in this paper jointly via backpropagation. For all following experiments, we use the Adam optimizer with a starter learning rate of 10^{-3} and an exponential decay rate of .95. We perform 50 000 to 100 000 iterations (i.e., mini-batch update steps)

¹Code and simulation scripts for all current applications are available at <https://github.com/stefanradev93/cINN>.

Algorithm 1 Amortized Bayesian inference with the BayesFlow method

```

1: Training phase (online learning with batch size M):
2: repeat
3:   Sample number of observations  $N \sim \mathcal{U}(N_{min}, N_{max})$ .
4:   for  $m = 1, \dots, M$  do
5:     Sample model parameters from prior:  $\theta^{(m)} \sim p(\theta)$ .
6:     for  $i = 1, \dots, N$  do
7:       Sample a noise instance:  $\xi_i \sim p(\xi)$ .
8:       Run the simulation (cf. Eq.1) to create a synthetic observation:  $x_i^{(m)} = g(\theta^{(m)}, \xi_i)$ .
9:       Pass the dataset  $\mathbf{x}_{1:N}^{(m)}$  through the summary network:  $\tilde{\mathbf{x}}^{(m)} = h_\psi(\mathbf{x}_{1:N}^{(m)})$ .
10:      Pass  $(\theta^{(m)}, \tilde{\mathbf{x}}^{(m)})$  through the inference network in forward direction:  $\mathbf{z}^{(m)} = f_\phi(\theta^{(m)}; \tilde{\mathbf{x}}^{(m)})$ .
11:      Compute loss according to Eq.17 from the training batch  $\{(\theta^{(m)}, \tilde{\mathbf{x}}^{(m)}, \mathbf{z}^{(m)})\}_{m=1}^M$ .
12:      Update neural network parameters  $\phi, \psi$  via backpropagation.
13: until convergence to  $\hat{\phi}, \hat{\psi}$ 
14:
15: Inference phase (given observed or test data  $\mathbf{x}_{1:N}^o$ ):
16: Pass the observed dataset through the summary network:  $\tilde{\mathbf{x}}^o = h_{\hat{\psi}}(\mathbf{x}_{1:N}^o)$ .
17: for  $l = 1, \dots, L$  do
18:   Sample a latent variable instance:  $\mathbf{z}^{(l)} \sim \mathcal{N}_D(\mathbf{0}, \mathbb{I})$ .
19:   Pass  $(\tilde{\mathbf{x}}^o, \mathbf{z}^{(l)})$  through the inference network in inverse direction:  $\theta^{(l)} = f_{\hat{\phi}}^{-1}(\mathbf{z}^{(l)}; \tilde{\mathbf{x}}^o)$ .
20: Return  $\{\theta^{(l)}\}_{l=1}^L$  as a sample from  $p(\theta | \mathbf{x}_{1:N}^o)$ 

```

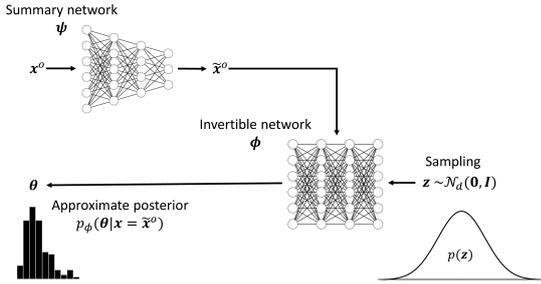


Fig. 2: Inference with pre-trained summary and inference networks. The posterior is approximated given real observed data via independent samples from a learned pushforward distribution. Thus, knowledge about the mapping between data and parameters (the inverse model) is compactly encoded within the weights of the two networks.

for each experiment, and report the results obtained by the converged networks. Note, that we did not perform an extensive search for optimal values of network hyperparameters, but use a default BayesFlow with 5 to 10 ACBs and a summary vector of size 128 for all examples in this paper (see **Appendix C** for more details on summary network architectures). All networks were implemented in Python using the *TensorFlow* library [1] and trained on a single-GPU machine equipped with NVIDIA[®] GTX1060 graphics card.

Regarding the data generation step, we take an approach which incorporates ideas from *online learning* [36] where data are simulated by Eq.1 on demand. Correspondingly, a dataset $\mathbf{x}_{1:N}$, or a batch of M datasets $\{\mathbf{x}_{1:N}^{(i)}\}_{i=1}^M$, is generated

on the fly and then passed through the neural network. This training approach has the advantage that the network never *experiences* the same input data twice. Moreover, training can continue as long as the network keeps improving (i.e., the loss keeps decreasing), since overfitting in the classical sense is nearly impossible. However, if simulations are computationally expensive and researchers need to experiment with different networks or training hyperparameters, it might be beneficial to store and re-use simulations, since simulation and training in online learning are tightly intertwined.

Once the networks have converged, we store the trained networks and use them to perform amortized inference on a separate validation set of datasets. The pre-trained networks can also be shared among a research community so that multiple researchers/labs can benefit from the amortization of inference.

B. Performance Validation

To evaluate the performance of BayesFlow in the following application examples, we consider a number of different metrics:

- Normalized root mean squared error (NRMSE) - to assess accuracy of point-estimates in recovering ground-truth parameter values;
- Coefficient of determination (R^2) - to assess the proportion of variance in ground-truth parameters that is captured by the point estimates;
- Re-simulation error (Err_{sim}) - to assess the predictive mismatch between the true data distribution and the data distribution generated with the estimated parameters (i.e., posterior predictive check);
- Calibration error (Err_{cal} , [21]) - to assess the coverage of the approximate posteriors (i.e., whether credibility intervals are indeed credible);

- Simulation-based calibration (SBC, [46]) - to visually detect systematic biases in the approximate posteriors; Details for computing all metrics are given in **Appendix B**.

C. Proof of Concept: Multivariate Normal Distribution

As a proof-of-concept, we apply the BayesFlow method to recover the posterior mean vector of a toy multivariate normal (MVN) example. For a single D -dimensional MVN vector, the forward model is given by:

$$\boldsymbol{\mu}^{(m)} \sim \mathcal{N}_D(\boldsymbol{\mu} | \mathbf{0}, \mathbb{I}) \quad (27)$$

$$\mathbf{x}^{(m)} \sim \mathcal{N}_D(\mathbf{x} | \boldsymbol{\mu}^{(m)}, \boldsymbol{\Sigma}) \quad (28)$$

where in this illustrative case we assume a single D -dimensional sample per observation ($N = 1$). If the covariance matrix $\boldsymbol{\Sigma}$ is known, the posterior of the mean vector $\boldsymbol{\mu}$ has a closed-form which is also a MVN $p(\boldsymbol{\mu} | \mathbf{x}, \boldsymbol{\Sigma}) = \mathcal{N}_d(\boldsymbol{\mu} | \mathbf{m}, \boldsymbol{\Lambda})$ with posterior precision matrix given by $\boldsymbol{\Lambda}^{-1} = \mathbb{I} + \boldsymbol{\Sigma}^{-1}$ and posterior mean given by $\mathbf{m} = \boldsymbol{\Lambda} \boldsymbol{\Sigma}^{-1} \mathbf{x}$ [7]. We can thus generate multiple batches of the form $\{(\mathbf{x}^{(m)}, \boldsymbol{\mu}^{(m)})\}_{m=1}^M$ and pass them directly through an invertible network. Since the ground-truth posterior is Gaussian, we can compute the KL divergence as a measure of mismatch between the true and approximate posteriors in closed form.

We run three experiments with $D \in \{5, 50, 500\}$ where the size of the ACB blocks was doubled for each successive D . To assess results, we compute the R^2 and NRMSE between approximate and true means as well as the KL divergence between approximate and true distributions on 100 test datasets. To compute the approximate covariance matrix, we draw 5000 samples from the approximate posteriors for $D = 5$ and $D = 50$ and 50000 samples for $D = 500$.

The KL divergence for the 5-D and 50-D MVNs reached essentially 0 after 2-3 epochs of 1000 iterations indicating that this is an easy problem for BayesFlow, and almost perfect recovery of the true posteriors is possible. The KL divergence for the 500-D MVN model reached 0.37 after 50 epochs, which represents a negligible increase in entropy relative to the true posterior (0.05% nats) and indicates decent approximation in light of the high dimensionality of the problem.

D. Multimodal Posterior - Gaussian Mixture Model

In order to test whether the BayesFlow method can recover multimodal posteriors, we apply it to a generative Gaussian mixture model (GMM). Multimodal posteriors arise in practice, for example, when forward models are defined as mixtures between different processes, or when models exhibit large multivariate trade-offs in their parameter space (e.g., there are multiple separate regions of posterior density with plausible parameter values). Therefore, it is important to show that our method is able to capture such behavior and does not suffer from mode collapse.

Following [2], we construct a scenario in which the observed data \mathbf{x} is a one-hot encoded vector representing one of the *hard* labels *red*, *green*, *blue*, or *yellow* (i.e., a single observation, thus $N = 1$). The parameters $\boldsymbol{\theta} = (\theta_1, \theta_2)$ are the 2D coordinates ²

²Note that this is not the typical GMM setup, as we construct the example such that the mixture assignments (labels) are observed and the data coordinates are the latent parameters.

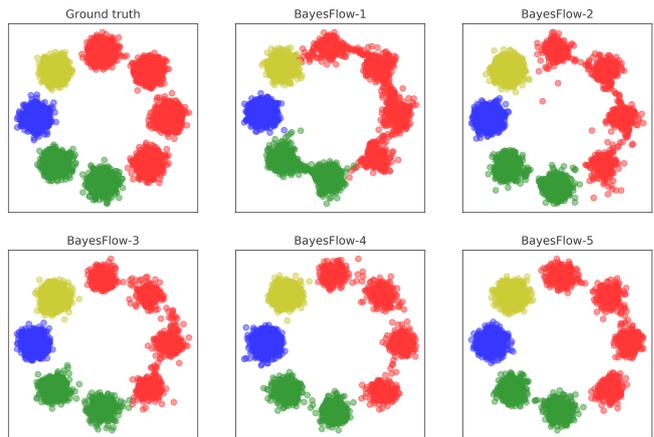


Fig. 3: Results on the GMM toy example with colors indicating cluster assignments. Approximation of the multimodal posterior become closer to the ground truth distribution with increasing depth (number of ACBs) of the conditional invertible network.

of points drawn from a mixture of eight Gaussian clusters with centers distributed around the origin in a clockwise manner and unit variance (see Figure 3, upper left). The first four clusters are assigned the label *red*, the next two the label *green*, and the remaining two the labels *blue* and *yellow*. The posterior $p(\boldsymbol{\theta} | \mathbf{x})$ is composed of the clusters indexed by the corresponding label. We perform the experiment multiple times by increasing the depth of the BayesFlow starting from 1 ACB block up to 5 ACB blocks. In this way, we can investigate the effects of cINN depth on the quality of the approximate multimodal posteriors. We train each BayesFlow for 50 epochs and draw 8000 samples from the approximate posteriors obtained by the trained models.

Results for all BayesFlows are depicted in Figure 3. We observe that approximations profit from having a deeper cINN chain, with cluster separation becoming clearer when using more ACBs. This confirms that our method is capable of recovering multimodal posteriors.

E. Stochastic Time-Series Model - The Ricker Model

In the following, we estimate the parameters of a well-known discrete stochastic population dynamics model [51]. With this example, we are pursuing several goals: First, we want to demonstrate that the BayesFlow method is able to accurately recover the parameters of an actual model with intractable likelihood by learning summary statistics from raw data. Second, we show that BayesFlow can deal adequately with parameters that are completely unrelated to the data by reducing estimates to the corresponding parameters' prior. Third, we compare the global performance of the BayesFlow method to that of related methods capable of amortized likelihood-free inference. Finally, we demonstrate the desired posterior contraction and improvement in estimation with increasing number of observations.

Discrete population dynamics models describe how the number of individuals in a population changes over discrete units of time [51]. In particular, the Ricker model describes

the number of individuals x_t in generation t as a function of the expected number of individuals in the previous generation by the following non-linear equations:

$$x_t \sim \text{Pois}(\rho N_t) \quad (29)$$

$$\xi_t \sim \mathcal{N}(0, \sigma^2) \quad (30)$$

$$N_{t+1} = r N_t e^{-N_t + \xi_t} \quad (31)$$

for $t = 1, \dots, T$ where N_t is the expected number of individuals at time t , r is the growth rate, ρ is a scaling parameter and ξ_t is random Gaussian noise. The likelihood function for the Ricker model is not available in closed form, and the model is known to exhibit chaotic behavior [33]. Thus, it is a suitable candidate for likelihood-free inference. The parameter estimation task consists of recovering $\theta = (\rho, r, \sigma)$ from the observed one-dimensional time-series data $x_{1:T}$ where each $x_t \in \mathbb{N}$.

What if the data does not contain any information about a particular parameter? In this case, any good estimation method should detect this, and return the prior of the particular parameter. To test this, we append a random uniform variable $u \sim \mathcal{U}(0, 1)$ to the parameter vector θ and train BayesFlow with this additional dummy parameter. We expect that the networks ignore this dummy parameter, that is, we assume that the estimated posterior of u resembles the uniform prior.

We compare the performance of BayesFlow to the following recent methods capable of amortized likelihood-free inference: conditional variational autoencoder (cVAE) [35], cVAE with autoregressive flow (cVAE-IAF) [26], *DeepInference* with heteroscedastic loss [41], approximate Bayesian computation with an LSTM neural network for learning informative summary statistics (ABC-NN) [22] and quantile random forest (ABC-RF) [43]. For training the models, we simulate time-series from the Ricker model with varying lengths. The number of time points T is drawn from a uniform distribution $T \sim \mathcal{U}(100, 500)$ at each training iteration.

All neural network methods were trained for 100 epochs with 1000 iterations each on simulated data from the Ricker model. The ABC-RF method was fitted on a reference table with 200 000 datasets, since the method does not allow for online learning and increasing the reference table did not seem to improve performance. In order to avoid using hand-crafted summary statistics for the ABC-RF method, we input summary vectors obtained by applying the summary network trained jointly with the cINN. Thus, the ABC-RF method has the advantage of using maximally informative statistics as input. We validate the performance of all methods on an independent test set of 500 datasets generated with $T = 500$. We report performance metrics for each method and each parameter in Table I.

Parameters r and ρ seem to be well recoverable by all methods considered here. The σ parameter turns out to be harder to estimate, with BayesFlow and the ABC-NN method performing best. Further, BayesFlow performs very well across all parameters and metrics. Importantly, the calibration error Err_{cal} obtained by BayesFlow is always low, indicating that the shape of the approximate posterior closely matches that of the true posteriors. Variational methods (cVAE, cVAE-IAF) experience some problems recovering the posterior of σ .

The ABC-NN and ABC-RF methods seem to recover point estimates with high accuracy but the approximate posteriors of the former exhibit relatively high calibration error. The ABC-RF method can only estimate posterior quantiles, so no comparable calibration metric could be computed.

Further results are depicted in Figure 4. Inspecting the full posteriors obtained by all methods on an example test dataset, we note that only BayesFlow and the ABC-NN methods are able to recover the uninformative posterior distribution of the dummy noise variable u (Figure 4a). Moreover, the importance of a Bayesian treatment of the Ricker model becomes clear when looking at the posteriors of σ . On most test datasets, the posterior density spreads over the entire prior range (high posterior variance) indicating large uncertainty in the obtained estimates. Moreover, the shapes of the marginal parameter posteriors vary widely across validation datasets, which highlights the importance of avoiding *ad hoc* restrictions on allowed posterior shapes (see Figure S5 for examples). We also observe that parameter estimation with BayesFlow becomes increasingly accurate when more time points are available (Figure 4b). Parameter recovery is especially good with the maximum number of time points (see Figure 4c). Finally, (Figure 4d) reveals a notable posterior contraction across increasing number of time points available to the summary network.

F. A Model of Perceptual Decision Making - The Lévy-Flight Model

In the following, we estimate the parameters of a stochastic differential equation model of human decision making. We perform the first Bayesian treatment of the recently proposed Lévy-Flight Model (LFM), as its intractability has so far rendered traditional non-amortized Bayesian inference methods prohibitively slow [49].

With this example, we first want to show empirically that BayesFlow is able to deal with *i.i.d.* datasets of variable size arising from N independent runs of a complex stochastic simulator. For this, we inspect global performance of BayesFlow over a wide range of dataset sizes. Additionally, we want to show the advantage of amortized inference compared to case-based inference in terms of efficiency and recovery. For this, we apply BayesFlow along with four other recent methods for likelihood-free inference to a single dataset and show that in some cases the speed advantage of amortized inference becomes noticeable even after as few as 5 datasets. Crucially, researchers often fit the same models to different datasets, so if a pre-trained model exists, it would present a huge advantage in terms of efficiency and productivity.

We focus on the family of evidence accumulator models (EAMs) which describe human decision making by a set of neurocognitively motivated parameters [42]. EAMs are most often applied to choice reaction times (RT) data to obtain an estimate of the underlying processes governing categorization and (perceptual) decision making. In its most general formulation, the forward model of EAMs takes the form of a stochastic ordinary differential equation (ODE):

$$dx = vdt + \xi\sqrt{dt} \quad (32)$$

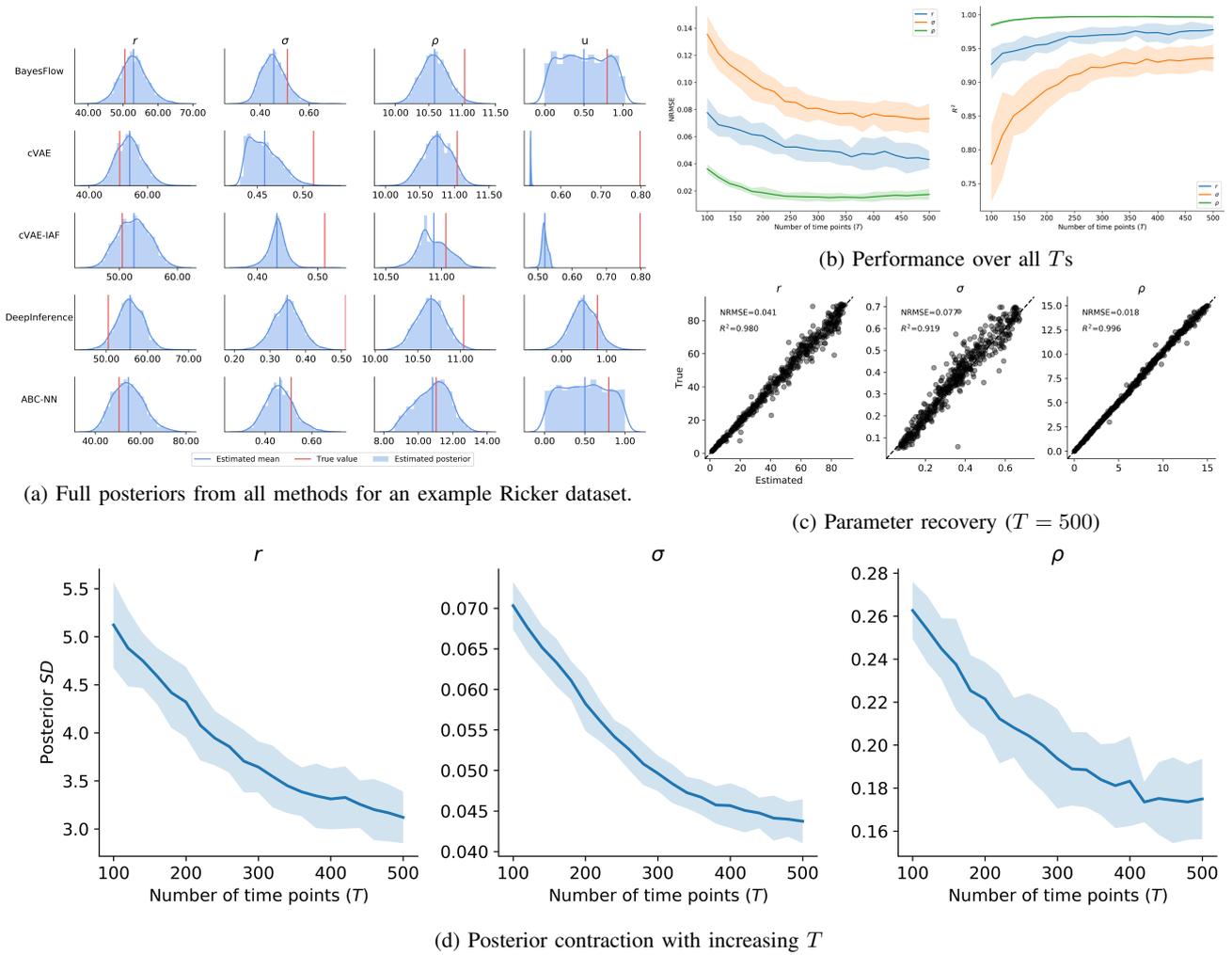


Fig. 4: Results on the Ricker model. (a) Approximate posteriors obtained by all implemented methods on a single Ricker dataset. Note that only BayesFlow and ABC-NN are able to approximate the uniform posterior of u ; (b) NRMSE and R^2 performance metrics over all T s obtained by the BayesFlow method. We observe that parameter estimation remains good over all T s, and becomes progressively better as more data is available (shaded regions indicate bootstrap 95% CIs); (c) Parameter recovery with BayesFlow for the maximum number of generations used during training ($T = 500$); (d) Posterior contraction in terms of posterior standard deviation for each parameter across increasing number of available generations (shaded regions indicate bootstrap 95% CIs).

TABLE I: Performance results on the Ricker model across all estimation methods

		BayesFlow	cVAE	cVAE-IAF	DeepInference	ABC-NN	ABC-RF
Err_{cal}	r	0.017 ± 0.007	0.014 ± 0.007	0.058 ± 0.017	0.122 ± 0.016	0.164 ± 0.015	-
	σ	0.013 ± 0.007	0.419 ± 0.011	0.382 ± 0.013	0.184 ± 0.021	0.119 ± 0.014	-
	ρ	0.084 ± 0.018	0.121 ± 0.017	0.188 ± 0.018	0.111 ± 0.019	0.283 ± 0.012	-
$NRMSE$	r	0.041 ± 0.002	0.047 ± 0.004	0.047 ± 0.006	0.052 ± 0.003	0.053 ± 0.003	0.044 ± 0.004
	σ	0.077 ± 0.005	0.137 ± 0.004	0.124 ± 0.006	0.108 ± 0.004	0.077 ± 0.004	0.081 ± 0.005
	ρ	0.018 ± 0.001	0.016 ± 0.002	0.019 ± 0.002	0.019 ± 0.002	0.033 ± 0.002	0.021 ± 0.001
R^2	r	0.980 ± 0.003	0.973 ± 0.005	0.973 ± 0.007	0.968 ± 0.005	0.966 ± 0.004	0.977 ± 0.004
	σ	0.919 ± 0.011	0.745 ± 0.020	0.792 ± 0.020	0.841 ± 0.014	0.919 ± 0.010	0.912 ± 0.011
	ρ	0.996 ± 0.001	0.997 ± 0.001	0.996 ± 0.001	0.996 ± 0.001	0.986 ± 0.002	0.994 ± 0.001
Err_{sim}	-	0.038 ± 0.001	0.041 ± 0.001	0.042 ± 0.001	0.041 ± 0.001	0.048 ± 0.002	0.041 ± 0.002

Note: For each parameter, bootstrapped means (± 1 standard error) of different performance metrics are displayed for all tested methods. For each metric and each parameter, the best performance across methods is printed in bold font.

where dx denotes a change in activation of an accumulator, v denotes the average speed of information accumulation (often termed the drift rate), and ξ represents a stochastic additive component, usually modeled as coming from a Gaussian distribution centered around 0: $\xi \sim \mathcal{N}(0, c^2)$.

EAMs are particularly amenable for likelihood-free inference, since the likelihood of most interesting members of this model family turn out to be intractable [34]. This intractability has precluded many interesting applications and empirically driven model refinements. Here, we apply BayesFlow to estimate the parameters of the recently proposed Lévy-Flight Model (LFM) [49]. The LFM assumes an α -stable noise distribution of the evidence accumulation process which allows to model discontinuities in the decision process. However, the inclusion of α -stable noise (instead of the typically assumed Gaussian noise) leads to a model with intractable likelihood:

$$dx = vdt + \xi dt^{1/\alpha} \quad (33)$$

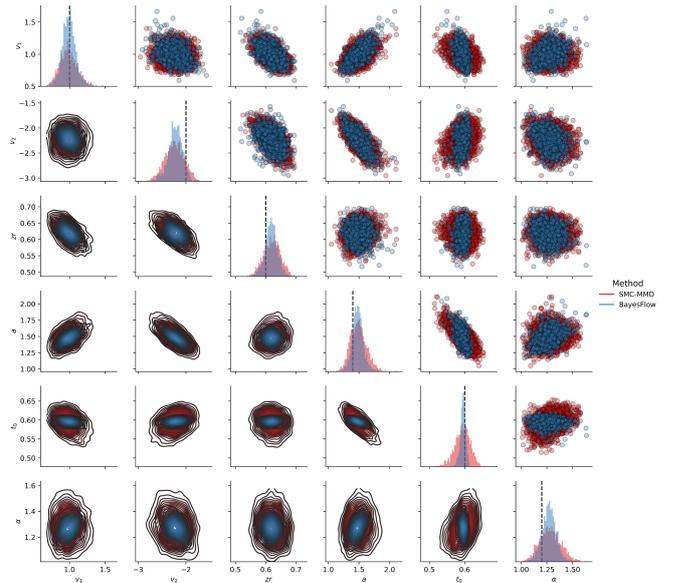
$$\xi \sim \text{AlphaStable}(\alpha, 0, 1, 0) \quad (34)$$

where α controls the probability of outliers in the noise distribution. The LFM has three additional parameters: the threshold a determining the amount of evidence needed for the termination of a decision process; a relative starting point, zr , determining the amount of starting evidence available to the accumulator before the actual decision alternatives are presented; and an additive non-decision time t_0 .

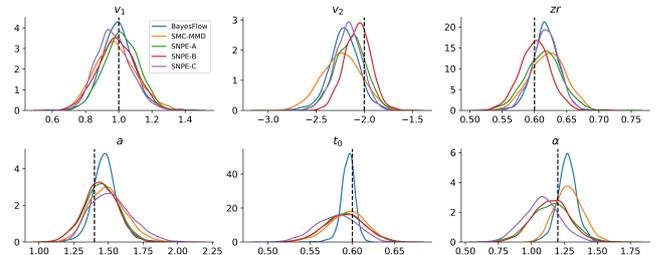
During training of the networks, we simulate response times data from two experimental conditions with two different drift rates, since such a design is often encountered in psychological research. The parameter estimation task is thus to recover the parameters $\theta = (v_0, v_1, a, t_0, zr, \alpha)$ from two-dimensional *i.i.d.* RT data $\mathbf{x}_{1:N}$ where each $\mathbf{x}_i \in \mathbb{R}^2$ represents RTs obtained in the two conditions. The number of trials is drawn from a uniform distribution $N \sim \mathcal{U}(100, 1000)$ at each training iteration. Training the networks took a little less than a day with the online learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approximately 7.39 seconds.

In order to investigate whether amortized inference is advantageous for this model, we additionally apply a version of the SMC-ABC algorithm available in the *pyABC* package [27] to a single dataset with $N = 500$. Since no sufficient summary statistics are available for EAM data, we apply the maximum mean discrepancy (MMD) metric as a distance between the full raw empirical RT distributions, in order to prevent information loss [39]. Since the MMD is expensive to compute, we use a GPU implementation to ensure that computation of MMD is not a bottleneck for the comparison. In order to achieve good approximation with 2000 samples from the SMC-MMD approximate posterior, we run the algorithm for 20 populations with a final rejection threshold $\epsilon = 0.04$. We also draw 2000 samples from the approximate posterior obtained by applying our pre-trained BayesFlow networks to the same dataset.

Along SMC-MMD, we apply three recent methods for neural density estimation, SNPE-A [37], SNPE-B [30], and SNPE-C ([16], also dubbed APT). Since these methods all depend on summary statistics of the data, we compute the first 6 moments



(a) Joint posteriors from BayesFlow and SMC-MMD



(b) Marginal posteriors from all methods

Fig. 5: Comparison results on the LFM model. (a) Marginal and bivariate posteriors obtained by BayesFlow and SMC-MMD on the single validation dataset. We observe markedly better sharpness in the BayesFlow posteriors; (b) Marginal posteriors obtained from all methods under comparison.

of each empirical response time distribution as well as the fractions of correct/wrong responses. We train each method for a single round with 100 epochs and 5000 simulated datasets, in order to keep running time at a minimum. Also, we did not observe improvement in performance when training for more than one round. For each model, we sample 2000 samples from the approximate joint posterior to align the number of samples with those obtained via SMC-MMD.

The comparison results are depicted in Figure 5. We first focus on the comparison with SMC-MMD on the single dataset. Figure 5a depicts marginal and bivariate posteriors obtained by BayesFlow and SMC-MMD. The approximate posteriors of BayesFlow appear noticeably sharper. Observing the SCB plots (Figure 6b), we can conclude that the approximate posteriors of BayesFlow mirror the sharpness of the true posterior, since otherwise the SCB plots would show marked deviations from uniformity. Further, Figure 5b depicts the marginal posteriors obtained from the application of each method. Noticeably, performance and sharpness varies across the methods and parameters, with all methods yielding good point-estimate

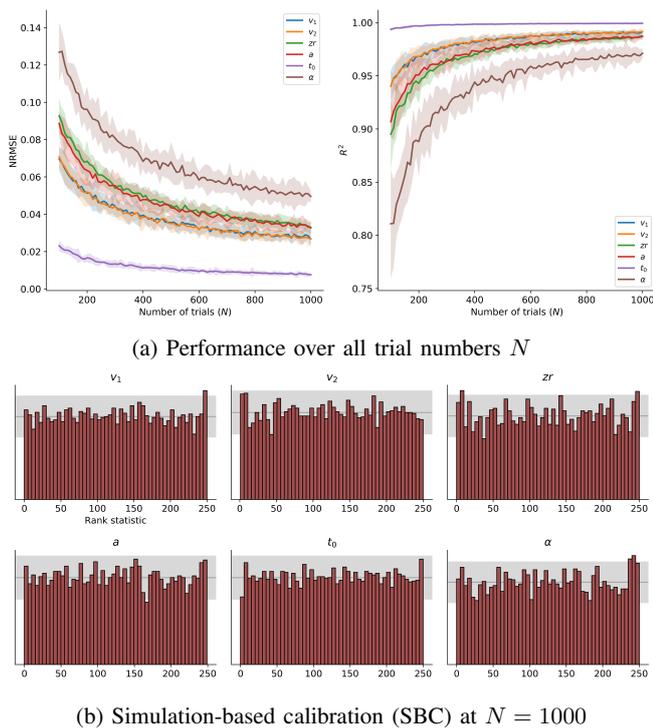


Fig. 6: BayesFlow results obtained on the LFM model.

TABLE II: Speed of inference and break-even for amortized inference for the LFM model

	Upfront Training	Inference (1 dataset)	Inference (500 datasets)	Break-even after
BayesFlow	23.2 h	60 ms	3.7 s	-
SMC-MMD	-	5.5 h	2700 h	5 datasets
SNPE-A	-	0.65 h	325 h	37 datasets
SNPE-B	-	0.65 h	325 h	37 datasets
SNPE-C	-	0.35 h	175 h	75 datasets

Note: Inference times for 500 datasets as well as the number of datasets for break-even with BayesFlow for the SMC-MMD, SNPE-A, SNPE-B, and SNPE-C methods are extrapolated from the wall-clock running time on a single dataset, so these are approximate quantities.

recovery via posterior means in terms of the NRMSE and R^2 metrics.

Importantly, Table II summarizes the advantage of amortized inference for the LFM model in terms of efficiency. For instance, compared to SMC-MMD, the extra effort of learning a global BayesFlow model upfront is worthwhile even after as few as 5 datasets, as inference with SMC-MMD would have taken more than a day to finish. On the other hand, the break-even for SNPE-C/APT occurs after 75 datasets, so in cases where only a few dozens of datasets are considered, case-based inference might be preferable. However, the difficulties in manually finding meaningful and efficiently computable summary statistics may eat up possible savings even in this situation. We acknowledge that our choices in this respect might be sub-optimal, so performance comparisons should be treated with some caution.

We note, that after a day of training, the pre-trained networks of BayesFlow take less than 5 seconds to perform inference on 500 datasets even with maximum number of trials $N = 1000$. Using the case-based SMC-MMD algorithm, 500 inference

runs would have taken more than half a year to complete. We also note, that parallelizing separate inference threads across multiple cores or across nodes of a (GPU) computing cluster can dramatically increase the wall-clock speed of the case-based methods considered here. However, the same applies to BayesFlow training, since its most expensive part, the simulation from the forward model, would profit the most from parallel computing.

The global performance of BayesFlow over all validation datasets and all trial sizes N is depicted in (Figure 6). First, we observe excellent recovery of all LFM parameters with NRMSEs ranging between 0.008 and 0.048 and R^2 between 0.972 and 0.99 for the maximum number of trials. Importantly, estimation remains very good across all trial numbers, and improves as more trials become available (Figure 6a). The parameter α appears to be most challenging to estimate, requiring more data for good estimation, whereas the non-decision time parameter t_0 can be recovered almost perfectly for all trial sizes. Last, the SCB histograms indicate no systematic deviations across the marginal posteriors (Figure 6b).

G. Stochastic Differential Equations - The SIR Epidemiology Model

With this example, we want to further corroborate the excellent global performance and probabilistic calibration observed for the LFM model on a non-*i.i.d.* stochastic ODE model. For this, we study a compartmental model from epidemiology, whose output comprises variable-sized multidimensional and inter-dependent time-series. It is therefore of interest to investigate how our method performs when applied to data which is the direct output of an ODE simulator.

Compartmental models in epidemiology describe the stochastic dynamics of infectious diseases as they spread over a population of individuals [23], [20]. The parameters of these models encode important characteristics of diseases, such as infection and recovery rates. The stochastic SIR model describes the transition of N individuals between three discrete states – susceptible (S), infected (I), and recovered (R) – whose dynamics follow the equations:

$$\Delta S = -\Delta N_{SI} \quad (35)$$

$$\Delta I = \Delta N_{SI} - \Delta N_{IR} \quad (36)$$

$$\Delta R = \Delta N_{IR} \quad (37)$$

$$\Delta N_{SI} \sim \text{Binomial}(S, 1 - \exp\left(-\beta \frac{I}{N} \Delta t\right)) \quad (38)$$

$$\Delta N_{IR} \sim \text{Binomial}(I, 1 - \exp(-\gamma \Delta t)) \quad (39)$$

where $S + I + R = N$ give the number of susceptible, infected, and recovered individuals, respectively. The parameter β controls the transition rate from being susceptible to infected, and γ controls the transition rate from being infected to recovered. The above listed stochastic system has no analytic solution and thus requires numerical simulation methods for recovering parameter values from data. Cast as a parameter estimation task, the challenge is to recover $\theta = \{\beta, \gamma\}$ from three dimensional time-series data $\mathbf{x}_{1:T}$ where each $\mathbf{x}_t \in \mathbb{N}^3$ is a triple containing the number of susceptible (S), number of infected (I), and recovered (R) individuals at time t .

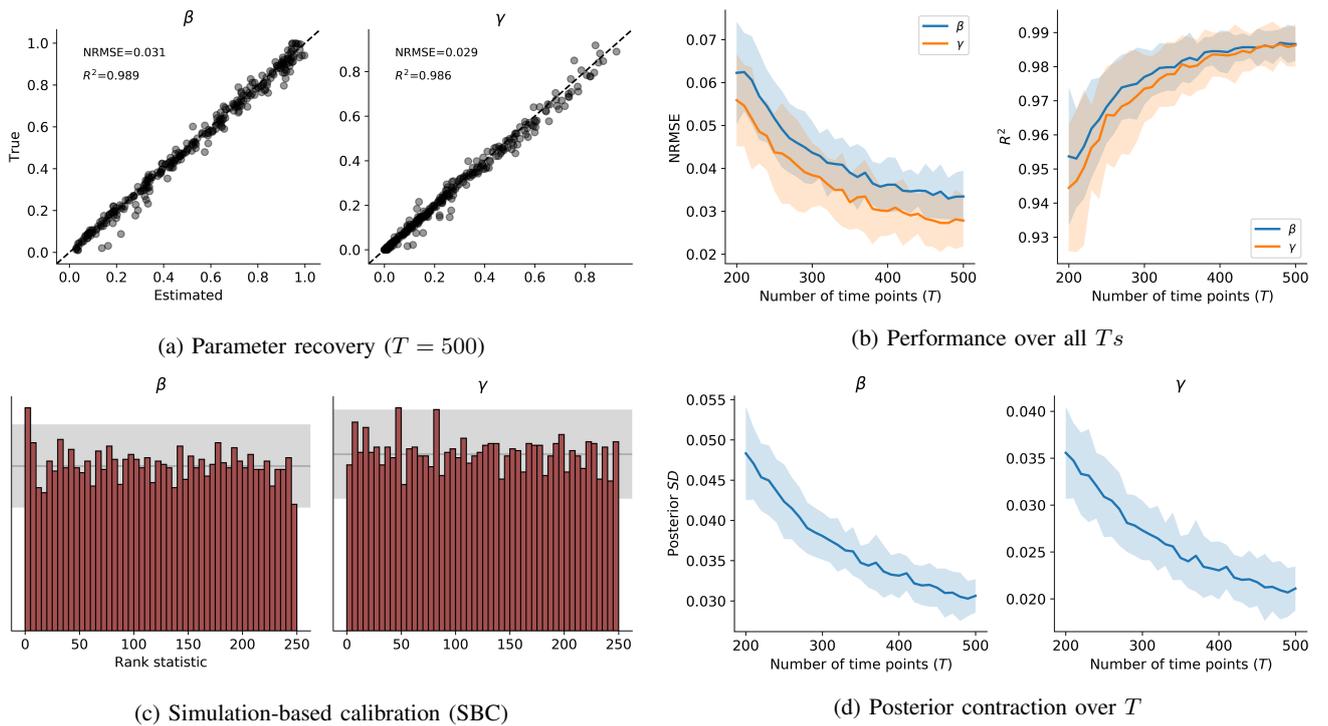


Fig. 7: Results obtained on the stochastic SIR model.

During training of the networks, we simulate time-series from the stochastic SIR model with varying lengths. The number of time points T is drawn from a uniform distribution $T \sim \mathcal{U}(200, 500)$ at each training iteration. For small T , the system has not yet reached an equilibrium (i.e., not all individuals have transitioned from being I to R). It is especially interesting to see if BayesFlow can recover the rate parameters, while the process dynamics are still unfolding over time. Training the networks took approximately two hours with the online learning approach. Inference on 1000 datasets with 2000 posterior samples per parameter took approximately 1.1 seconds.

The results on the SIR model are depicted in Figure 7. In line with the previous examples, we observe very good recovery of the true parameters, with NRMSE at $T = 500$ around 0.03, and R^2 s around 0.99. We observe decent performance even at smaller T s and the expected improvements as T increases. Specifically, the posterior variance shrinks as T increases. The SCB plots indicate that the approximate posteriors are well calibrated, with the approximate posterior mean of β slightly overestimating the true parameter values in the lower range.

H. Learned vs. Hand-Crafted Summaries: The Lotka-Volterra Population Model

See **Appendix A**

IV. DISCUSSION

In the current work, we proposed and explored a novel method which uses invertible neural networks to perform globally amortized approximate Bayesian inference. The method, which we named BayesFlow, requires only simulations from

a forward model to learn an efficient probabilistic mapping between data and parameters. We demonstrated the utility of BayesFlow by applying it to models and data from various research domains. Further, we explored an online learning approach with variable number of observations per iteration. We demonstrated that this approach leads to excellent parameter estimation throughout the examples considered in the current work. In theory, BayesFlow is applicable to any mathematical forward model which can be implemented as a computer simulation. In the following, we highlight the main advantages of BayesFlow.

First, the introduction of separate summary and inference networks renders the method independent of the shape or the size of the observed data. The summary network learns a fixed-size vector representation of the data in an automatic, data-driven manner. Since the summary network is optimized jointly with the inference network, the learned data representation is encouraged to be maximally informative for inferring the parameters' posterior. This is particularly useful in settings where appropriate summary statistics are not known and, as a consequence, relevant information is lost through the choice of sub-optimal summary functions. However, if *sufficient* statistics are available in a given domain, one might omit the summary network altogether and feed these statistics directly to the invertible network.

Second, we showed that BayesFlow generates samples from the correct posterior under perfect convergence without distributional assumptions on the shape of the posterior. This is in contrast to variational methods which optimize a lower-bound on the posterior [26], [24], and oftentimes assume

Gaussian approximate posteriors. Additionally, we also showed throughout all examples that the posterior means generated by the BayesFlow method are mostly excellent estimates for the true values. Beyond this, the fact that the BayesFlow method recovers the full posterior over parameters does not necessitate the usage of point estimates or summary statistics of the posterior. Further, we observe the desired posterior contraction (posterior variance decreases with increasing number of observations) and better recovery with increasing number of observations. These are indispensable properties of any Bayesian parameter estimation method, since they mirror the decrease in epistemic uncertainty and the simultaneous increase in information due to availability of more data.

Third, the largest computational cost of BayesFlow is paid during the training phase. Once trained, the networks can efficiently compute the posterior for any observed dataset arising from the forward model. This is similar to the recently introduced *prepaid method* [33]. However, this method memorizes a large database of pre-computed summary statistics for fast nearest-neighbor inference, whereas a BayesFlow’s network weights define an abstract representation of the relationship between data and parameters over the whole space of hidden parameters. Traditionally, abstract representations like this only existed for analytically invertible model families, whereas more complex forward models required case-based inference, that is, expensive re-training for each observed dataset. Amortized inference as realized by BayesFlow is thus especially advantageous for exploring, testing and comparing competing scientific hypotheses in research domains where an intractable model needs to be fit to multiple independent datasets.

Finally, all computations in the BayesFlow method benefit from a high degree of parallelism and can thus utilize the advantages of modern GPU acceleration.

These advantages notwithstanding, limitations of the proposed method should also be mentioned. Although we could provide a theoretical guarantee that BayesFlow samples from the true joint posterior under perfect convergence, this might not be achieved in practice. Therefore, it is essential that proper calibration of point estimates and estimated joint posteriors is performed for each application of the method. Fortunately, validating a trained BayesFlow architecture is easy due to amortized inference. Below, we discuss potential challenges and limitations of the method.

First, the design of the summary network and inference networks is a crucial choice for achieving optimal performance of the method. As already mentioned, the summary network should be able to represent the observed data without losing essential information and the invertible network should be powerful enough to capture the behavior of the forward model. Nevertheless, in some real-world scenarios, there might be little guidance on how to actually construct suitable summary networks. Recent work on probabilistic symmetry [6] and algorithmic alignment [52] as well as our current experiments do, however, provide some insights about the design of summary networks. For instance, *i.i.d.* data induce a permutation invariant distribution which is well modeled with a deep invariant network [6]. Data with temporal or spatial dependencies are best

modeled with recurrent [21], or convolutional [41] networks. When pairwise or multi-way relationships are particularly informative, attention [48] or graph networks [52] appear as reasonable choices. On the other hand, the depth of the invertible network should be tailored to the complexity of the mathematical model of interest. More ACBs will enable the network to encode more complex distributions but will increase training time. Very high-dimensional problems might also require very large networks with millions of parameters, up to a point where estimation becomes practically unfeasible. However, most mathematical models in the life sciences prioritize parsimony and interpretability, so they do not contain hundreds or thousands of latent parameters. In any case, future applications might require novel network architectures and solutions which go beyond our initial recommendations.

Another potential issue is the large number of neural network and optimization hyperparameters that might require fine-tuning by the user for optimal performance on a given task. We observe that excellent performance is often achieved with default settings. Using larger networks consisting of 5 to 10 ACBs does not seem to hurt performance or destabilize training, even if the model to be learned is relatively simple. Based on our results, we expect that a single architecture should be able to perform well on models from a given domain. Future research should investigate this question of generality by applying the method to different or even competing models within different research domains. Future research should investigate the impact of modern hyperparameter optimization methods such as Bayesian optimization [14].

Finally, even though modern deep learning libraries allow for rapid and relatively straightforward development of various neural network architectures, the implementational burden associated with the method is non-trivial. Thus, we are currently developing a general user-friendly software, which will abstract away most intricacies from the users of our method.

We hope that the new BayesFlow method will enable researchers from a variety of fields to accelerate model-based inference and will further prove its utility beyond the examples considered in this paper.

ACKNOWLEDGMENT

We thank Paul Bürkner, Manuel Haussmann, Jeffrey Rouder, Raphael Hartmann, David Izydorczyk, Hannes Wendler, Chris Wendler, and Karin Prillinger for their invaluable comments and suggestions that greatly improved the manuscript. We also thank Francis Tuerlinckx and Stijn Verdonck for their support and thought-provoking ideas.

REFERENCES

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation*, pages 265–283, 2016.
- [2] Lynton Ardizzone, Jakob Kruse, Sebastian Wirkert, Daniel Rahner, Eric W Pellegrini, Ralf S Klessen, Lena Maier-Hein, Carsten Rother, and Ullrich Köthe. Analyzing inverse problems with invertible neural networks. In *Intl. Conf. on Learning Representations*, 2019.

- [3] Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. Guided image generation with conditional invertible neural networks. *arXiv:1907.02392*, 2019.
- [4] Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- [5] Daniele Bigoni, Olivier Zahm, Alessio Spantini, and Youssef Marzouk. Greedy inference with layers of lazy maps. *arXiv:1906.00031*, 2019.
- [6] Benjamin Bloem-Reddy and Yee Whye Teh. Probabilistic symmetry and invariant neural networks. *arXiv:1901.06082*, 2019.
- [7] William M Bolstad and James M Curran. *Introduction to Bayesian statistics*. John Wiley & Sons, 2016.
- [8] Laming Chen, Guoxin Zhang, and Eric Zhou. Fast greedy map inference for determinantal point process to improve recommendation diversity. In *Advances in Neural Information Processing Systems*, pages 5622–5633, 2018.
- [9] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *arXiv:1911.01429*, 2019.
- [10] Katalin Csilléry, Michael GB Blum, Oscar E Gaggiotti, and Olivier François. Approximate Bayesian Computation (ABC) in Practice. *Trends in Ecology & Evolution*, 25(7):410–418, 2010.
- [11] Matthew C Deans. Maximally informative statistics for localization and mapping. In *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, volume 2, pages 1824–1829. IEEE, 2002.
- [12] Gianluca Detommaso, Jakob Kruse, Lynton Ardizzone, Carsten Rother, Ullrich Köthe, and Robert Scheichl. HINT: hierarchical invertible neural transport for general and sequential bayesian inference. *arXiv:1905.10687*, 2019.
- [13] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. *arXiv:1605.08803*, 2016.
- [14] Katharina Eggensperger, Matthias Feurer, Frank Hutter, James Bergstra, Jasper Snoek, Holger Hoos, and Kevin Leyton-Brown. Towards an empirical foundation for assessing bayesian optimization of hyperparameters. In *NIPS workshop on Bayesian Optimization in Theory and Practice*, volume 10, page 3, 2013.
- [15] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with LSTM. *Neural Computation*, 12(19):2451–2471, 2000.
- [16] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior transformation for likelihood-free inference. *arXiv:1905.07488*, 2019.
- [17] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773, 2012.
- [18] Aditya Grover, Manik Dhar, and Stefano Ermon. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [19] John R Hershey and Peder A Olsen. Approximating the kullback leibler divergence between gaussian mixture models. In *2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07*, volume 4, pages IV–317. IEEE, 2007.
- [20] Herbert W Hethcote. The mathematics of infectious diseases. *SIAM review*, 42(4):599–653, 2000.
- [21] Seong Jae Hwang, Zirui Tao, Won Hwa Kim, and Vikas Singh. Conditional recurrent flow: Conditional generation of longitudinal samples with applications to neuroimaging. *arXiv:1811.09897*, 2018.
- [22] Bai Jiang, Tung-yu Wu, Charles Zheng, and Wing H Wong. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.
- [23] Matt J Keeling and Pejman Rohani. *Modeling infectious diseases in humans and animals*. Princeton University Press, 2011.
- [24] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *stat*, 1050:1, 2014.
- [25] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pages 10215–10224, 2018.
- [26] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pages 4743–4751, 2016.
- [27] Emmanuel Klinger, Dennis Rickert, and Jan Hasenauer. pyabc: distributed, likelihood-free inference. *Bioinformatics*, 34(20):3591–3593, 2018.
- [28] Hongzhou Lin and Stefanie Jegelka. Resnet with one-neuron hidden layers is a universal approximator. In *Advances in Neural Information Processing Systems*, pages 6169–6178, 2018.
- [29] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pages 3431–3440, 2015.
- [30] Jan-Matthis Lueckmann, Pedro J Goncalves, Giacomo Bassetto, Kaan Öcal, Marcel Nonnenmacher, and Jakob H Macke. Flexible statistical inference for mechanistic models of neural dynamics. In *Advances in Neural Information Processing Systems*, pages 1289–1299, 2017.
- [31] Ulf Kai Mertens. *Deep learning methods for likelihood-free inference: approximating the posterior distribution with convolutional neural networks*. PhD thesis, Heidelberg University, 2019.
- [32] Ulf Kai Mertens, Andreas Voss, and Stefan Radev. Abrox—a user-friendly python module for approximate bayesian computation with a focus on model comparison. *PLoS one*, 13(3):e0193981, 2018.
- [33] Merijn Mestdagh, Stijn Verdonck, Kristof Meers, Tim Loossens, and Francis Tuerlinckx. Prepaid parameter estimation without likelihoods. *PLoS computational biology*, 15(9):e1007181, 2019.
- [34] Steven Miletić, Brandon M Turner, Birte U Forstmann, and Leendert van Maanen. Parameter recovery for the leaky competing accumulator model. *Journal of Mathematical Psychology*, 76:25–50, 2017.
- [35] Ashish Mishra, Shiva Krishna Reddy, Anurag Mittal, and Hema A Murthy. A generative model for zero shot learning using conditional variational autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 2188–2196, 2018.
- [36] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [37] George Papamakarios and Iain Murray. Fast ϵ -free inference of simulation models with bayesian conditional density estimation. In *Advances in Neural Information Processing Systems*, pages 1028–1036, 2016.
- [38] George Papamakarios, David C Sterratt, and Iain Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. *arXiv:1805.07226*, 2018.
- [39] Mijung Park, Wittawat Jitkrittum, and Dino Sejdinovic. K2-ABC: approximate bayesian computation with kernel embeddings. In *Intl. Conf. Artificial Intelligence and Statistics*, pages 398–407, 2016.
- [40] Matthew D Parno and Youssef M Marzouk. Transport map accelerated markov chain monte carlo. *SIAM/ASA Journal on Uncertainty Quantification*, 6(2):645–682, 2018.
- [41] Stefan T Radev, Ulf K Mertens, Andreas Voss, and Ullrich Köthe. Towards end-to-end likelihood-free inference with convolutional neural networks. *British Journal of Mathematical and Statistical Psychology*, 2019.
- [42] Roger Ratcliff and Gail McKoon. The diffusion decision model: theory and data for two-choice decision tasks. *Neural computation*, 20(4):873–922, 2008.
- [43] Louis Raynal, Jean-Michel Marin, Pierre Pudlo, Mathieu Ribatet, Christian P Robert, and Arnaud Estoup. ABC random forests for Bayesian parameter inference. *Bioinformatics*, 35(10):1720–1728, 2018.
- [44] Scott A Sisson and Yanan Fan. *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839], 2011.
- [45] Mikael Sunnåker, Alberto Giovanni Busetto, Elina Numminen, Jukka Corander, Mathieu Foll, and Christophe Dessimoz. Approximate bayesian computation. *PLoS computational biology*, 9(1):e1002803, 2013.
- [46] Sean Talts, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. Validating bayesian inference algorithms with simulation-based calibration. *arXiv:1804.06788*, 2018.
- [47] Brandon M Turner and Per B Sederberg. A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review*, 21(2):227–250, 2014.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [49] Andreas Voss, Veronika Lerche, Ulf Mertens, and Jochen Voss. Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models. *Psychonomic bulletin & review*, pages 1–20, 2019.
- [50] Darren J Wilkinson. *Stochastic modelling for systems biology*. CRC press, 2011.
- [51] Simon N Wood. Statistical inference for noisy nonlinear ecological dynamic systems. *Nature*, 466(7310):1102, 2010.
- [52] Keyulu Xu, Jingling Li, Mzhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *arXiv:1905.13211*, 2019.

APPENDIX

A. Learned vs. Hand-Crafted Summaries: The Lotka-Volterra Population Model

With this final example, we want to compare the performance of our method with an LSTM summary network vs. performance obtained with a standard set of hand-crafted summary statistics. For this, we focus on the well-studied Lotka-Volterra (LV) model. The LV model describes the dynamics of biological systems in which a population of predators interacts with a population of prey [50]. It involves a pair of first order, non-linear, differential equations given by:

$$\frac{d}{dt} = \alpha u - \beta uv \quad (1)$$

$$\frac{d}{dt} = -\gamma v + \delta \beta uv \quad (2)$$

where u denotes the number of preys, v denotes the number of predators, and the parameter vector controlling the interaction between the species is $\theta = (\alpha, \beta, \gamma, \delta)$.

During training of the networks, we set the initial conditions as $u_0 = 10$ and $v_0 = 5$ and consider an interval $I_T = 15$ of discrete time units with $T = 500$ time steps (samples) in between. Each sample \mathbf{x}_t in each LV time-series $\mathbf{x}_{1:T}$ is thus a 2-dimensional vector containing the number of prey and predators in the population at time unit t .

We train two invertible neural networks. The first is trained jointly with an LSTM summary network which outputs a 9-dimensional learned summary statistic $h_\psi(\mathbf{x}_{1:T})$. The second uses a set of 9 typically used, hand-crafted summary statistics [37], [38], which include: the mean of the time series; the log variance of the time-series; the auto-correlation of each timeseries at lags 0.2 and 0.4 time units; the cross-correlation between the two time series. The same cINN architecture with 5 ACBs is used for both training scenarios. For each scenario, we perform the same number of iterations and epochs. Online learning for each training scenario took approximately 4 hours in total wall-clock time.

The results obtained on the LV model are depicted in Figure S1. We observe notably better recovery of the true parameter estimates when performing inference with the learned summary statistics. The approximate posteriors are also better calibrated when conditioned on the set of 9 learned summary statistics. These results highlight the advantages of using a summary networks when no sufficient summary statistics are available. Finally, Figure S1e and Figure S1f depict the posteriors obtained by the two different INNs on a single dataset with ground-truth parameters $\theta = (1, 1, 1, 1)$. Evidently, learning the summary statistics leads to much sharper posteriors and better point-estimate recovery.

B. Computation of Validation Metrics

Normalized Root Mean Squared Error: The normalized root mean squared error (NRMSE) between a sample of true parameters $\{\theta^{(m)}\}_{m=1}^M$ and a sample of estimated parameters $\{\hat{\theta}^{(m)}\}_{m=1}^M$ is given by:

$$NRMSE = \sqrt{\sum_{m=1}^M \frac{(\theta^{(m)} - \hat{\theta}^{(m)})^2}{\theta_{max} - \theta_{min}}} \quad (3)$$

Due to the normalization factor $\theta_{max} - \theta_{min}$, the NRMSE is scale-independent, and thus suitable for comparing the recovery across parameters with different numerical ranges. The NRMSE equals zero when the estimates are exactly equal to the true values.

Coefficient of Determination : The coefficient of determination R^2 measures the proportion of variance in a sample of true parameters $\{\theta^{(m)}\}_{m=1}^M$ that is explained by a sample of estimated parameters $\{\hat{\theta}^{(m)}\}_{m=1}^M$. It is computed as:

$$R^2 = 1 - \sum_{m=1}^M \frac{(\theta^{(m)} - \hat{\theta}^{(m)})^2}{(\theta^{(m)} - \bar{\theta}^{(m)})^2} \quad (4)$$

where $\bar{\theta}$ denotes the mean of the true parameter samples. When R^2 equals 1, the estimates are perfect reconstructions of the true parameters.

Re-simulation Error: To compute the re-simulation error Err_{sim} , we first obtain an estimate of the true parameter value given an observed (validations) dataset $\mathbf{x}_{1:N}^o$ by computing the mean of the approximate posterior $\hat{\theta}$. Then, we run the mathematical model to obtain a simulated dataset $\mathbf{x}_{1:N}^s = g(\hat{\theta}, \xi)$. Finally, we compute the maximum mean discrepancy (MMD, [17]) between the observed and the simulated dataset $MMD(\mathbf{x}_{1:N}^o, \mathbf{x}_{1:N}^s)$. The MMD is a kernel-based metric which estimates the mismatch between two distributions given samples from the distributions by comparing all of their moments. It equals zero when the two distributions are equal almost everywhere [17]). Thus, a low MMD indicates that the distribution of $\mathbf{x}_{1:N}^s$ is close to the distribution of $\mathbf{x}_{1:N}^o$. Conversely, a high MMD indicates that the distribution of $\mathbf{x}_{1:N}^s$ is far from the distribution of $\mathbf{x}_{1:N}^o$. We report the median MMD computed over all validation datasets.

Calibration Error: The calibration error Err_{cal} quantifies how well the coverage of an approximate posterior matches the coverage of an unknown true posterior. Let α_θ be the fraction of true parameter values lying in a corresponding α -credible interval of the approximate posterior. Thus, for a perfectly calibrated approximate posterior, α_θ should equal α for all $\alpha \in (0, 1)$. We compute the calibration error for each marginal posterior as the median absolute deviation $|\alpha_\theta - \alpha|$ for 100 equally spaced values of $\alpha \in (0, 1)$. Therefore, the calibration error ranges between 0 and 1 with 0 indicating perfect calibration and 1 indicating complete miscalibration of the approximate posterior.

Kullback-Leibler Divergence: The Kullback-Leibler divergence (KL) quantifies the increase in entropy incurred by approximating a target probability distribution P with a distribution Q . Its general form for absolutely continuous distributions is given by

$$KL(P || Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx \quad (5)$$

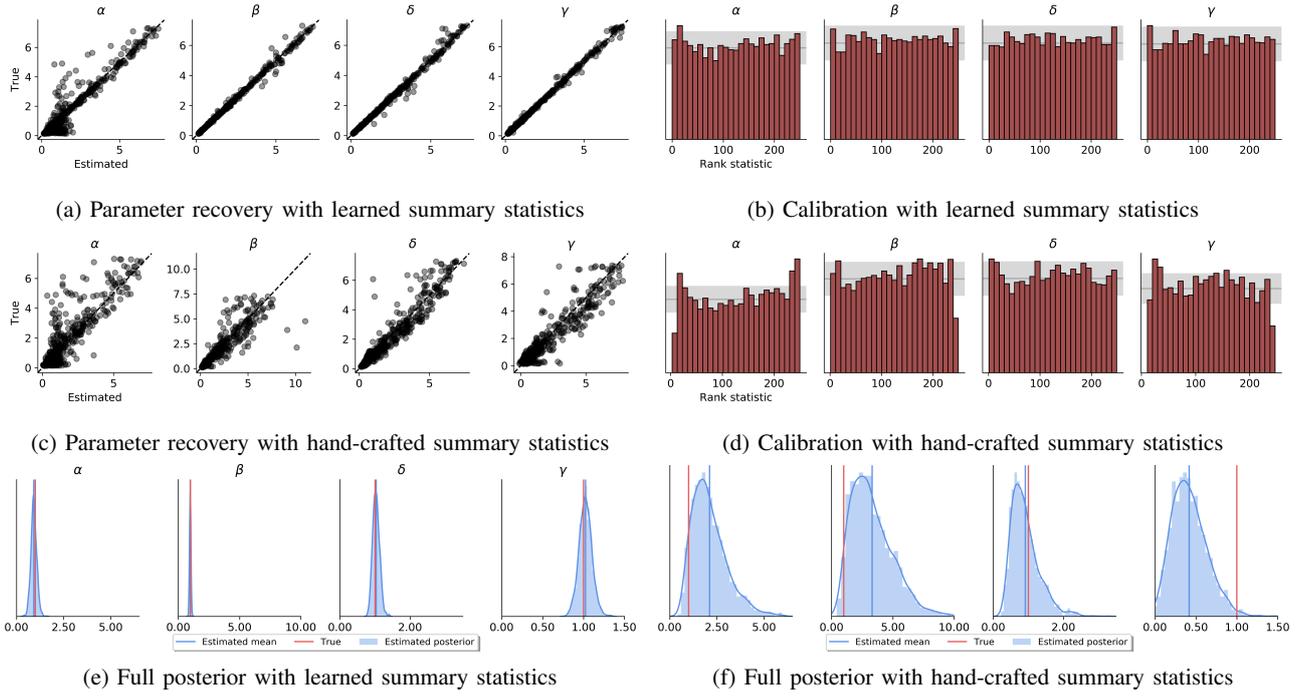


Fig. S1: Comparison of recovery/calibration on the LV model with learned vs. hand-crafted summary statistics **(a)** Simulation-based calibration (SBC) with learned summary statistics; **(b)** Parameter recovery with learned summary statistics; **(c)** Parameter recovery with hand-crafted summary statistics; **(d)** Simulation-based calibration (SBC) with hand-crafted summary statistics; **(e)** Example full posteriors obtained on a single dataset with ground-truth rate parameters $\theta = (1, 1, 1, 1)$ obtained with learned summaries; **(f)** The posterior obtained from the same dataset using hand-crafted summary statistics.

where p and q denote the pdfs of P and Q . In the case where P and Q are both multivariate Gaussian distributions, the KL divergence can be computed in closed form [19]:

$$\mathbb{KL}(P \parallel Q) = \frac{1}{2} \left[\log \frac{\det \Sigma_q}{\det \Sigma_p} + \text{Tr}(\Sigma_q^{-1} \Sigma_p) - d + (\mu_p - \mu_q)^T \Sigma_q^{-1} (\mu_p - \mu_q) \right] \quad (6)$$

where Σ_p and Σ_q denote the covariance matrices of p and q , μ_p and μ_q the respective mean vectors, and d the number of dimensions of the Gaussian. In the case of diagonal Gaussian distributions, Eq.6 reduces to:

$$\mathbb{KL}(P \parallel Q) = \sum_{i=1}^d \left(\log \frac{\sigma_{q,i}}{\sigma_{p,i}} + \frac{\sigma_{p,i}^2 + (\mu_{q,i} - \mu_{p,i})^2}{2\sigma_{q,i}^2} - \frac{1}{2} \right) \quad (7)$$

Even though the KL divergence is not a proper distance metric, as it is not symmetric in its arguments, it can be used to quantify the error of approximation when a closed-form solution is available.

Simulation-Based Calibration

Simulation-based calibration is a method to detect systematic biases in any Bayesian posterior sampling method [46]. It is based on the *self-consistency* of the Bayesian joint distribution.

Given a sample from the prior distribution $\tilde{\theta} \sim p(\theta)$ and a sample from the forward model $\tilde{x} \sim p(x | \tilde{\theta})$, one can integrate $\tilde{\theta}$ and \tilde{x} out of the joint distribution and recover back the prior of θ :

$$p(\theta) = \int p(\theta, \tilde{\theta}, \tilde{x}) d\tilde{x} d\tilde{\theta} \quad (8)$$

$$= \int p(\theta, \tilde{x} | \tilde{\theta}) p(\tilde{\theta}) d\tilde{x} d\tilde{\theta} \quad (9)$$

$$= \int p(\theta | \tilde{x}) p(\tilde{x} | \tilde{\theta}) p(\tilde{\theta}) d\tilde{x} d\tilde{\theta} \quad (10)$$

If the Bayesian sampling method produces samples from the exact posterior, the equality implied by Eq.10 should hold regardless of the particular form of the posterior. Thus, any violation of this equality indicates some error incurred by the sampling method. The authors of [46] propose **Algorithm 2** for visually detecting such violations:

Algorithm 2 is correct, since Eq.10 implies that the rank statistic defined in line 5 should be uniformly distributed. Hence, any deviations from uniformity indicate some interpretable error in the approximate posterior [46].

C. Model Details

The Ricker Model: Summary Network. We use a bidirectional long short-term memory (LSTM) recurrent neural network [15] for the raw Ricker time-series. The LSTM network architecture is a reasonable choice for this example, as it is able to capture

Algorithm 2 Simulation-based calibration (SBC) for a single parameter θ

- 1: **for** $m = 1, \dots, M$ **do**
- 2: Sample $\tilde{\theta}^{(m)} \sim p(\theta)$
- 3: Simulate a dataset $\mathbf{x}_{1:N}^{(m)} = g(\tilde{\theta}^{(m)}, \boldsymbol{\xi})$
- 4: Draw posterior samples $\{\theta^{(l)}\}_{l=1}^L \sim p_\phi(\theta | \mathbf{x}_{1:N}^{(m)})$
- 5: Compute rank statistic $r^{(m)} = \sum_{l=1}^L \mathbb{1}_{[\theta^{(l)} < \tilde{\theta}^{(m)}]}$
- 6: Store $r^{(m)}$
- 7: **end for**
- 8: Create a histogram of $\{r^{(i)}\}_{m=1}^M$ and inspect it for uniformity

long-term dependencies in datasets with temporal or spatial autocorrelations. LSTMs can also easily deal with variable-length time-series.

Simulation. We place the following uniform priors over the Ricker model parameters:

$$\rho \sim \mathcal{U}(0, 15) \quad (11)$$

$$r \sim \mathcal{U}(1, 90) \quad (12)$$

$$\sigma \sim \mathcal{U}(0.05, 0.7) \quad (13)$$

These ranges appear to be very broad, as datasets generated by extreme parameter values appear implausible in real-world scenarios. Nevertheless, we stick to broad priors for training, even though parameter recovery might degrade at the extremes.

Figure S2 depicts different simulated Ricker timeseries generated via draws from the prior.

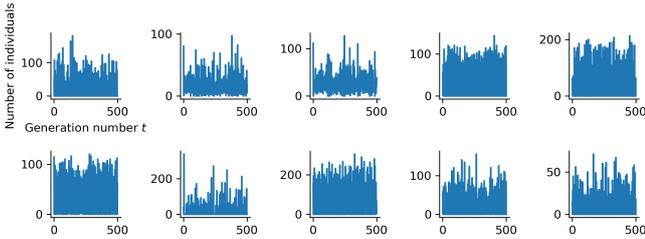


Fig. S2: Example Ricker datasets generated with different parameters.

The Lévy-Flight Model: Summary Network. We use a permutation invariant neural network [6] for the *i.i.d.* reaction times (RT) data. Similarly to the toy Regression example, each response in an RT dataset is assumed to be independent of all others, so permutations of the dataset must lead to the same parameter estimates.

Simulation. We place the following uniform priors over the LFM parameters, since they are broad enough to cover the range of realistic RT distributions encountered in empirical choice RT scenarios:

$$v_0 \sim \mathcal{U}(0, 6) \quad (14)$$

$$v_1 \sim \mathcal{U}(-6, 0) \quad (15)$$

$$zr \sim \mathcal{U}(0.3, 0.7) \quad (16)$$

$$a \sim \mathcal{U}(0.6, 3) \quad (17)$$

$$t_0 \sim \mathcal{U}(0.3, 1) \quad (18)$$

$$\alpha \sim \mathcal{U}(1, 2) \quad (19)$$

Figure S3 depicts different simulated RT distributions generated via draws from the prior.

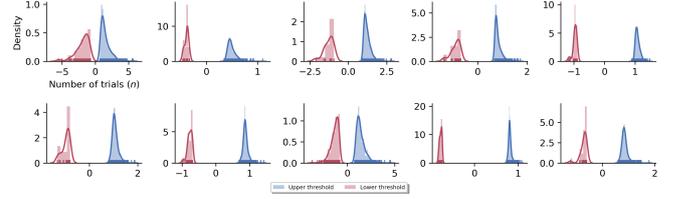


Fig. S3: Example RT distributions generated with different parameters.

The Stochastic SIR Model: Summary Network. We use a 1D fully convolutional neural network [29] for the raw SIR time-series into fixed-size vectors. Here, we choose a convolutional network architecture over the previously mentioned LSTM, as convolutional networks are more computationally efficient. Further, we wanted to underline the utility of 1D convolutional networks for multidimensional time-series data. Finally, convolutional networks can also deal with variable input sizes.

Simulation. We place the following uniform priors over the two rate parameters of the stochastic SIR model:

$$\beta \sim \mathcal{U}(0.01, 1) \quad (20)$$

$$\gamma \sim \mathcal{U}(0.01, \beta) \quad (21)$$

These ranges were chosen based on empirical plausibility of the generated SIR time-series.

Figure S4 depicts different SIR timeseries generated via draws from the prior.

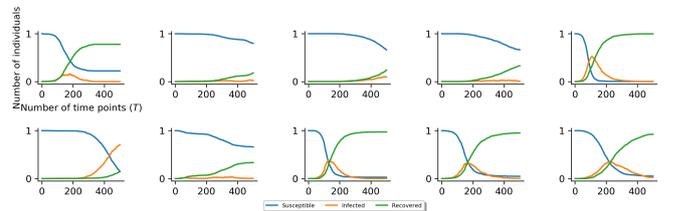


Fig. S4: Example SIR timeseries generated with different parameters.

The Lotka-Volterra Model: Summary Network. We use a bidirectional long short-term memory (LSTM) recurrent neural network [15] for the raw LV time-series (as in the Ricker example).

Simulation. We place the following broad uniform priors over the LV parameters. Some of the parameter combinations produced divergent simulations, which we removed during online learning.

$$\alpha \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (22)$$

$$\beta \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (23)$$

$$\gamma \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (24)$$

$$\delta \sim \mathcal{U}(\exp(-2), \exp(2)) \quad (25)$$

D. Example Posteriors on Ricker Datasets

Marginal posteriors from ten validation datasets simulated from the Ricker model are depicted in Figure S5. We observe widely different posterior shapes, highlighting the importance of working with arbitrary posterior shapes.

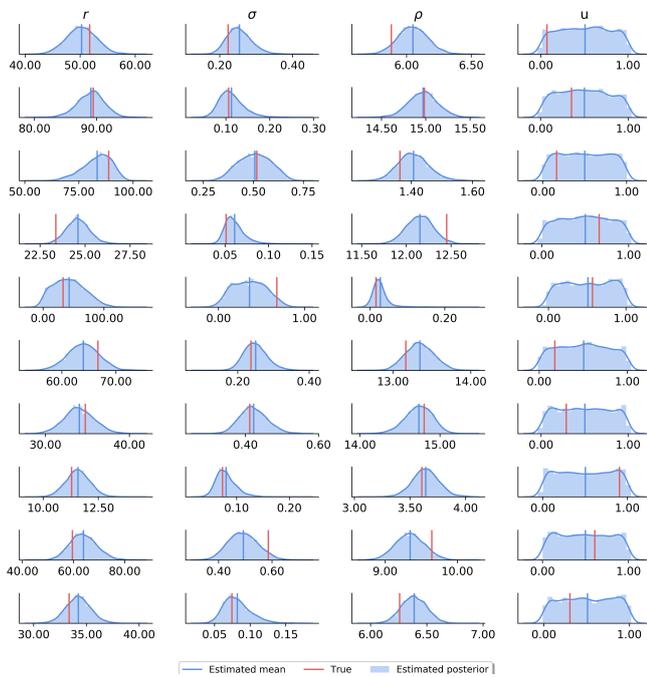


Fig. S5: Ten example Ricker marginal posteriors

APPENDIX A2 - MANUSCRIPT 2

Manuscript 2: Amortized Bayesian model comparison with evidential deep learning.

Amortized Bayesian Model Comparison with Evidential Deep Learning

Stefan T. Radev, Marco D’Alessandro, Ulf K. Mertens, Andreas Voss, Ullrich Köthe, Paul-Christian Bürkner

Abstract—Comparing competing mathematical models of complex natural processes is a shared goal among many branches of science. The Bayesian probabilistic framework offers a principled way to perform model comparison and extract useful metrics for guiding decisions. However, many interesting models are intractable with standard Bayesian methods, as they lack a closed-form likelihood function or the likelihood is computationally too expensive to evaluate. With this work, we propose a novel method for performing Bayesian model comparison using specialized deep learning architectures. Our method is purely simulation-based and circumvents the step of explicitly fitting all alternative models under consideration to each observed dataset. Moreover, it requires no hand-crafted summary statistics of the data and is designed to amortize the cost of simulation over multiple models, datasets, and dataset sizes. This makes the method especially effective in scenarios where model fit needs to be assessed for a large number of datasets, so that case-based inference is practically infeasible. Finally, we propose a novel way to measure epistemic uncertainty in model comparison problems. We demonstrate the utility of our method on toy examples and simulated data from non-trivial models from cognitive science and single-cell neuroscience. We show that our method achieves excellent results in terms of accuracy, calibration, and efficiency across the examples considered in this work. We argue that our framework can enhance and enrich model-based analysis and inference in many fields dealing with computational models of natural processes. We further argue that the proposed measure of epistemic uncertainty provides a unique proxy to quantify absolute evidence even in a framework which assumes that the true data-generating model is within a finite set of candidate models.

I. INTRODUCTION

Researchers from various scientific fields face the problem of selecting the most plausible theory for an empirical phenomenon among multiple alternative theories. These theories are often formally stated as mathematical models which describe how observable quantities arise from unobservable (latent) parameters. Focusing on the level of mathematical models, the problem of theory selection then becomes one of *model selection*.

For instance, neuroscientists might be interested in comparing different models of spiking patterns given *in vivo* recordings of neural activity [24]. Epidemiologists, on the other hand, might consider different models for predicting the spread and dynamics of an unfolding infectious disease [57]. Crucially, the preference for one model over alternative models in these examples can have important consequences for research projects or social policies.

Accounting for complex natural phenomena often requires specifying complex models which entail some degree of randomness. Inherent stochasticity, incomplete description, or epistemic ignorance all call for some form of uncertainty

awareness. To make matters worse, empirical data on which models are fit are necessarily finite and can only be acquired with finite precision. Finally, the plausibility of many non-trivial models throughout various branches of science can be assessed only approximately, through expensive simulation-based methods [46], [9], [53], [37], [24], [8].

Ideally, a method for approximate model comparison should meet the following desiderata:

- 1) *Theoretical guarantee*: Model probability estimates should be, at least in theory, calibrated to the true model probabilities induced by an empirical problem;
- 2) *Accurate approximation*: Model probability estimates should be accurate even for finite or small sample sizes;
- 3) *Occam’s razor*: Preference for simpler models should be expressed by the model probability estimates;
- 4) *Scalability*: The method should be applicable to complex models with implicit likelihood within reasonable time limits;
- 5) *Efficiency*: The method should enable fully amortized inference- over arbitrarily many models, datasets and different dataset sizes;
- 6) *Maximum data utilization*: The method should capitalize on all information contained in the data and avoid information loss through insufficient summary statistics of the data;

In the current work, we address these desiderata with a novel method for Bayesian model comparison based on evidential deep neural networks. Our method works in a purely simulation-based manner and circumvents the step of separately fitting all alternative models to each dataset. To this end, for any particular model comparison problem, we propose to train a *specialized expert network* which encodes global information about the generative scope of each model family. In this way, Bayesian model comparison is amortized over multiple models, datasets, and dataset sizes, which makes our method applicable in scenarios where case-based inference is way too costly to perform with standard methods (cf. Figure 1).

In addition, we propose to avoid hand-crafted summary statistics (a feature on which standard methods for simulation-based inference heavily rely) by utilizing novel deep learning architectures which are aligned to the probabilistic structure of the raw data (e.g., permutation invariant networks [3], recurrent networks [14]).

Finally, we explore a novel way to measure epistemic uncertainty in model comparison problems, following the pioneering work of [48] on image classification. We argue that this measure of epistemic uncertainty provides a unique proxy to quantify absolute evidence even in an \mathcal{M} -closed

framework, which assumes that the true data-generating model is within the candidate set [61].

II. BACKGROUND

A. Bayesian Inference

A consistent mathematical framework for describing uncertainty and quantifying model plausibility is offered by the Bayesian view on probability theory [25]. In a Bayesian setting, we start with a collection of J competing generative models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_J\}$. Each \mathcal{M}_j is associated with a generative mechanism g_j , typically realized as a Monte Carlo simulation program, and a corresponding parameter space Θ_j . Ideally, each g_j represents a theoretically plausible (potentially noisy) mechanism by which observable quantities \mathbf{x} arise from hidden parameters $\boldsymbol{\theta}$ and independent noise $\boldsymbol{\xi}$:

$$\mathbf{x} = g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi}) \text{ with } \boldsymbol{\theta}_j \in \Theta_j \quad (1)$$

where Θ_j is the corresponding parameter space of model g_j and the subscript j explicates that each model might be specified over a different parameter space. We assume that the functional or algorithmic form of each g_j is known and that we have a sample (dataset) $\{\mathbf{x}_i\}_{i=1}^N := \mathbf{x}_{1:N}$ of N (multivariate) observations $\mathbf{x}_n \in \mathcal{X}$ generated from an unknown process p^* . The task of Bayesian model selection is to choose the model in \mathcal{M} that best describes the observed data by balancing simplicity (sparsity) and predictive performance.

B. The Likelihood

A central object in Bayesian inference is the *likelihood function*, denoted as $p(\mathbf{x} | \boldsymbol{\theta}_j, \mathcal{M}_j)$. Loosely speaking, the likelihood returns the relative probability of an observation observation \mathbf{x} (or a sequence of observations $\mathbf{x}_{1:N}$) given a parameter configuration $\boldsymbol{\theta}_j$ and model assumptions \mathcal{M}_j . When the parameters are systematically varied and the data held constant, the likelihood can be used to quantify how well each model instantiation fits the data.

If the likelihood of a generative model can be associated with a known probability density function (e.g., Gaussian), the model can be formulated entirely in terms of the likelihood and the likelihood can be evaluated analytically or numerically for any pair $(\mathbf{x}, \boldsymbol{\theta})$. On the other hand, if the likelihood is unknown or intractable, as is the case when dealing with complex models, one can still generate random samples from the model by running the simulation program with a random configuration of its parameters.

This is due to the fact that each stochastic model, viewed as a Monte Carlo simulator, defines an implicit likelihood given by the relationship:

$$p(\mathbf{x} | \boldsymbol{\theta}_j, \mathcal{M}_j) = \int_{\Xi} \delta(\mathbf{x} - g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi})) p(\boldsymbol{\xi} | \boldsymbol{\theta}_j) d\boldsymbol{\xi} \quad (2)$$

where $\delta(\cdot)$ is the Dirac delta function and the integral runs over all possible execution paths of the stochastic simulation for a fixed $\boldsymbol{\theta}_j$. For most complex models, this integral is analytically intractable or too expensive to approximate numerically, so it is much easier to specify the model directly in terms of the simulation program g_j instead of deriving the likelihood

$p(\mathbf{x} | \boldsymbol{\theta}_j, \mathcal{M}_j)$. Importantly, we can still *sample* from the likelihood by running the simulator with different Monte Carlo realizations of $\boldsymbol{\xi}$, that is, for a fixed $\boldsymbol{\theta}_j$, we have the following equivalence:

$$\mathbf{x}_n \sim p(\mathbf{x} | \boldsymbol{\theta}_j, \mathcal{M}_j) \iff \mathbf{x}_n = g_j(\boldsymbol{\theta}_j, \boldsymbol{\xi}_n) \text{ with } \boldsymbol{\xi}_n \sim p(\boldsymbol{\xi}) \quad (3)$$

C. Bayes Factors

How does one assign preferences to competing models using a Bayesian toolkit? The canonical measure of evidence for a given model is the *marginal likelihood*:

$$p(\mathbf{x}_{1:N} | \mathcal{M}_j) = \int_{\Theta_j} p(\mathbf{x}_{1:N} | \boldsymbol{\theta}_j, \mathcal{M}_j) p(\boldsymbol{\theta}_j | \mathcal{M}_j) d\boldsymbol{\theta}_j \quad (4)$$

which is, in general, intractable to compute for non-trivial models. Importantly, the dependence on the prior over model \mathcal{M}_j 's parameters introduces a probabilistic version of Occam's razor, which expresses our preference for a simpler model over a more complex one when both models can account for the data equally well. The marginal likelihood thus focuses on *prior predictions* and penalizes the prior complexity of a model (i.e., the prior acts as a weight on the likelihood). This is in contrast to *posterior predictions*, which require marginalization over the parameter posterior $p(\boldsymbol{\theta}_j | \mathbf{x}_{1:N}, \mathcal{M}_j)$ and can be used to select the model which best predicts new data.

Provided that the marginal likelihood can be efficiently approximated, one can compute the ratio of marginal likelihoods for two models \mathcal{M}_j and \mathcal{M}_k via

$$\text{BF}_{jk} = \frac{p(\mathbf{x}_{1:N} | \mathcal{M}_j)}{p(\mathbf{x}_{1:N} | \mathcal{M}_k)}. \quad (5)$$

This famous ratio is called a Bayes factor (BF) and is used in Bayesian settings for quantifying relative model preference. Thus, a $\text{BF}_{jk} > 1$ indicates preference for model j over model k , given a set of observations $\mathbf{x}_{1:N}$. Alternatively, one can directly focus on the (marginal) posterior probability of a model \mathcal{M}_j ,

$$p(\mathcal{M}_j | \mathbf{x}_{1:N}) \propto p(\mathbf{x}_{1:N} | \mathcal{M}_j) p(\mathcal{M}_j) \quad (6)$$

which equips the model space itself with a prior distribution $p(\mathcal{M})$ over the considered model space encoding potential preferences for certain models before collecting any data. Such a prior might be useful if a model embodies extraordinary claims (e.g., telekinesis) and thus requires extraordinary evidence supporting it. However, if no prior reasons can be given for favoring some models over others (i.e., one prefers not to prefer), a uniform model prior $p(\mathcal{M}) = 1/J$ can be assumed.

The ratio of posterior model probabilities is called the *posterior odds* and is connected to the Bayes factor via the corresponding model priors:

$$\frac{p(\mathcal{M}_j | \mathbf{x}_{1:N})}{p(\mathcal{M}_k | \mathbf{x}_{1:N})} = \frac{p(\mathbf{x}_{1:N} | \mathcal{M}_j)}{p(\mathbf{x}_{1:N} | \mathcal{M}_k)} \times \frac{p(\mathcal{M}_j)}{p(\mathcal{M}_k)} \quad (7)$$

If two models are equally likely *a priori*, the posterior odds equal the Bayes factor. In this case, if the Bayes factor, or, equivalently, the posterior odds equal one, the observed data provide no decisive evidence for one of the models over the

other. However, a relative evidence of one does not allow to distinguish whether the data are equally likely or equally unlikely under both models, as this is a question of absolute evidence. Needless to say, the distinction between relative and absolute evidence is of paramount importance for model comparison, so we address it in the next section on model comparison frameworks.

D. \mathcal{M} -Frameworks

In Bayesian inference, the relationship between the true generative process p^* and the model list \mathcal{M} can be classified into three categories: \mathcal{M} -closed, \mathcal{M} -complete and \mathcal{M} -open [61]. Closely related to the distinction between relative and absolute evidence is the distinction between \mathcal{M} -closed and \mathcal{M} -complete frameworks. Under an \mathcal{M} -closed framework, the true model is assumed to be in the predefined set of competing models \mathcal{M} , so relative evidence is identical to absolute evidence. Under an \mathcal{M} -complete framework, a true model is assumed to exist but is not necessarily assumed to be a member of \mathcal{M} . However, one still focuses on the models in \mathcal{M} due to computational or conceptual limitations¹.

Deciding on the particular \mathcal{M} -framework under which a model comparison problem is tackled is often a matter of prior theoretical considerations. However, since in most non-trivial research scenarios \mathcal{M} is a finite set and candidate models in \mathcal{M} are often simpler approximations to the true model, there will be *uncertainty* as to whether the observed data could have been generated by one of these models. In the following, we will refer to this uncertainty as *epistemic uncertainty*. Our method utilizes a data-driven way to calibrate its epistemic uncertainty in addition to the model probabilities through simulations under an \mathcal{M} -closed framework.

Consequently, given real observed data, a researcher can obtain a measure of uncertainty with regard to whether the generative model of the data is likely to be in \mathcal{M} or not. From this perspective, our method lies somewhere between an \mathcal{M} -closed and an \mathcal{M} -complete framework as it provides information from both viewpoints. In this way, our approach to model misspecification differs from *likelihood-tempering* methods, which require an explicit evaluation of a *tilted* likelihood (raised to a power $0 < t < 1$) in order to prevent overconfident Bayesian updates [18].

III. RELATED WORK

Bayesian methods for model comparison can be categorized as either posterior predictive or prior predictive approaches [12], with our method falling into the latter category. Posterior predictive approaches are concerned with predicting new data using models trained on the current data. In prior predictive approaches, models are conditioned only on prior information but not on the current data. Accordingly, all current data counts as new data for the purpose of prior predictive methods.

Naturally, cross-validation (CV) procedures are the main approach for posterior predictive comparisons [56]. Examples

for widely applied methods that fall into this category are approximate cross-validation procedures using Pareto-smoothed importance sampling [55], [6], information criterion approaches such as the widely applicable information criterion (WAIC; [59]), or stacking of posterior predictive distributions [61].

All of these methods require not only the ability to evaluate the likelihood of each model for each observation during parameter estimation, but also for new observations during prediction. What is more, if application of exact CV methods is required because approximations are insufficient or unavailable, models need to be estimated several times based on different datasets or subsets of the original dataset. This renders such methods practically infeasible when working with complex simulators for which estimating models even once is already very slow. Thus, even a single intractable model in the model set suffices to disproportionately increase the difficulty of performing model comparison.

In contrast, our proposed method circumvents explicit parameter estimation and focuses directly on the efficient approximation of Bayes factors (or posterior model probabilities). Moreover, it overcomes two major sources of intractability that stand in the way of Bayesian model comparison via Bayes factors: the likelihood (Eq.3) and the marginal likelihood (Eq.4).

When the likelihood can be computed in closed-form, sophisticated algorithms for efficiently approximating the (intractable) marginal likelihood have been proposed in the Bayesian universe, such as *bridge sampling* and *path sampling* [13], [17]. However, these methods still depend on the ability to evaluate the likelihood $p(x|\theta_j, \mathcal{M}_j)$ for each candidate model. If, in addition, the likelihood itself is intractable, as is the case with complex simulators, researchers need to resort to expensive simulation-based methods [52], [53], [37], [40].

A standard set of tools for Bayesian simulation-based inference is offered by approximate Bayesian computation (ABC) methods [50], [38]. ABC methods approximate the model posterior by repeatedly sampling parameters from each proposal (prior) distribution and then simulating multiple datasets by running each simulator with the sampled parameters. A pre-defined similarity criterion determines whether a simulated dataset (or a summary statistic thereof) is sufficiently similar to the actually observed dataset. The model that most frequently generates synthetic observations matching those in the observed dataset is the one favored by ABC model comparison.

Despite being simple and elegant, standard ABC methods involve a crucial trade-off between accuracy and efficiency. In other words, stricter similarity criteria yield more accurate approximations of the desired posteriors at the price of higher and oftentimes intolerable rejection rates. What is more, most ABC methods require multiple *ad hoc* decisions from the method designer, such as the choice of similarity criterion or the summary statistics of the data (e.g., moments of empirical distributions) [37]. However, there is no guarantee that hand-crafted summaries extract all relevant information and model comparison with insufficient summary statistics can dramatically deteriorate the resulting model posteriors [47]. More scalable developments from the ABC family (ABC-SMC, ABC-MCMC, ABC neural networks and the recently proposed ABC random forests) offer great efficiency boosts but still rely

¹In this work, we delegate the discussion of whether the concept of a true model has any ontological meaning to philosophy. See also [61] for discussion of an \mathcal{M} -open framework, in which no true model is assumed to exist.

on hand-crafted summary statistics [37], [26], [49].

Recently, a number of promising innovations from the machine learning and deep learning literature have entered the field of simulation-based inference [7]. For instance, the sequential neural likelihood (SNL, [40]), the automatic posterior transformation (APT, [16]), the amortized ratio estimation [20] or the BayesFlow method [44] all implement powerful neural density estimators to overcome the shortcomings of standard ABC methods. Moreover, these methods involve some degree of *amortization*, which ensures extremely efficient inference after a potentially costly upfront training phase. However, neural density estimation focuses solely on efficient Bayesian parameter estimation instead of scaling up Bayesian model comparison. With certain caveats, neural density estimators can be adapted for Bayesian model comparison by post-processing the samples from an approximate posterior/likelihood over each model’s parameters. However, such an approach will involve training a separate neural estimator for each model in the candidate set and has not yet been systematically investigated. In addition, most of these methods also rely on fixed summary statistics [40]) and few applications using raw data directly exist [44], [16].

Alongside advancements in simulation-based inference, there has been an upsurge in the development of methods for uncertainty quantification in deep learning applications. For instance, much work has been done on the efficient estimation of Bayesian neural networks [21], [33], [35] since the pioneering work of [36]. Parallel to the establishment of novel variational methods [30], [29], these ideas have paved the way towards more interpretable and trustworthy neural network inference. Moreover, the need for distinguishing between different sources of uncertainty and the overconfidence of deep neural networks in classification and regression tasks has been demonstrated quite effectively [28], [48]. Our current work draws on recent methods for evidence and uncertainty representation in classification tasks [48]. However, our goal is to efficiently approximate Bayes factors between competing mechanistic models using non-Bayesian neural networks, not to estimate neural network parameters (e.g., weights) via Bayesian methods.

Our method combines latest ideas from simulation-based inference and uncertainty quantification for building efficient and uncertainty-aware estimators for amortized Bayesian model comparison. As such, it is intended to complement the toolbox of simulation-based methods for parameter estimation with crucial model comparison capabilities and incorporates some unique features beyond the scope of standard ABC methods. In the following, we describe the building blocks of our method.

IV. EVIDENTIAL NETWORKS FOR BAYESIAN MODEL COMPARISON

A. Model Selection as Classification

In line with previous simulation-based approaches to model selection, we will utilize the fact that we can generate arbitrary amounts of data via Eq.3 for each considered model \mathcal{M}_j . Following [43], [37], we cast the problem of model comparison as a probabilistic classification task. In other words, we seek

a parametric mapping $f_\phi : \mathcal{X}^N \rightarrow \Delta^J$ from an arbitrary data space \mathcal{X}^N to a probability simplex Δ^J containing the posterior model probabilities $p(\mathcal{M} | \mathbf{x}_{1:N})$. Previously, different learning algorithms (e.g., random forests [37]) have been employed to tackle model comparison as classification. Following recent developments in algorithmic alignment and probabilistic symmetry [60], [3], our method parameterizes f_ϕ via a specialized neural network with trainable parameters ϕ which is aligned to the probabilistic structure of the observed data. See the **Network Architectures** section in **Appendix A** for a detailed description of the employed networks’ structure.

In addition, our method differs from previous classification approaches to model comparison in the following aspects. First, it requires no hand-crafted summary statistics, since the most informative summary statistics are learned directly from data. Second, it uses online learning (i.e., on-the-fly simulations) and requires no storage of large reference tables or data grids. Third, the addition of new competing models does not require changing the architecture or re-training the network from scratch, since the underlying data domain remains the same. In line with the transfer learning literature, only the last layer of a pre-trained network needs to be changed and training can be resumed from where it had stopped. Last, our method is uncertainty-aware, as it returns a higher-order distribution over posterior model probabilities. From this distribution, one can extract both absolute and relative evidences, as well as quantify the model selection uncertainty implied by the observed data.

To set up the model classification task, we run **Algorithm 1** repeatedly to construct training batches with B simulated datasets of size N and B model indices of the form $\mathcal{D}_N^{(B)} := \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$. We then feed each batch to a neural network which takes as input simulated data with variable sizes and returns a distribution over posterior model probabilities. The neural network parameters are optimized via standard backpropagation. Upon convergence, we can apply the pre-trained network to arbitrarily many datasets of the form $\mathbf{x}_{1:N}^{(obs)}$ to obtain a vector of probabilities $p_\phi(\mathbf{m} | \mathbf{x}_{1:N}^{(obs)})$ which approximates the true model posterior $p(\mathcal{M} | \mathbf{x}_{1:N}^{(obs)})$.

Note, that this procedure incurs no memory overhead, as the training batches need not be stored in memory all at once. Intuitively, the connection between data and models is encoded in the network’s weights. Once trained, the evidential network can be reused to perform instant model selection on multiple real observations. As mentioned above, the addition of new models requires simply adjusting the pre-trained network, which requires much less time than re-training the network from scratch. We now describe how model probabilities and evidences are represented by the evidential network.

B. Evidence Representation

In order to obtain a measure of absolute evidence by considering a finite number of competing models, we place a Dirichlet distribution over the estimated posterior model probabilities [48]. This corresponds to modeling second-order probabilities in terms of the theory of subjective logic (SL) [27]. These second-order probabilities represent an uncertainty measure over quantities which are themselves probabilities. We use

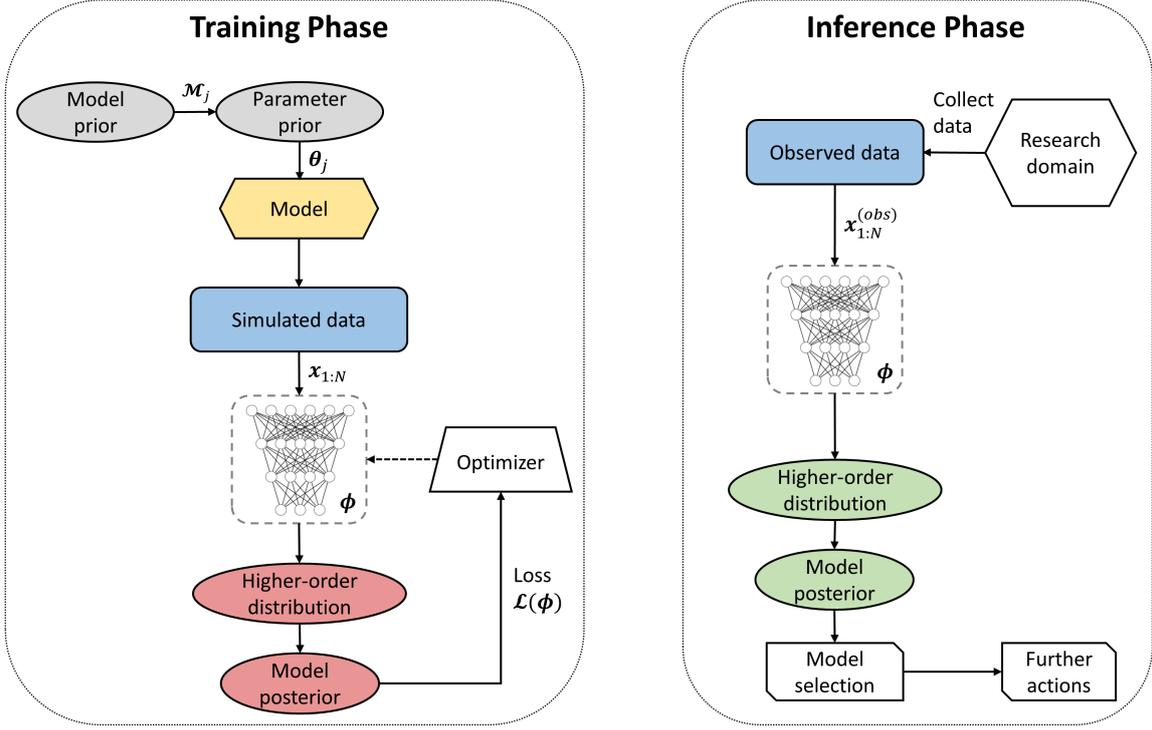


Fig. 1: Left panel: The simulation-based training phase of our evidential method. Right panel: The inference phase with real data and a pre-trained evidential network.

Algorithm 1 Monte Carlo generation of synthetic datasets for model comparison

Require: $p(\mathcal{M})$ - prior over models, $\{p(\theta_j | \mathcal{M}_j)\}$ - list of priors over model parameters, $\{g_j\}$ - list of stochastic simulators, $\{p_j(\xi)\}$ - list of noise distributions (RNGs), $p(N)$ - distribution over dataset sizes, B - number of datasets to generate (batch size)

- 1: Draw dataset size: $N \sim p(N)$
 - 2: **for** $b = 1, \dots, B$ **do**
 - 3: Draw model index from model prior: $\mathcal{M}_j^{(b)} \sim p(\mathcal{M})$
 - 4: Draw model parameters from prior: $\theta_j^{(b)} \sim p(\theta_j | \mathcal{M}_j^{(b)})$
 - 5: **for** $n = 1, \dots, N$ **do**
 - 6: Sample noise instance: $\xi_n \sim p_j(\xi)$
 - 7: Run simulator j to obtain n -th synthetic observation: $x_n = g_j(\theta_j^{(b)}, \xi_n)$
 - 8: **end for**
 - 9: Encode model index as a one-hot-encoded vector: $\mathbf{m}^{(b)} = \text{OneHotEncode}(\mathcal{M}_j^{(b)})$
 - 10: Store pair $(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})$ in $\mathcal{D}_N^{(B)}$
 - 11: **end for**
 - 12: **return** mini-batch $\mathcal{D}_N^{(B)} := \{\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)}\}_{b=1}^B$
-

the second-order probabilities to capture epistemic uncertainty about whether the observed data has been generated by one of the candidate models considered during training.

The probability density function (PDF) of a Dirichlet distribution is given by:

$$\text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{j=1}^J \pi_j^{\alpha_j - 1} \quad (8)$$

where $\boldsymbol{\pi}$ belongs to the unit $J - 1$ simplex (i.e., $\boldsymbol{\pi} \in \Delta^J := \{\boldsymbol{\pi} | \sum_{j=1}^J \pi_j = 1\}$) and $B(\boldsymbol{\alpha})$ is the multivariate beta function. The Dirichlet density is parameterized by a

vector of *concentration parameters* $\boldsymbol{\alpha} \in \mathbb{R}_+^J$ which can be interpreted as evidences in the ST framework [27]. The sum of the individual evidence components $\alpha_0 = \sum_{j=1}^J \alpha_j$ is referred to as the Dirichlet strength, and it affects the precision of the higher-order distribution in terms of its variance. Intuitively, the Dirichlet strength governs the *peakedness* of the distribution, with larger values leading to more peaked densities (i.e., most of the density being concentrated in a smaller region of the simplex). We can use the mean of the Dirichlet distribution,

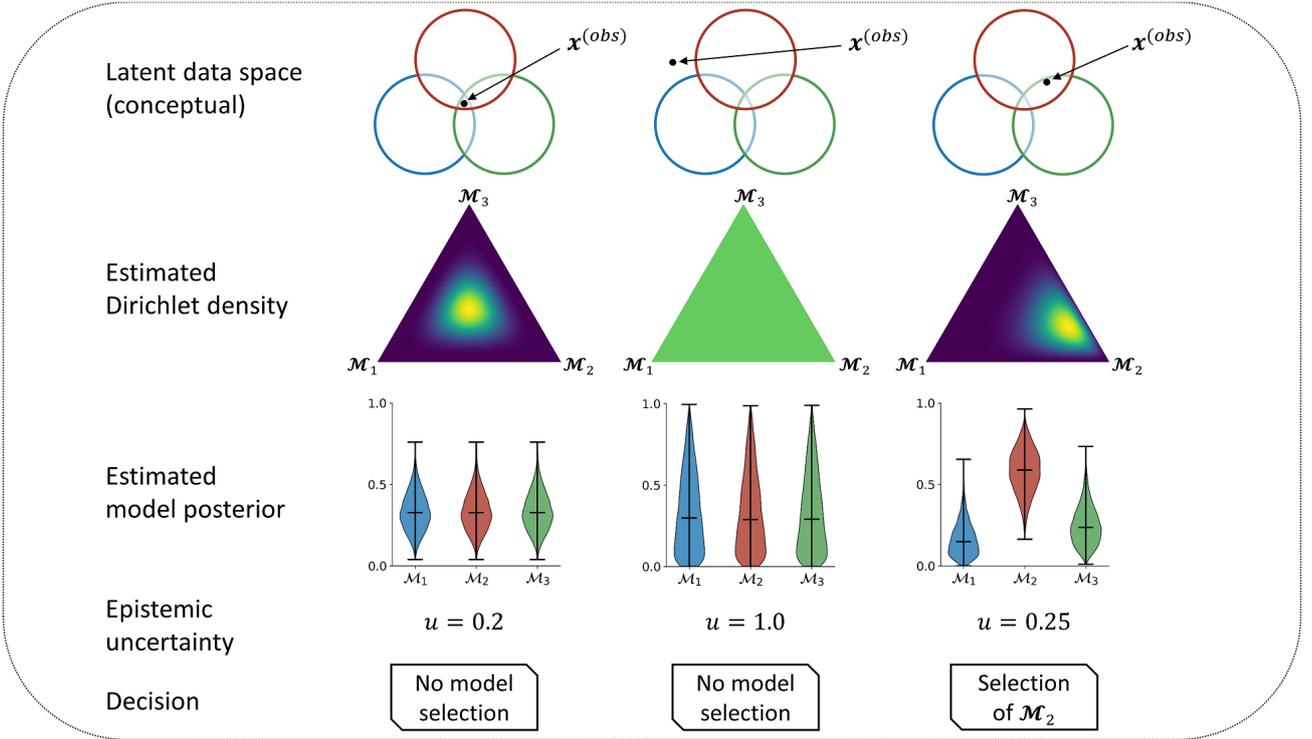


Fig. 2: Three different hypothetical model comparison scenarios with different observations. The first column depicts observing a dataset which is equally probable under all models. In this case, the best candidate model cannot be selected and the Dirichlet density peaks in the middle of the simplex. The second column depicts a dataset which is beyond the generative scope of all models and no model selection decision is possible. The Dirichlet density in this case is flat which indicates total uncertainty. The third column illustrates an observed dataset which is most probable under model 2, so the Dirichlet simplex is peaked towards the corner encoding model 2, and the corresponding model posterior for model 2 is highest.

which is a vector of probabilities given by:

$$\mathbb{E}_{\pi \sim \text{Dir}(\alpha)}[\pi] = \frac{\alpha}{\alpha_0} \quad (9)$$

to approximate the posterior model probabilities $p(\mathcal{M} | \mathbf{x}_{1:N})$, as will become clearer later in this section. A crucial advantage of such a Dirichlet representation is that it allows to look beyond model probabilities by inspecting the vector of computed evidences. For instance, imagine a scenario with three possible models. If $\alpha = (5, 5, 5)$, the data provides equally strong evidence for all models (Figure 2, first column) – all models explain the data well. If, on the other hand, $\alpha = (1, 1, 1)$, then the Dirichlet distribution reduces to a uniform on the simplex indicating no evidence for any of the models (Figure 2, second column) – no model explains the observations well. Note that in either case one cannot select a model on the basis of the data, because posterior model probabilities are equal, yet the interpretation of the two outcomes is very different: The second-order Dirichlet distribution allows one to distinguish between *equally likely* (first case) and *equally unlikely* (second case) models. The last column of Figure 2 illustrates a scenario with $\alpha = (2, 7, 3)$ in which case one can distinguish between all models (see also Figure 6 for a scenario with data simulated from an actual model).

We can further quantify this distinction by computing an uncertainty score given by:

$$u = \frac{J}{\alpha_0} \quad (10)$$

where J is the number of candidate models. Importantly, in our framework, individual concentration parameters (resp. neural network outputs) are lower bounded by 1. Thus, the uncertainty score ranges between 0 (total certainty) and 1 (total uncertainty) and has a straightforward interpretation. Accordingly, total uncertainty is given when $\alpha_0 = J$, which would mean that the data provide no evidence for any of the J candidate models. On the other hand, $u \ll 1$ implies a large Dirichlet strength $\alpha_0 \gg J$, which would read that the data provide plenty of evidence for one or more models in question. The uncertainty score corresponds to the concept of *vacuity* (i.e., epistemic uncertainty) in the terminology of SL [27]. We argue that epistemic uncertainty should be a crucial aspect in model selection, as it quantifies the strength of evidence, and, consequently, the strength of the theoretical conclusions we can draw given the observed data.

Consequently, model comparison in our framework consists in inferring the parameters of a Dirichlet distribution given an

observed dataset. The problem of inferring posterior model probabilities can be formulated as:

$$p(\mathcal{M} | \mathbf{x}_{1:N}) \approx p_\phi(\mathbf{m} | \mathbf{x}_{1:N}) = \mathbb{E}_{\boldsymbol{\pi} \sim \text{Dir}(f_\phi(\mathbf{x}_{1:N}))}[\boldsymbol{\pi}] \quad (11)$$

where f_ϕ is a neural network with positive outputs greater than one, $f_\phi : \mathcal{X}^N \rightarrow [1, \infty]^J$. Additionally, we can also obtain a measure of absolute model evidence by considering the uncertainty encoded by the full Dirichlet distribution (Eq.10).

C. Learning Evidence in an \mathcal{M} -Closed Framework

How do we ensure that the outputs of the neural network match the true unknown model posterior probabilities? Consider, for illustrational purposes, a dataset with a single observation, that is $N = 1$ such that $\mathbf{x}_{1:N} = \mathbf{x}$. As per **Algorithm 1**, we have unlimited access to training samples from $p(\mathcal{M}, \mathbf{x}) = \int p(\mathcal{M})p(\boldsymbol{\theta} | \mathcal{M})p(\mathbf{x} | \boldsymbol{\theta}, \mathcal{M})d\boldsymbol{\theta}$. We use the mean of the Dirichlet distribution $p_\phi(\mathbf{m} | \mathbf{x})$ parameterized by an evidential neural network with parameters ϕ to approximate $p(\mathcal{M} | \mathbf{x})$. To optimize the parameters of the neural network, we can minimize some loss \mathcal{L} in expectation over all possible datasets:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{(\mathbf{m}, \mathbf{x}) \sim p(\mathcal{M}, \mathbf{x})} [\mathcal{L}(p_\phi(\mathbf{m} | \mathbf{x}), \mathbf{m})] \quad (12)$$

where \mathbf{m} is a one-hot encoded vector of the true model index \mathcal{M}_j . We also require that \mathcal{L} be a *strictly proper loss* [15]. A loss function is strictly proper if and only if it attains its minimum when $p_\phi(\mathbf{m} | \mathbf{x}) = p(\mathcal{M} | \mathbf{x})$ [15]. When we choose the Shannon entropy $\mathbb{H}(p_\phi(\mathbf{m} | \mathbf{x})) = -\sum_j p_\phi(\mathbf{m} | \mathbf{x})_j \log p_\phi(\mathbf{m} | \mathbf{x})_j$ for \mathcal{L} , we obtain the strictly proper logarithmic loss:

$$\mathcal{L}(p_\phi(\mathbf{m} | \mathbf{x}), \mathbf{m}) = -\sum_{j=1}^J m_j \log p_\phi(\mathbf{m} | \mathbf{x})_j \quad (13)$$

$$= -\sum_{j=1}^J m_j \log \left(\frac{f_\phi(\mathbf{x})_j}{\sum_{j'=1}^J f_\phi(\mathbf{x})_{j'}} \right) \quad (14)$$

where $m_j = 1$ when j is the true model index and 0 otherwise. Thus, in order to estimate ϕ , we can minimize the expected logarithmic loss over all simulated datasets where $f_\phi(\mathbf{x})_j$ denotes the j -th component of the Dirichlet density given by the evidential neural network. Since we use a strictly proper loss, the evidential network yields the true model posterior probabilities over all possible datasets when perfectly converged.

Intuitively, the logarithmic loss encourages high evidence for the true model and low evidences for the alternative models. Correspondingly, if a dataset with certain characteristics can be generated by different models, evidence for these models will jointly increase. Additionally, the model which generates these characteristics most frequently will accumulate the most evidence and thus be preferred. However, we also want low evidence, or, equivalently, high epistemic uncertainty, for datasets which are implausible under all models. We address this problem in the next section.

D. Learning Absolute Evidence through Regularization

We now propose a way to address the scenario in which no model explains the observed data well. In this case, we want the evidential network to estimate low evidence for all models in the candidate set. In order to attenuate evidence for datasets which are implausible under all models considered, we incorporate a Kullback-Leibler (KL) divergence into the criterion in Eq.13. We compute the KL divergence between the Dirichlet density generated by the neural network and a uniform Dirichlet density implying total uncertainty. Thus, the KL shrinks evidences which do not contribute to correct model assignments during training, so an implausible dataset at inference time will lead to low evidence under all models. This type of regularization has been used for capturing out-of-distribution (OOD) uncertainty in image classification [48]. Accordingly, our modified optimization criterion becomes:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \mathbb{E}_{(\mathbf{m}, \mathbf{x}) \sim p(\mathcal{M}, \mathbf{x})} [\mathcal{L}(p_\phi(\mathbf{m} | \mathbf{x}), \mathbf{m}) + \lambda \Omega(\tilde{\boldsymbol{\alpha}})] \quad (15)$$

with $\Omega(\tilde{\boldsymbol{\alpha}}) = \mathbb{KL}[\text{Dir}(\tilde{\boldsymbol{\alpha}}) || \text{Dir}(\mathbf{1})]$. The term $\tilde{\boldsymbol{\alpha}} = \mathbf{m} + (1 - \mathbf{m}) \odot \boldsymbol{\alpha}$ represents the estimated evidence vector after removing the evidence for the true model. This is possible, because we know the true model during simulation-based training. For application on real datasets after training, knowing the ground truth is not required anymore as ϕ has been obtained already at this point. The KL penalizes evidences for the false models and drives these evidences towards unity. Equivalently, the KL acts as a ground-truth preserving prior on the higher-order Dirichlet distribution which preserves evidence for the true model and attenuates misleading evidences for the false models. The hyperparameter λ controls the weight of regularization and encodes the tolerance of the algorithm to accept implausible (out-of-distribution) datasets during inference. With large values of λ , it becomes possible to detect cases where all models are deficient; with $\lambda = 0$, only relative evidence is generated. Note, that in the latter case, we recover our original proper criterion without penalization. The KL weight λ can be selected through prior empirical considerations on how well the simulations cover the plausible set of real-world datasets.

Importantly, the introduction of the KL regularizer renders the loss no longer *strictly proper*. Therefore, a large regularization weight λ would lead to poorer calibration of the approximate model posteriors, as the regularized loss is no longer minimized by the true model posterior. However, since the KL prior is ground-truth preserving, the accuracy of recovering the true model should not be affected. Indeed, we observe this behavior throughout our experiments.

To make optimization tractable, we utilize the fact that we can easily simulate batches of the form $\mathcal{D}_N^{(B)} = \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ via **Algorithm 1** and approximate Eq.15 via standard backpropagation by minimizing the following loss:

$$\mathcal{L}(\phi) = \frac{1}{B} \sum_{b=1}^B \left[-\sum_{j=1}^J m_j^{(b)} \log \left(\frac{f_\phi(\mathbf{x}_{1:N}^{(b)})_j}{\sum_{j'=1}^J f_\phi(\mathbf{x}_{1:N}^{(b)})_{j'}} \right) + \lambda \Omega(\tilde{\boldsymbol{\alpha}}^{(b)}) \right] \quad (16)$$

Algorithm 2 Training phase and inference phase for amortized Bayesian model comparison

Require: f_ϕ - evidential neural network, $\{\mathbf{x}_{1:N_i}^{(obs)}\}_{i=1}^I$ - list of I observed datasets for inference, λ - regularization weight, B - number of simulations at each iteration (batch size)

- 1: *Training phase:*
 - 2: **repeat**
 - 3: Generate a training batch $\mathcal{D}_N^{(B)} = \{(\mathbf{m}^{(b)}, \mathbf{x}_{1:N}^{(b)})\}_{b=1}^B$ via **Algorithm 1**
 - 4: Compute evidences for each simulated dataset in $\mathcal{D}_N^{(B)}$: $\alpha^{(b)} = f_\phi(\mathbf{x}_{1:N}^{(b)})$
 - 5: Compute loss according to Eq.16
 - 6: Update neural network parameters ϕ via backpropagation
 - 7: **until** convergence to $\hat{\phi}$
 - 8: *Amortized inference phase:*
 - 9: **for** $i = 1, \dots, I$ **do**
 - 10: Compute model evidences $\alpha_i^{(obs)} = f_{\hat{\phi}}(\mathbf{x}_{1:N_i}^{(obs)})$
 - 11: Compute uncertainty $u_i = J / \sum_{j=1}^J \alpha_{i,j}^{(obs)}$
 - 12: Approximate true model posterior probabilities $p(\mathcal{M} | \mathbf{x}_{1:N_i}^{(obs)})$ via $p_\phi(\mathbf{m} | \mathbf{x}_{1:N_i}) = \alpha_i^{(obs)} / \sum_{j=1}^J \alpha_{i,j}^{(obs)}$
 - 13: **end for**
 - 14: Choose further actions
-

over multiple batches to converge at a Monte Carlo estimator $\hat{\phi}$ of ϕ^* . In practice, convergence can be determined as the point at which the loss stops decreasing, a criterion similar to *early stopping*. Alternatively, the network can be trained for a pre-defined number of epochs. Note, that, at least in principle, the network can be trained arbitrarily long, since we assume that we can access the joint distribution $p(\mathcal{M}, \mathbf{x}, N)$ through simulation (cf. Figure 1, left panel).

E. Implicit Preference for Simpler Models

Remembering that $p_\phi(\mathbf{m} | \mathbf{x}_{1:N}) \propto p(\mathbf{x}_{1:N} | \mathcal{M})p(\mathcal{M})$, we note that perfect convergence implies that preference for simpler models (Bayesian Occam’s razor) is automatically encoded by our method. This is due to the fact that we are approximating an expectation over all possible datasets, parameters, and models. Accordingly, datasets generated by a simpler model tend to be more similar compared to those from a more complex competitor. Therefore, during training, certain datasets which are plausible under both models will be generated more often by the simpler model than by the complex model. Thus, a perfectly converged evidential network will capture this behavior by assigning higher posterior probability to the simpler model (assuming equal prior probabilities). Therefore, at least in theory, our method captures complexity differences arising purely from the generative behavior of the models and does not presuppose an *ad hoc* measure of complexity (e.g., number of parameters).

F. Putting it All Together

The essential steps of our evidential method are summarized in **Algorithm 2**. Note, that steps 2-7 and 9-13 can be executed in parallel with GPU support in order to dramatically accelerate convergence and inference. In sum, we propose to cast the problem of model selection as evidence estimation and learn a Dirichlet distribution over posterior model probabilities directly via simulations from the competing models. To this end, we

train an evidential neural network which approximates posterior model probabilities and further quantifies the epistemic uncertainty as to whether an observed dataset is within the generative scope of the candidate models. Moreover, once trained on simulations from a set of models, the network can be reused and extended with new models across a research domain, essentially *amortizing* the model comparison process. Accordingly, if the priors over model parameter do not change, multiple researchers can reuse the same network for multiple applications. If the priors over model parameters change or additional models need to be considered, the parameters of a pre-trained network can be adjusted or the network extended with additional output nodes for the new models.

V. EXPERIMENTS

In this section, we demonstrate the utility of our method on a toy example and relevant models from chemistry, cognitive science and neurobiology. A further toy example with 400 models as well as details for neural network training, architectures, performance metrics, and forward models are to be found in the **Appendix**.

A. Experiment 1: Beta-Binomial Model with Known Analytical Marginal Likelihood

As a basic proof-of-concept for our evidential method, we apply it to a toy model comparison scenario with an analytically tractable marginal likelihood. Here, we pursue the following goals. First, we want to demonstrate that the estimated posterior probabilities closely approximate the analytic model posteriors. To show this, we compare the analytically computed vs. the estimated Bayes factors. In addition, we want to show that accuracy of true model recovery matches closely the accuracy obtained by the analytic Bayes factors across all N . For this, we consider the simple beta-binomial model given by:

$$\theta \sim \text{Beta}(\alpha, \beta) \quad (17)$$

$$x_j \sim \text{Bernoulli}(\theta) \text{ for } j = 1, \dots, N \quad (18)$$

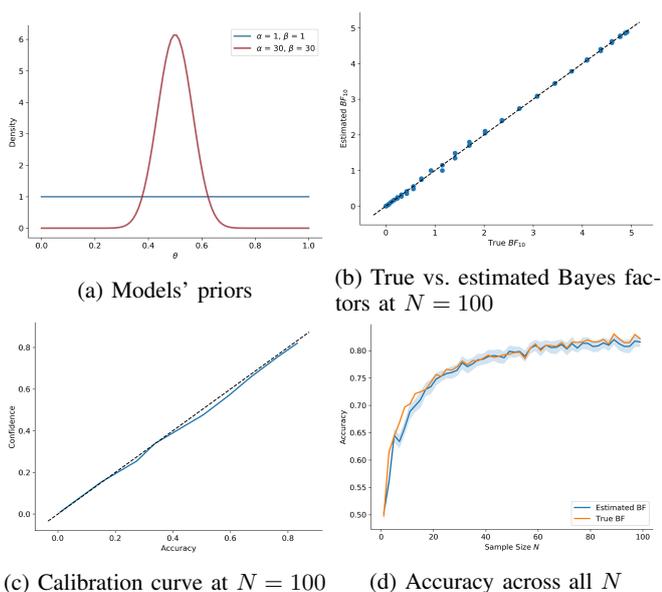


Fig. 3: **(a)** Prior densities of the theta parameter for both models of **Experiment 1**; **(b)** True vs. estimated Bayes factors obtained from the network-induced Dirichlet distribution at $N = 100$; **(c)** Calibration curve at $N = 100$ indicating very good calibration (dotted line represents perfect calibration); **(d)** Accuracy at all N achieved with both the analytic and the estimated Bayes factors (the shaded region represents a 95% bootstrap confidence interval around the accuracies of the evidential network).

The analytical marginal likelihood of the beta-binomial model is:

$$p(x_{1:N}) = \binom{N}{K} \frac{\text{Beta}(\alpha + K, \beta + N - K)}{\text{Beta}(\alpha, \beta)} \quad (19)$$

where K denotes the number of successes in the N trials. For this example, we will consider a model comparison scenario with two models, one with a flat prior $\text{Beta}(1, 1)$ on the parameter θ , and one with a sharp prior $\text{Beta}(30, 30)$. The two prior densities are depicted in Figure 3a.

We train a small permutation invariant evidential network with batches of size $B = 64$ until convergence. For each batch, we draw the samples size from a discrete uniform distribution $N \sim \mathcal{U}_D(1, 100)$ and input the raw binary data to the network. We validate the network on 5000 separate validation datasets for each N . Convergence took approximately 15 minutes, whereas inference on all 5000 validation datasets took less than 2 seconds.

Our results demonstrate that the estimated Bayes factors closely approximate the analytic Bayes factors (Figure 3b). We also observe no systematic under- or overconfidence in the estimated Bayes factors, which is indicated by the calibration curve resembling a straight line (Figure 3c). Finally, the accuracy of recovery achieved with the estimated Bayes factors closely matches that of the analytic Bayes factors across all sample sizes N (Figure 3d).

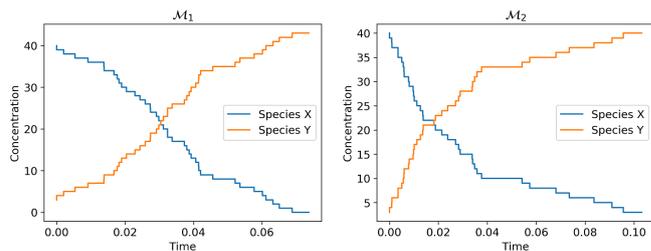


Fig. 4: Observed concentration time-series from both Markov jump models of **Experiment 2** with $\theta = 2.0$.

B. Experiment 2: Markov Jump Process of Stochastic Chemical Reaction Kinetics

In this experiment, we apply our evidential method to a simple model of non-exchangeable chemical molecule concentration data. Further, we demonstrate the efficiency benefits of our amortized learning method compared to the standard non-amortized ABC-SMC algorithm.

We define two Markov jump process models \mathcal{M}_1 and \mathcal{M}_2 for conversion of (chemical) species z to species y :



Each model has a single rate parameter θ_i . We use the Gillespie simulator to generate simulated time-series from the two models with an upper time limit of 0.1 seconds. Both models start with initial concentrations $x_0 = 40$ and $y_0 = 3$, so they only differ in terms of their reaction kinetics. The input time-series $x_{1:N}$ consist of a time vector $t_{1:N}$ as well as two vectors of molecule concentrations for each species at each time step, $z_{1:N}$ and $y_{1:N}$, which we stack together. We place a wide uniform prior over each rate parameter: $\theta_i \sim \mathcal{U}(0, 100)$.

We train an evidential sequence network for 50 epochs of 1000 mini-batch updates and validate its performance on 500 previously unobserved time-series. Wall-clock training time was approximately 52.3 minutes. In contrast, wall-clock inference time on the 500 validation time-series was 254 ms, leading to dramatic amortization gains. The bootstrap accuracy of recovery was 0.98 ($SD = 0.01$) over the entire validation set.

We also apply the ABC-SMC algorithm available from the *pyABC* [31] library to a single data-set $x_{1:N}^{(obs)}$ generated from model 1 (\mathcal{M}_1) with rate parameter $\theta_1 = 2.0$. Figure 4 depicts the observed data generated from model 1 (left panel) as well as observed data generated from model 2 with $\theta_2 = 2.0$. Notably, the differences in the data-generating processes defined by the two models are subtle and not straightforward to explicitly quantify.

For the ABC-SMC method, we set the minimum rejection threshold ϵ to 0.7 and the maximum number of populations to 15, as these settings lead to perfect recovery of the true model. As a distance function, we use the L_2 norm between the raw concentration times-series of species z , evaluated at 20 time points.²

²These settings were picked from the original *pyABC* documentation available at https://pyabc.readthedocs.io/en/latest/examples/chemical_reaction.html

The convergence of the ABC-SMC algorithm on the single dataset took 12.2 minutes wall-clock time. Thus, inference on the 500 validation data-sets would have taken more than 4 days to complete. Accordingly, we see that the training effort with our method is worthwhile even after as few as 5 data-sets. As for recovery on the single test dataset, ABC-SMC selects the true model with a probability of 1, whereas our evidential networks outputs a probability of 0.997 which results in a negligible difference of 0.003 between the results from two methods.

C. Experiment 3: Stochastic Models of Decision Making

In this experiment, we apply our evidential method to compare several non-trivial nested stochastic *evidence accumulator models* (EAMs) from the field of human decision making [54], [45]. With this experiment, we want to demonstrate the performance of our method in terms of accuracy and posterior calibration on exchangeable data obtained from complex cognitive models. Additionally, we want to demonstrate how our regularization scheme can be used to capture absolute evidence by artificially rendering the data implausible under all models.

1) *Model Comparison Setting*: EAMs describe the dynamics of decision making via different neurocognitively plausible parameters (i.e., speed of information processing, decision threshold, bias/pre-activation, etc.). EAMs are most often applied to choice reaction times (RT) data to infer neurocognitive processes underlying generation of RT distributions in cognitive tasks. The most general form of an EAM is given by a stochastic differential equation:

$$dx = vdt + cd\xi \quad (22)$$

where dx denotes a change in activation of an accumulator, v denotes the average speed of information accumulation (often termed the drift rate), and $d\xi$ represents a stochastic additive component with $d\xi \sim \mathcal{N}(0, c^2)$.

Multiple flavors of the above stated basic EAM form exist throughout the literature [54], [45], [53], [4]. Moreover, most EAMs are intractable with standard Bayesian methods [4], so model selection is usually hard and computationally cumbersome. With this example, we pursue several goals. First, we want to demonstrate the utility of our method for performing model selection on multiple nested models. Second, we want to empirically show that our method implements Occam’s razor. Third, we want to show that our method can indeed provide a proxy for absolute evidence.

To this end, we start with a very basic EAM defined by four parameters $\theta = (v_1, v_2, a, t_0)$ with v_i denoting the speed of information processing (drift rate) for two simulated RT experimental tasks $i \in \{1, 2\}$, a denoting the decision threshold, and t_0 denoting an additive constant representing the time required for non-decisional processes like motor reactions. We then define five more models with increasing complexity by successively *freeing* the parameters z_r (bias), α (heavy-tailness of noise distribution), s_{t_0} (variability of non-decision time), s_v (variability of drift-rate), and s_{z_r} (variability of bias). Note, that

the inclusion of non-Gaussian diffusion noise renders an EAM model intractable, since in this case no closed-form likelihood is available (see [58] for more details). Table S1 lists the priors over model parameters as well as fixed parameter values.

The task of model selection is thus to choose among six nested EAM models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3, \mathcal{M}_4, \mathcal{M}_5, \mathcal{M}_6\}$, each able to capture increasingly complex behavioral patterns. Each model j is able to account for all datasets generated by the previous models $i < j$, since the previous models are nested within the j -th model. For instance, model \mathcal{M}_6 can generate all datasets possible under the other models at the cost of increased functional and parametric complexity. Therefore, we need to show that our method encodes Occam’s razor purely through the generative behavior of the models.

In order to show that our regularization method can be used as a proxy to capture absolute evidence, we perform the following experiment. We define a temporal shifting constant $K \in (0, 10)$ (in units of seconds) and apply the shift to each response time in each validation dataset. Therefore, as K increases, each dataset becomes increasingly implausible under all models considered. For each K , we compute the average uncertainty over all shifted validation datasets and plot is as a function of K . Here, we only consider the maximum number of trials $N = 300$.

We train three evidential neural networks with different KL weights: $\lambda \in \{0.0, 0.1, 1.0\}$ in order to investigate the effects of λ on accuracy, calibration, and uncertainty. All networks were trained with variable number of trials $N \sim \mathcal{U}_D(1, 300)$ per batch for a total of 50000 iterations. The training of each network took approximately half a day on a single computer. In contrast, performing inference on 5000 datasets with a pre-trained network took less than 10 seconds.

2) *Validation Results*: To quantify the global performance of our method, we compute the accuracy of recovery as a function of the number of observations (N) for each of the models. We also compute the epistemic uncertainty as a function of N . To this end, we generate 500 new datasets for each N and compute the accuracy of recovery and average uncertainty. These results are depicted in Figure 5.

Accuracy. We observe that accuracy of recovery increases with increasing sample size and begins to flatten out around $N = 100$, independently of the regularization weight λ (Figure 5a). This behavior is desirable, as selecting the true model should become easier when more information is available. Further, since the models are nested, perfect recovery is not possible, as the models exhibit a large shared data space.

Calibration. Figure 5d depicts calibration curves for each model and each regularization value. The unregularized network appears to be very well calibrated, whereas the regularized networks become increasingly underconfident with increasing regularization weight. This is due to the fact that the regularized networks are encouraged to generate zero evidence for the wrong models, so model probabilities become miscalibrated. Importantly, none of the networks shows overconfidence.

Occam’s Razor. We also test Occam’s razor by generating 500 datasets from each model with $N = 300$ and compute the average predicted model posterior probabilities by the unregularized network. Thus, all datasets generated by model

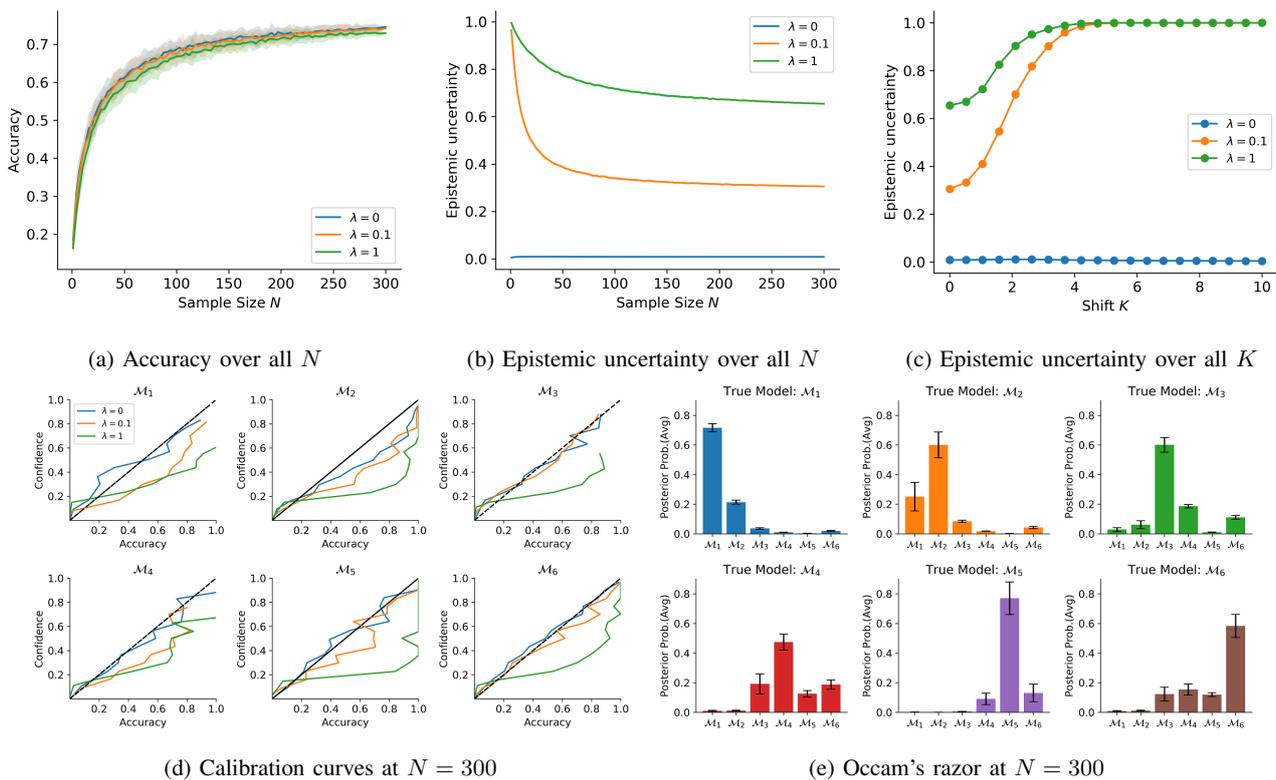


Fig. 5: Detailed validation results from **Experiment 3**.

j are plausible under the remaining models $\mathcal{M}_i, i > j$. These average model probabilities are depicted in Figure 5e. Even though data-set generated by the nested simpler models are plausible under the more complex models, we observe that Occam's razor is encoded by the behavior of the network, which, on average, consistently selects the simpler model when it is the true data-generating model. We also observe that this behavior is independent of regularization (results for $\lambda = 0.1$ and $\lambda = 1$ are not depicted in Figure 5e)

Epistemic uncertainty and absolute evidence. Epistemic uncertainty over different trial numbers (N) is zero when no KL regularization is applied ($\lambda = 0$). On the other hand, both small ($\lambda = 0.1$) or large ($\lambda = 1.0$) regularization weights lead to non-zero uncertainty over all possible N (Figure 5b). This pattern reflects a reduction in epistemic uncertainty with increasing amount of information and mirrors the inverse of the recovery curve. Note, that the value at which epistemic uncertainty begins to flatten out is larger for the highly regularized model, as it encodes more cautiousness with respect to the challenging task of selecting a true nested model. Finally, results on shifted data-sets are depicted in Figure 5c. Indeed, we observe that the regularized networks are able to detect implausible datasets and output total uncertainty around $K > 4$ for all manipulated datasets. Uncertainty increases faster for high regularization. On the other hand, the unregularized model does not have any way of signaling impossibility of a decision, so its uncertainty remains at 0 over all K .

D. Experiment 4: Stochastic Models of Single-Neuron Activity

In this experiment, we apply our evidential method to complex nested spiking neuron models describing the properties of biological cells in the nervous system. The purpose of this experiment is threefold. First, we want to assess the ability of our method to classify models deploying a variety of spiking patterns which might account for different cortical and sub-cortical neuronal activity. Second, we want to challenge the network's ability to detect biologically implausible data patterns as accounted by epistemic uncertainty. Finally, we compare our method with other viable neural network architectures that are able to perform amortized model comparison as classification. To this aim, we rely on a renowned computational model of biological neuronal dynamics.

1) Model Comparison Setting: In computational neuroscience, mathematical modeling of neuroelectric dynamics serves as a basis to understand the functional organization of the brain from both single-cell and large-scale network processing perspectives [22], [5], [23], [2], [10]. A plurality of different neural models have been proposed during the last decades, spanning from completely abstract to biologically detailed models. The former offer a simplified mathematical representation able to account for the main functional properties of spiking neurons, the latter provide a detailed analogy between models' state variables and ion channels in biological neurons [41]. Importantly, these computational models differ in their capability to reproduce firing patterns observed in real cortical neurons [24].

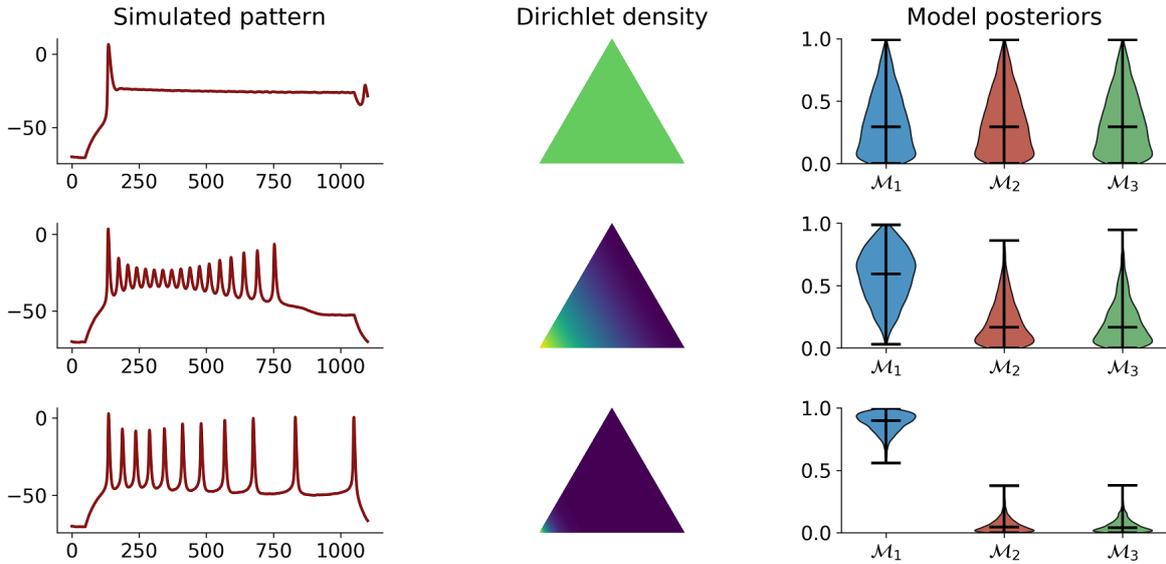


Fig. 6: Three simulated firing patterns, corresponding estimated Dirichlet densities and model posteriors. Each row illustrates a different value of the parameter \bar{g}_K : $\bar{g}_K = 0.1$, $\bar{g}_K = 0.5$, and $\bar{g}_K = 0.75$, respectively. An increase in the parameter \bar{g}_K is accompanied by a decrease in epistemic uncertainty (as measured via Eq.10). An implausible value of \bar{g}_k (first row) results in a flat density as an index of total uncertainty (uniform green areas). As the parameter value surpasses the plausible boundary (second and third rows), the Dirichlet simplex becomes peaked towards the lower left edge encoding \mathcal{M}_1 .

The model family we consider here is a Hodgkin-Huxley type model of cerebral cortex and thalamic neurons [42], [22]. The forward model is formulated as a set of five ordinary differential equations (ODEs) describing how the neuron membrane potential $V(t)$ unfolds in time as a function of an injected current $I_{inj}(t)$, and ion channels properties. See **Appendix D** for more details regarding the forward simulation process.

To set up the model comparison problem, we treat different types of conductance, $g_L, \bar{g}_{Na}, \bar{g}_K$ and \bar{g}_M , as free parameters, and formulate three neural models based on different parameter configurations. In particular, we consider three models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3\}$ defined by the parameter sets $\theta_1 = (\bar{g}_{Na}, \bar{g}_K)$, $\theta_2 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M)$, and $\theta_3 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M, g_L)$. When not treated as free parameters, we set \bar{g}_M and g_L to default values, such that $\bar{g}_M = 0.07$ and $g_L = 0.1$.

We compare the performance of our evidential method to the following methods: a standard softmax classifier, a classifier with Monte Carlo dropout [11], and two variational classifiers using a Kullback-Leibler [29] and a maximum mean discrepancy (MMD, [63]) latent space regularizer, respectively. We also train three evidential networks with $\lambda = 0$ (no regularization), $\lambda = 0.5$, and $\lambda = 1.0$ to better quantify the effects of performing regularized model comparison. Here, we do not consider non-amortized methods, such as ABC or ABC-MCMC, as implemented in [38], since they would have taken an infeasible amount of time to validate on hundreds of datasets.

2) *Validation Results*: In order to assess performance, we train an unregularized evidential network for 60 epochs resulting in 60000 mini-batch updates. For each batch, we draw

TABLE I: Comparison results from **Experiment 4**

Neural Architecture	Accuracy	Calibration Error
Evidential ($\lambda = 0$)	0.919 ± 0.004	0.078 ± 0.010
Evidential ($\lambda = 0.5$)	0.917 ± 0.006	0.097 ± 0.012
Evidential ($\lambda = 1.0$)	0.900 ± 0.006	0.095 ± 0.011
Softmax classifier	0.913 ± 0.006	0.105 ± 0.015
MC Dropout classifier	0.885 ± 0.005	0.087 ± 0.009
MMD-VAE classifier	0.906 ± 0.006	0.091 ± 0.012
VAE classifier	0.924 ± 0.005	0.096 ± 0.012

a random input current duration $T \sim \mathcal{U}_D(100, 400)$ (in units of milliseconds), with the same constant input current, I_{inj} , for each dataset simulation. Here, T reflects the physical time window in which biological spiking patterns can unfold. Since the sampling rate of membrane potential is fixed ($dt = 0.2$), T affects both the span of observable spiking behavior and the number of simulated data points. The entire training phase with online learning took approximately 2.5 hours. On the other hand, model comparison on 5000 validation time series took approximately 0.7 seconds, which highlights the extreme efficiency gains obtainable via globally amortized inference.

Regarding model selection, we observe accuracies above 0.92 across all T , with no gains in accuracy for increasing T , which shows that even short input currents are sufficient for performing reliable model selection for these complex models. Further, mean bootstrap calibration curves and accuracies on 5000 validation datasets are depicted in Figure S4d. We observe good calibration for all three models, with calibration errors less than 0.1. Notably, overconfidence was 0 for all three models. The normalized *confusion matrix* is depicted in Figure S4b.

In order to assess how well we can capture epistemic uncertainty for biologically implausible firing patterns, we train another evidential network with a gradually increasing regularization weight up to $\lambda = 1.0$. We then fix the parameter $\bar{g}_{Na} = 4.0$ of model \mathcal{M}_1 and gradually increase its second parameter \bar{g}_K from 0.1 to 2.0. Since spiking patterns observed with low values of \bar{g}_K are quite implausible and have not been observed during training, we expect uncertainty to gradually decrease. Indeed, Figure S4c shows this pattern. Three example firing patterns and the corresponding posterior estimates are depicted in Figure 6. On the other hand, changing the sign of the output membrane potential, which results in biologically implausible firing patterns, leads to a trivial selection of \mathcal{M}_3 . This is contrary to expectations, and shows that absolute evidence is also relative to what the evidential network has learned during training.

Finally, Table I presents the comparison results in terms of accuracy and calibration error (all methods achieved 0 overconfidence). We train each neural network method for 30 epochs with identical optimizer settings and the same recurrent network architecture for ease of comparison. We then compute validation metrics on 3000 simulated neural firing patterns and report means and standard errors. Our unregularized evidential network ($\lambda = 0$) achieves the lowest calibration error, followed by the MC dropout classifier. In terms of accuracy, the KL variational classifier performs slightly better than our unregularized evidential network (but still within one standard error). Overall, the performance of all amortized methods considered in this experiment is similar, which highlights the viability of the approach to Bayesian model comparison in general. Note, that training of each method took less than 1.5 hours, and bootstrap validation on 3000 less than a minute. The latter would have been impossible to achieve within a reasonable time-frame using non-amortized methods.

VI. DISCUSSION

In the current work, we introduced a novel simulation-based method for approximate Bayesian model comparison based on specialized evidential neural networks. We demonstrated that our method can successfully deal with both exchangeable and non-exchangeable (time-dependent) sequences with variable numbers of observations without relying on fixed summary statistics. Further, we presented a way to amortize the process of model comparison for a given family of models by splitting it into a potentially costly global training phase and a cheap inference phase. In this way, pre-trained evidential networks can be stored, shared, and reused across multiple datasets and model comparison applications. Finally, we demonstrated a way to obtain a measure of absolute evidence in spite of operating in an \mathcal{M} -closed framework during the simulation phase. In the following, we reiterate the main advantages of our method.

Theoretical guarantee. By using a strictly proper loss [15], we showed that our method can closely approximate analytic model posterior probabilities and Bayes factors in theory and practice. In other words, posterior probability estimates are perfectly calibrated to the true model posterior probabilities when the strictly proper logarithmic loss is globally minimized.

Indeed, our experiments confirm that the network outputs are well calibrated. However, when optimizing the regularized version of the logarithmic loss, we are no longer working with a strictly proper loss, so calibration declines at the cost of capturing implausible datasets. However, we demonstrated that the accuracy of recovery (i.e., selecting the most plausible model in the set of considered models) does not suffer when training with regularization. In any case, perfect convergence is never guaranteed in finite-sample scenarios, so validation tools such as calibration and accuracy curves are indispensable in practical applications.

Amortized inference. Following ideas from *inference compilation* [32] and *pre-paid parameter estimation* [39], our method avoids fitting each candidate model to each dataset separately. Instead, we cast the problem of model comparison as a supervised learning of absolute evidence and train a specialized neural network architecture to assign model evidences to each possible dataset. This requires only the specification of plausible priors over each model's parameters and the corresponding forward process, from which simulations can be obtained on the fly. During the upfront training, we use online learning to avoid storage overhead due to large simulated grids or reference tables [37], [39]. Importantly, the separation of model comparison into a costly upfront training phase and a cheap inference phase ensures that subsequent applications of the pre-trained networks to multiple observed datasets are very efficient. Indeed, we showed in our experiments that inference on thousands of datasets can take less than a second with our method. Moreover, by sharing and applying a pre-trained network for inference within a particular research domain, results will be highly reproducible, since the *settings* of the method will be held constant in all applications.

Raw data utilization and variable sample size. The problem of insufficient summary statistics has plagued the field of approximate Bayesian computation for a long time, so as to deserve being dubbed the *curse of insufficiency* [37]. Using sub-optimal summary statistics can severely compromise the quality of posterior approximations and validity of conclusions based on these approximations [47]. Our method avoids using hand-crafted summaries by aligning the architecture of the evidential neural network to the inherent probabilistic symmetry of the data [3]. Using specialized neural network architectures, such as permutation invariant networks or a combination of recurrent and convolutional networks, we also ensure that our method can deal with datasets containing variable numbers of observations. Moreover, by minimizing the strictly proper version of the logarithmic loss, we ensure that perfect convergence implies maximal data utilization by the network.

Absolute evidence and epistemic uncertainty. Besides point estimates of model posterior probabilities, our evidential networks yield a full higher-order probability distribution over the posterior model probabilities themselves. By choosing a Dirichlet distribution, we can use the mean of the Dirichlet distribution as the best approximation of model posterior probabilities. Beyond that, following ideas from the study of subjective logic [27] and uncertainty quantification in classification tasks [48], we can extract a measure of *epistemic uncertainty*. We employ epistemic uncertainty to quantify the

impossibility of making a model selection decision based on a dataset, which is classified as implausible under all candidate models. Therefore, the epistemic uncertainty serves as a proxy to measure absolute evidence, in contrast to relative evidence, as given by Bayes factors or posterior odds. This is an important practical advantage, as it allows us to conclude that all models in the candidate set are a poor approximation of the data-generating process of interest. Indeed, our initial experiments confirm that our measure of epistemic uncertainty increases when datasets no longer lie within the range of the considered models. However, extensive validation is needed in order to explore and understand which aspects of an observed sample lead to model misfit. Further, exploring connections to approaches using auxiliary probabilistic classifiers for detecting model misspecifications, such as the recent CARMEN method [51], seems to be an interesting avenue for future research.

These advantageous properties notwithstanding, our proposed method has certain limitations. First, our regularized optimization criterion induces a trade-off between calibration and epistemic uncertainty, as confirmed by our experiments. This trade-off is due to the fact that we capture epistemic uncertainty via a special form of Kullback-Leibler (KL) regularization during the training phase, which renders the optimized loss function no longer strictly proper. We leave it to future research to investigate whether this trade-off is fundamental and whether there are more elegant ways to quantify absolute evidence from a simulation-based perspective.

Second, our method is intended for model comparison from a prior predictive (marginal likelihood) perspective. However, since we do not explicitly fit each model to data, we cannot perform model comparison/selection based on posterior predictive performance. We note that in certain scenarios, posterior predictive performance might be a preferable metric for model comparison, so in this case, simulation-based sampling methods should be employed (e.g., ABC or neural density estimation, [8], [40]).

Third, perfect convergence might be hard to achieve in real-world applications. In this case, approximation error will propagate into model posterior estimates. Therefore, it is important to use performance diagnostic tools, such as calibration curves, accuracy of recovery, and overconfidence bounds, in order to detect potential estimation problems. Finally, even though our method exhibits excellent performance on the domain examples considered in the current work, it might break down in high-dimensional parameter spaces. Future research should focus on applications to even more challenging model comparison scenarios, for instance, hierarchical Bayesian models with intractable likelihoods, or neural network models.

ACKNOWLEDGMENT

This research was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation; grant number GRK 2277 "Statistical Modeling in Psychology"). We thank the Technology Industries of Finland Centennial Foundation (grant 70007503; Artificial Intelligence for Research and Development) for partial support of this work. We also thank David Izydorczyk and Mattia Sensi for reading the paper and providing constructive feedback.

REFERENCES

- [1] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pages 265–283, 2016.
- [2] LF Abbott and Thomas B Kepler. Model neurons from hodgkin-huxley to hopfield. In *Statistical mechanics of neural networks*, pages 5–18. Springer, 1990.
- [3] Benjamin Bloem-Reddy and Yee Whye Teh. Probabilistic symmetry and invariant neural networks. *arXiv preprint arXiv:1901.06082*, 2019.
- [4] Scott D Brown and Andrew Heathcote. The simplest complete model of choice response time: Linear ballistic accumulation. *Cognitive psychology*, 57(3):153–178, 2008.
- [5] Anthony N Burkitt. A review of the integrate-and-fire neuron model: I. homogeneous synaptic input. *Biological cybernetics*, 95(1):1–19, 2006.
- [6] Paul-Christian Bürkner, Jonah Gabry, and Aki Vehtari. Approximate leave-future-out cross-validation for bayesian time series models. *Journal of Statistical Computation and Simulation*, 90(14):2499–2523, 2020.
- [7] Kyle Cranmer, Johann Brehmer, and Gilles Louppe. The frontier of simulation-based inference. *Proceedings of the National Academy of Sciences*, 117(48):30055–30062, 2020.
- [8] José Mir Justino da Costa, Helcio Rangel Barreto Orlando, and Wellington Betencurte da Silva. Model selection and parameter estimation in tumor growth models using approximate bayesian computation-abc. *Computational and Applied Mathematics*, 37(3):2795–2815, 2018.
- [9] Thomas S Deisboeck and Georgios S Stamatakos. *Multiscale cancer modeling*. CRC press, 2010.
- [10] Salvador Dura-Bernal, Benjamin A Suter, Pdraig Gleeson, Matteo Cantarelli, Adrian Quintana, Facundo Rodriguez, David J Kedziora, George L Chadderton, Cliff C Kerr, Samuel A Neymotin, et al. Netpyne, a tool for data-driven multiscale modeling of brain circuits. *Elife*, 8:e44494, 2019.
- [11] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [12] Andrew Gelman, John B. Carlin, Hal S. Stern, David B. Dunson, Aki Vehtari, and Donald B. Rubin. *Bayesian Data Analysis, Third Edition*. Chapman and Hall/CRC.
- [13] Andrew Gelman and Xiao-Li Meng. Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical science*, pages 163–185, 1998.
- [14] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with lstm. 1999.
- [15] Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [16] David S Greenberg, Marcel Nonnenmacher, and Jakob H Macke. Automatic posterior transformation for likelihood-free inference. *arXiv preprint arXiv:1905.07488*, 2019.
- [17] Quentin F Gronau, Alexandra Sarafoglou, Dora Matzke, Alexander Ly, Udo Boehm, Maarten Marsman, David S Leslie, Jonathan J Forster, Eric-Jan Wagenmakers, and Helen Steingroever. A tutorial on bridge sampling. *Journal of mathematical psychology*, 81:80–97, 2017.
- [18] Peter Grünwald, Thijs Van Ommen, et al. Inconsistency of bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4):1069–1103, 2017.
- [19] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1321–1330. JMLR. org, 2017.
- [20] Joeri Hermans, Volodimir Begy, and Gilles Louppe. Likelihood-free mcmc with amortized approximate ratio estimators. In *International Conference on Machine Learning*, pages 4239–4248. PMLR, 2020.
- [21] José Miguel Hernández-Lobato and Ryan Adams. Probabilistic backpropagation for scalable learning of bayesian neural networks. In *International Conference on Machine Learning*, pages 1861–1869. PMLR, 2015.
- [22] Alan L Hodgkin and Andrew F Huxley. A quantitative description of membrane current and its application to conduction and excitation in nerve. *The Journal of physiology*, 117(4):500–544, 1952.
- [23] Eugene M Izhikevich. Simple model of spiking neurons. *IEEE Transactions on neural networks*, 14(6):1569–1572, 2003.
- [24] Eugene M Izhikevich. Which model to use for cortical spiking neurons? *IEEE transactions on neural networks*, 15(5):1063–1070, 2004.
- [25] Edwin T Jaynes. *Probability theory: The logic of science*. Cambridge university press, 2003.

- [26] Bai Jiang, Tung-yu Wu, Charles Zheng, and Wing H Wong. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.
- [27] Audun Jsang. *Subjective Logic: A formalism for reasoning under uncertainty*. Springer Publishing Company, Incorporated, 2018.
- [28] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 5580–5590, 2017.
- [29] Diederik P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pages 2575–2583, 2015.
- [30] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.
- [31] Emmanuel Klinger, Dennis Rickert, and Jan Hasenauer. pyabc: distributed, likelihood-free inference. *Bioinformatics*, 34(20):3591–3593, 2018.
- [32] Tuan Anh Le, Atilim Gunes Baydin, and Frank Wood. Inference compilation and universal probabilistic programming. *arXiv preprint arXiv:1610.09900*, 2016.
- [33] Yingzhen Li and Yarin Gal. Dropout inference in bayesian neural networks with alpha-divergences. In *International conference on machine learning*, pages 2052–2061. PMLR, 2017.
- [34] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [35] Christos Louizos and Max Welling. Multiplicative normalizing flows for variational bayesian neural networks. In *International Conference on Machine Learning*, pages 2218–2227. PMLR, 2017.
- [36] David JC MacKay. Bayesian neural networks and density networks. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 354(1):73–80, 1995.
- [37] Jean-Michel Marin, Pierre Pudlo, Arnaud Estoup, and Christian Robert. *Likelihood-free model choice*. Chapman and Hall/CRC Press Boca Raton, FL, 2018.
- [38] Ulf Kai Mertens, Andreas Voss, and Stefan Radev. Abrox—a user-friendly python module for approximate bayesian computation with a focus on model comparison. *PLoS one*, 13(3):e0193981, 2018.
- [39] Merijn Mestdagh, Stijn Verdonck, Kristof Meers, Tim Loossens, and Francis Tuerlinckx. Prepaid parameter estimation without likelihoods. *PLoS computational biology*, 15(9):e1007181, 2019.
- [40] George Papamakarios, David C Sterratt, and Iain Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. *arXiv preprint arXiv:1805.07226*, 2018.
- [41] Martin Pospischil, Zuzanna Piwkowska, Thierry Bal, and Alain Destexhe. Comparison of different neuron models to conductance-based post-stimulus time histograms obtained in cortical pyramidal cells using dynamic-clamp in vitro. *Biological cybernetics*, 105(2):167, 2011.
- [42] Martin Pospischil, María Toledo-Rodriguez, Cyril Monier, Zuzanna Piwkowska, Thierry Bal, Yves Frégnac, Henry Markram, and Alain Destexhe. Minimal hodgkin–huxley type models for different classes of cortical and thalamic neurons. *Biological cybernetics*, 99(4-5):427–441, 2008.
- [43] Pierre Pudlo, Jean-Michel Marin, Arnaud Estoup, Jean-Marie Cornuet, Mathieu Gautier, and Christian P Robert. Reliable abc model choice via random forests. *Bioinformatics*, 32(6):859–866, 2015.
- [44] Stefan T Radev, Ulf K Mertens, Andreas Voss, Lynton Ardizzone, and Ullrich Köthe. Bayesflow: Learning complex stochastic models with invertible neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [45] Roger Ratcliff and Gail McKoon. The diffusion decision model: theory and data for two-choice decision tasks. *Neural computation*, 20(4):873–922, 2008.
- [46] Oliver Ratmann, Christophe Andrieu, Carsten Wiuf, and Sylvia Richardson. Model criticism based on likelihood-free inference, with an application to protein network evolution. *Proceedings of the National Academy of Sciences*, 106(26):10576–10581, 2009.
- [47] Christian P Robert, Jean-Marie Cornuet, Jean-Michel Marin, and Natesh S Pillai. Lack of confidence in approximate bayesian computation model choice. *Proceedings of the National Academy of Sciences*, 108(37):15112–15117, 2011.
- [48] Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification uncertainty. In *Advances in Neural Information Processing Systems*, pages 3179–3189, 2018.
- [49] Scott A Sisson and Yanan Fan. *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839], 2011.
- [50] Mikael Sunnåker, Alberto Giovanni Busetto, Elina Numminen, Jukka Corander, Matthieu Foll, and Christophe Dessimoz. Approximate bayesian computation. *PLoS computational biology*, 9(1):e1002803, 2013.
- [51] Owen Thomas and Jukka Corander. Diagnosing model misspecification and performing generalized bayes’ updates via probabilistic classifiers. *arXiv preprint arXiv:1912.05810*, 2019.
- [52] Brandon M Turner and Per B Sederberg. A generalized, likelihood-free method for posterior estimation. *Psychonomic bulletin & review*, 21(2):227–250, 2014.
- [53] Brandon M Turner, Per B Sederberg, and James L McClelland. Bayesian analysis of simulation-based models. *Journal of Mathematical Psychology*, 72:191–199, 2016.
- [54] Marius Usher and James L McClelland. The time course of perceptual choice: the leaky, competing accumulator model. *Psychological review*, 108(3):550, 2001.
- [55] Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical bayesian model evaluation using leave-one-out cross-validation and waic. *Statistics and computing*, 27(5):1413–1432, 2017.
- [56] Aki Vehtari, Janne Ojanen, et al. A survey of bayesian predictive methods for model assessment, selection and comparison. *Statistics Surveys*, 6:142–228, 2012.
- [57] Robert Verity, Lucy C Okell, Ilaria Dorigatti, Peter Winskill, Charles Whittaker, Natsuko Imai, Gina Cuomo-Dannenburg, Hayley Thompson, Patrick GT Walker, Han Fu, et al. Estimates of the severity of coronavirus disease 2019: a model-based analysis. *The Lancet Infectious Diseases*, 2020.
- [58] Andreas Voss, Veronika Lerche, Ulf Mertens, and Jochen Voss. Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models. *Psychonomic bulletin & review*, 26(3):813–832, 2019.
- [59] Sumio Watanabe. Waic and wbic are information criteria for singular statistical model evaluation. In *Proceedings of the Workshop on Information Theoretic Methods in Science and Engineering*, pages 90–94, 2013.
- [60] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? *arXiv preprint arXiv:1905.13211*, 2019.
- [61] Yuling Yao, Aki Vehtari, Daniel Simpson, Andrew Gelman, et al. Using stacking to average bayesian predictive distributions (with discussion). *Bayesian Analysis*, 13(3):917–1007, 2018.
- [62] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. In *Advances in neural information processing systems*, pages 3391–3401, 2017.
- [63] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Information maximizing variational autoencoders. *arXiv preprint arXiv:1706.02262*, 2017.

APPENDIX

A. Neural Network Architectures

As already discussed, we need specialized neural network architectures for dealing with datasets with variable numbers of observations N and different probabilistic symmetries (e.g., *i.i.d.* or temporal ordering). In the following, we describe the network architectures used to tackle the most common probabilistic symmetries observed in the social and life sciences, that is, exchangeable and temporal sequences.

Exchangeable Sequences: The most prominent feature of *i.i.d.* sequences is *permutation invariance*, since changing the order of individual elements does not change the associated likelihood function or posterior. Accordingly, if we denote an arbitrary permutation of N elements as $\mathbb{S}_N(\cdot)$, the following should hold for the associated model posteriors:

$$p(\mathbf{m} | \mathbf{x}_{1:N}) = p(\mathbf{m} | \mathbb{S}_N(\mathbf{x}_{1:N})) \quad (1)$$

We encode probabilistic permutation invariance by enforcing functional permutation invariance with respect to the outputs of the evidential network [3]. Following recent work on deep sets [62] and probabilistic symmetry [3], we implement a deep invariant network performing a series of equivariant and invariant transformations. An invariant transformation is characterized by:

$$f(\mathbb{S}_N(\mathbf{x}_{1:N})) = f(\mathbf{x}_{1:N}) \quad (2)$$

that is, permuting the input elements does not change the resulting output. Such a transformation is often referred to as a pooling operation. On the other hand, an equivariant transformation is characterized by:

$$f(\mathbb{S}_N(\mathbf{x}_{1:N})) = \mathbb{S}_N(f(\mathbf{x}_{1:N})) \quad (3)$$

that is, permuting the input is equivalent to permuting the output of the transformation. We parameterize a learnable invariant function via an *invariant module* performing a sequence of non-linear transformations followed by a pooling (sum) operation and another non-linear transformation:

$$\tilde{\mathbf{z}} = \Sigma_I(\mathbf{x}_{1:N}) = f_1 \left(\sum_{i=1}^N f_2(\mathbf{x}_i) \right) \quad (4)$$

where f_1 and f_2 can be arbitrary (non-linear) functions, which we parameterize via fully connected (FC) neural networks. Figure S1 (left panel) presents a graphical illustration of the invariant module.

We parameterize a learnable equivariant transformation via an *equivariant module* performing the following operations for each input element i :

$$\mathbf{z}_i = \Sigma_E(\mathbf{x}_i, \tilde{\mathbf{z}}) = f_3(\mathbf{x}_i, \tilde{\mathbf{z}}) \quad (5)$$

so that that f_3 is a combination of element-wise and invariant transforms (see Figure S1, center). Again, we parameterize the internal function f_3 via a standard FC neural network. Note, that an equivariant module also takes as an input the output of an invariant module, in order to increase the expressiveness of the learned transformation. Thus, each equivariant module contains a separate invariant module. Finally, we can stack

multiple equivariant modules followed by an invariant module, in order to obtain a deep invariant evidential network $f_\phi : \mathcal{X}^N \rightarrow [1, \infty)^M$:

$$\boldsymbol{\alpha} = f_\phi(\mathbf{x}_{1:N}) = (\Sigma_I \circ \Sigma_E^{(L)} \circ \Sigma_E^{(L-1)} \circ \dots \circ \Sigma_E^{(1)})(\mathbf{x}_{1:N}) \quad (6)$$

where ϕ denotes the vector of all learnable neural network parameters and the final invariant module implements a shifted ReLU output non-linearity:

$$\boldsymbol{\alpha} = \max \left(1, \Sigma_I(\mathbf{z}_{1:N}^{(L)}) \right) \quad (7)$$

in order to represent Dirichlet evidences ($0 < \alpha < 1$ is technically possible and valid but not of relevance for our purposes). The rightmost panel in Figure S1 provides a graphical illustration of a deep invariant network. We use this architecture in **Experiments 1, 2 and 4**.

Non-Exchangeable Sequences: One of the most common non-exchangeable sequences encountered in practice are time-series with arbitrarily long temporal dependencies. A natural choice for time series-data with variable length are LSTM networks [14], as recurrent networks are designed to deal with long sequences of variable size. Another reasonable choice might be 1D fully convolutional networks [34], which can also process sequences with variable length. A different type of frequently encountered non-exchangeable data are images, which have successfully been tackled via 2D convolutional networks.

Since, in this work, we apply our evidential method to time-series models, we will describe an architecture for processing data with temporal dependencies. We use a combination of a LSTM and a 1D convolutional network to reduce the observed time-series into fixed-size vector representations. We then concatenate these vectors and pass them through a standard fully connected network to obtain the final Dirichlet evidences. At a high level, our architecture performs the following operations on an input sequence $\mathbf{x}_{1:N}$. First, a many-to-one LSTM network reduces the input sequence to a vector \mathbf{u} of pre-defined size. Then, a 1D convolutional network reduces the input sequence to a matrix $\mathbf{V} \in \mathbb{R}^{N' \times W}$ where N' is the length of the filtered sequence and W is the number of filters in the final convolutional layer. In this way, only the *time* dimension of V depends on the length of the input. A mean pooling operator is then applied to the time dimension of \mathbf{V} to obtain a fixed-size representation \mathbf{v} of size W . Finally, the outputs of the LSTM and convolutional networks are concatenated and fed through a FC network f with a shifted ReLU output non-linearity, which yields the Dirichlet evidences $\boldsymbol{\alpha}$. Thus, the computational flow is as follows:

$$\mathbf{u} = \text{LSTM}(\mathbf{x}_{1:N}) \quad (8)$$

$$\mathbf{v} = \frac{1}{N'} \sum_{i=1}^{N'} \mathbf{V}_i, \text{ with } \mathbf{V} = \text{Conv1D}(\mathbf{x}_{1:N}) \quad (9)$$

$$\mathbf{z} = \text{Concat}(\mathbf{u}, \mathbf{v}) \quad (10)$$

$$\boldsymbol{\alpha} = \max(1, f(\mathbf{z})) \quad (11)$$

Figure S2 illustrates the computational flow of a sequence network. This architecture provides us with a powerful estimator

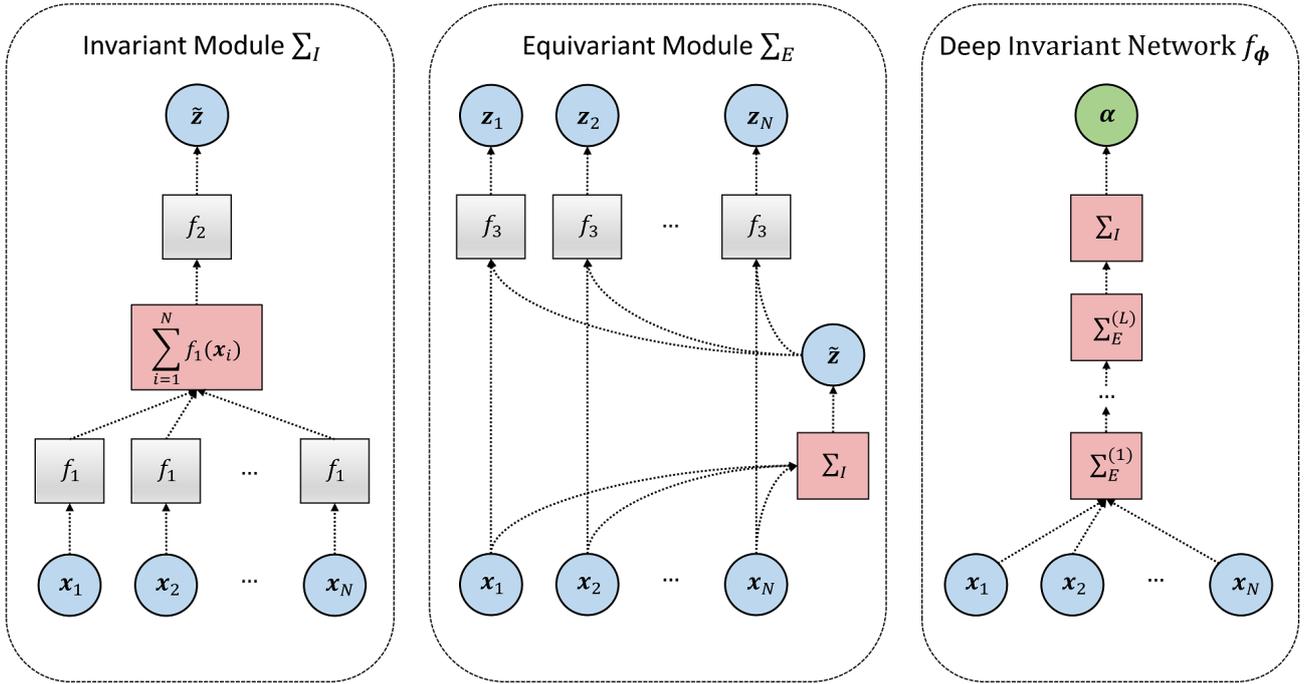


Fig. S1: The basic building blocks of a deep invariant evidential network encoding exchangeable sequences (adapted after [3]). The leftmost panel depicts an invariant module implementing an expressive permutation invariant function. The middle panel depicts an equivariant module implementing an expressive permutation equivariant function with a nested invariant module. The right panel depicts a deep invariant network consisting of a composition of equivariant modules followed by an invariant module with a ReLU output activation non-linearity.

for comparing models whose outputs consist of multivariate time-series. All neural network parameters ϕ are jointly optimized during the training phase. We use this architecture in **Experiments 3** and **5**.

B. Neural Network Training

We train all neural networks described in this paper via mini-batch gradient descent³. For all following experiments, we use the Adam optimizer with a starter learning rate of 10^{-4} and an exponential decay rate of .99. We did not perform an extensive search for optimal values of network hyperparameters. All networks were implemented in Python using the *TensorFlow* library [1] and trained on a single-GPU machine equipped with NVIDIA[®] GTX1060 graphics card. See **Appendix A** for details on the neural network architectures.

C. Performance Metrics

Throughout the following examples, we use a set of performance metrics to assess the overall performance of our method. To test how well the method is able to recover the true model (hard assignment of model indices), we compute the accuracy of recovery as the fraction of correct model assignments over the total number of test datasets, where model assignments are done

by selecting the model with the highest probability. To test how well the posterior probability estimates of the evidence network match the true model posteriors, we compute the expected calibration error (ECE, [19]). The ECE measures the gap between the confidence and the accuracy of a classifier and is an unbiased estimate of exact miscalibration [19]. In practice, we will report *calibration curves* for each model, as these are easier to interpret for multi-class classification problems. Finally, to ensure that the method does not exhibit overconfidence, we compute an *overconfidence metric* which is given by the difference between a high probability threshold T (e.g., $T = 0.95$) and the accuracy of the model above this threshold: $\text{overconfidence} = \max\{0, T - \frac{1}{|D_T|} \sum_{i \in D_T} \mathbb{1}_{[\hat{m}^{(i)} = m^{(i)}]}\}$ where D_T is the set of indices of predicted probabilities larger than T . Any deviation from zero would be indicative of overconfidence and thus lack of confidence in the method's estimates.

D. Bonus Experiment: 400 Gaussian Mixture Models

With this example, we want to show that our method is capable of performing model comparison on problems involving hundreds of competing models. Further, we want to corroborate the desired improvement in accuracy with increasing number of observations, as shown in the previous experiment.

To this aim, we construct a setting with 400 2D Gaussian Mixture Models (GMMs) with 2 mixture components. The

³Code and simulation scripts for all experiments are available at <https://github.com/stefanradev93/BayesFlow>.

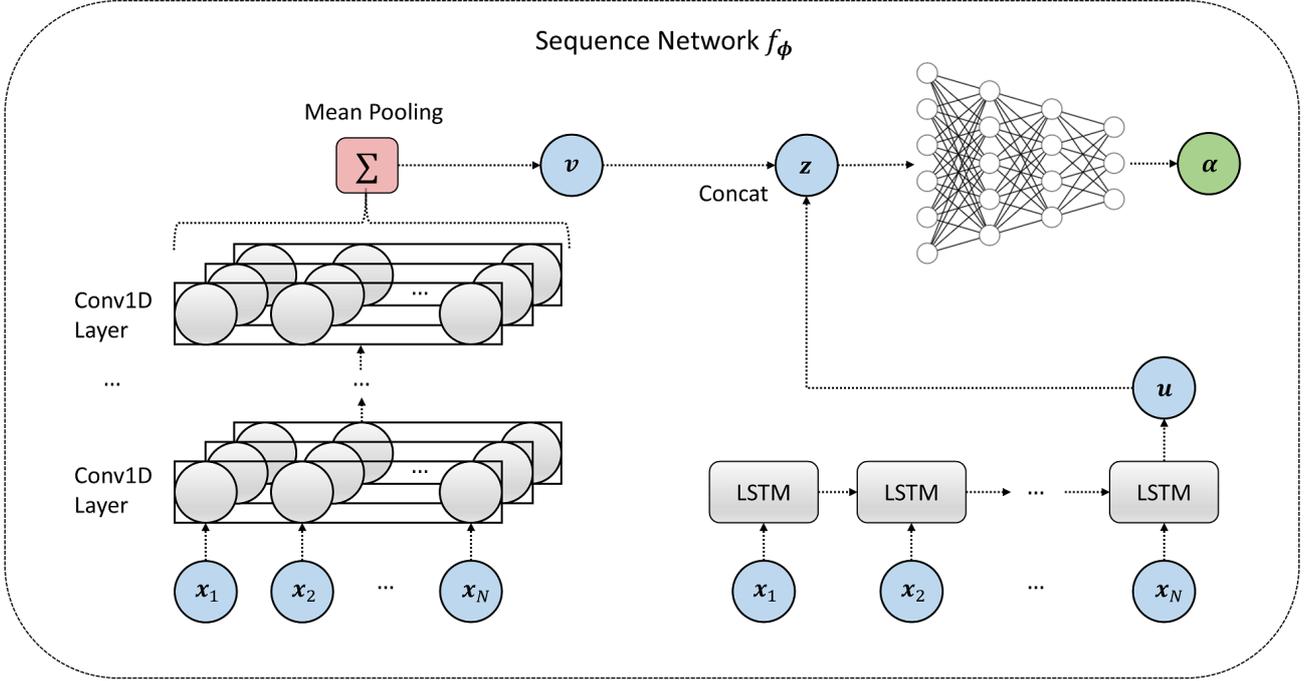


Fig. S2: The building blocks of an evidential sequence network. An input sequence generated from a model is fed through a number of convolutional layers in a hierarchical manner (left). The output of the final convolutional layer is passed through a mean pooling operator to obtain a fixed sized vector v . The input sequence is also fed through an LSTM many-to-one recurrent neural network (right) to obtain another fixed-size vector u . The vectors v and u are then concatenated into a vector z and passed through a standard fully connected network (upper right) to obtain the final Dirichlet evidences.

construction of the models and data proceeds as follows. We first specify the 400 mean vectors on two linearly-spaced 20×20 grids: $\mu_0^{(m)} \in [-10, 0] \times [-10, 0]$ and $\mu_1^{(m)} \in [0, 10] \times [0, 10]$ for $m = 1, \dots, 400$. We then generate data from each GMM model by:

$$N \sim \mathcal{U}_D(1, 250) \quad (12)$$

$$\pi^{(m)} \sim \text{Beta}(30, 30) \quad (13)$$

$$k \sim \text{Bernoulli}(\pi^{(m)}) \quad (14)$$

$$x_j^{(m)} \sim \mathcal{N}(\mu_k^{(m)}, \mathbb{I}) \text{ for } j = 1, \dots, N \quad (15)$$

where \mathcal{U}_D denotes the discrete uniform distribution and \mathbb{I} the identity matrix of appropriate dimension. Thus, each dataset consists of N *i.i.d.* samples from one of 400 GMMs with different component mean vectors. Figure S3a depicts simulated datasets from 9 GMM with linearly increasing X -coordinates showing that differences between the models are very subtle.

We train an invariant evidential network for 30 epochs for approximately 1.2 hours wall-clock time. We then validate the performance of the network on all possible N between 1 and 250 with 5000 previously unseen simulated datasets. Figure S3b depicts a heatmap of the normalized *confusion matrix* obtained on the validation datasets with $N = 250$. Importantly, we observe that the main diagonal of the heatmap indicates that the predicted model indices mostly match the true model indices, with mistakes occurring mainly between

models with very close means. Bootstrap mean accuracy at $N = 250$ was around 0.451 ($SD=0.007$), which is 180 times better than chance accuracy (0.0025) considering the dimensionality of the problem. Finally, Figure S3c depicts the accuracy of recovery over all N considered. Again, we observe that accuracy improves as more observations become available, with sub-linear scaling.

E. Priors for Experiment 4

The prior parameter distributions for all models used in Experiment 4 are listed in Table S1.

F. Details for Experiment 4

The forward model is formulated as a set of five ordinary differential equations (ODEs) describing how the neuron membrane potential $V(t)$ unfolds in time as a function of an injected current $I_{inj}(t)$, and ion channels properties. The change in membrane potential is defined by the membrane ODE:

$$C \frac{dV}{dt} = -I_L - I_{Na} - I_K - I_M + I_{inj} + \sigma\eta(t) \quad (16)$$

where C is the specific membrane capacitance, $\sigma\eta(t)$ the intrinsic neural noise, and the I_j s are the ionic currents flowing through channels, such that:

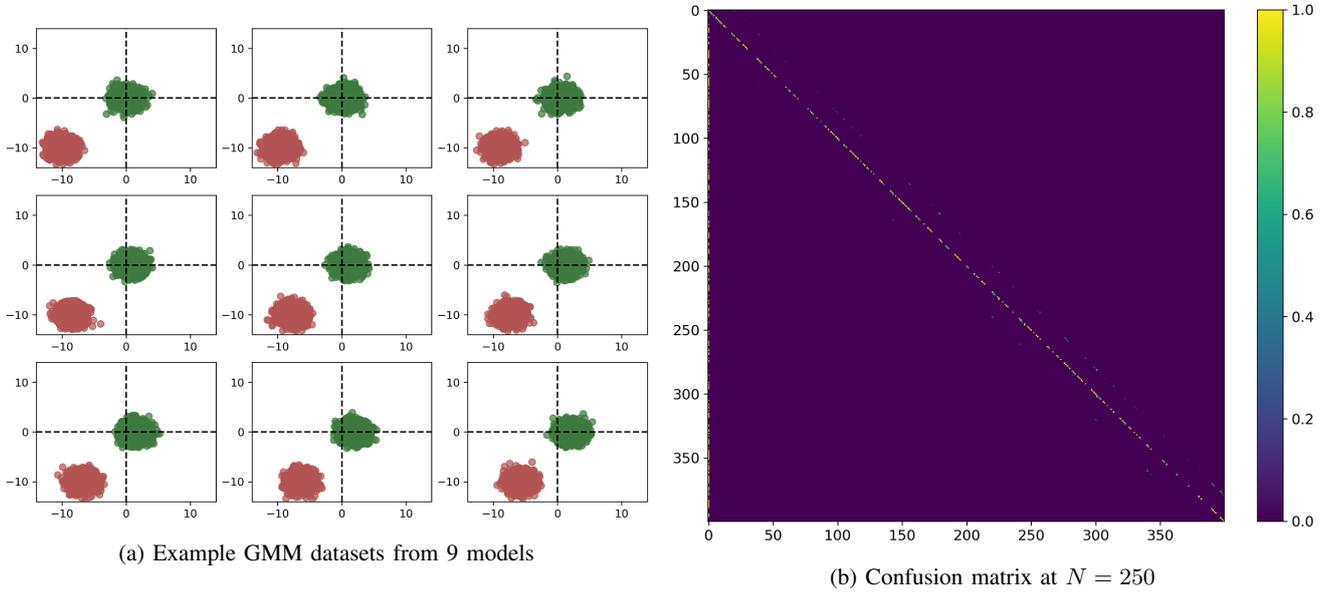


Fig. S3: Data and results from **Experiment 2**. (a) 9 example data-sets from 9 different GMM models. The X-coordinates of both clusters increase linearly left to right, and coordinate differences are very small; (b) Heatmap of the confusion matrix between true and predicted model indices. Accuracy of recovery is color-coded according to the colorbar depicted at the right. $N = 250$; (c) Accuracy of recovery as a function of sample size N .

TABLE S1: Parameter priors of the free parameters and numerical values of the fixed parameters of the six EAM models considered in **Experiment 4**.

	v_1	v_2	a	t_0	z_r	α	st_0	s_v	s_{zr}
\mathcal{M}_1	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	0.5	0	0	0	0
\mathcal{M}_2	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	$\mathcal{U}(0.3, 0.7)$	0	0	0	0
\mathcal{M}_3	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	$\mathcal{U}(0.3, 0.7)$	$\mathcal{U}(1, 2)$	0	0	0
\mathcal{M}_4	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	$\mathcal{U}(0.3, 0.7)$	$\mathcal{U}(1, 2)$	$\mathcal{U}(0, 0.4)$	0	0
\mathcal{M}_5	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	$\mathcal{U}(0.3, 0.7)$	$\mathcal{U}(1, 2)$	$\mathcal{U}(0, 0.4)$	$\mathcal{U}(0, 2)$	0
\mathcal{M}_6	$\mathcal{U}(0, 6)$	$\mathcal{U}(-6, 0)$	$\mathcal{U}(0.6, 3)$	$\mathcal{U}(0.2, 1.5)$	$\mathcal{U}(0.3, 0.7)$	$\mathcal{U}(1, 2)$	$\mathcal{U}(0, 0.4)$	$\mathcal{U}(0, 2)$	$\mathcal{U}(0, 0.6)$

$$I_L = g_L(V - E_L) \quad (17)$$

$$I_{Na} = \bar{g}_{Na} m^3 h (V - E_{Na}) \quad (18)$$

$$I_K = \bar{g}_K n^4 (V - E_K) \quad (19)$$

$$I_M = \bar{g}_M p (V - E_M). \quad (20)$$

Here, g_L is the leak conductance, while \bar{g}_{Na} , \bar{g}_K , \bar{g}_M are the sodium, potassium and M-type channel maximum conductances, respectively. E_L , E_{Na} and E_K denote the leak equilibrium potential, the sodium and potassium reversal potentials, respectively. In particular, g_L is assumed constant through time, whilst the other conductances vary over time. Consistently, (m, h, n, p) indicates the vector of the state variables accounting for ion channel gating kinetics evolving according to the following set of ODEs:

$$\frac{di}{dt} = \alpha_i(V)(1 - i) - \beta_i(V)i \quad (21)$$

$$\frac{dp}{dt} = \frac{p_\infty(V) - p}{\tau_p(V)} \quad (22)$$

where $i \in \{m, h, n\}$, and $\alpha_i(V)$, $\beta_i(V)$, $p_\infty(V)$ and $\tau_p(V)$ are nonlinear functions of the membrane potential (see [42] for details).

In our simulated experiment, we treat conductances g_L , \bar{g}_{Na} , \bar{g}_K and \bar{g}_M as free parameters, and consider different neuronal models based on different parameter configurations. It is also assumed that such configurations allow to affect the span of the possible firing patterns attainable by each model. In particular, we consider 3 models $\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3\}$ defined by the parameter sets $\theta_1 = (\bar{g}_{Na}, \bar{g}_K)$, $\theta_2 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M)$ and $\theta_3 = (\bar{g}_{Na}, \bar{g}_K, \bar{g}_M, g_L)$. We place the following priors over the parameters \bar{g}_{Na} and \bar{g}_K :

$$\bar{g}_{Na} \sim \mathcal{U}(1.5, 30) \quad (23)$$

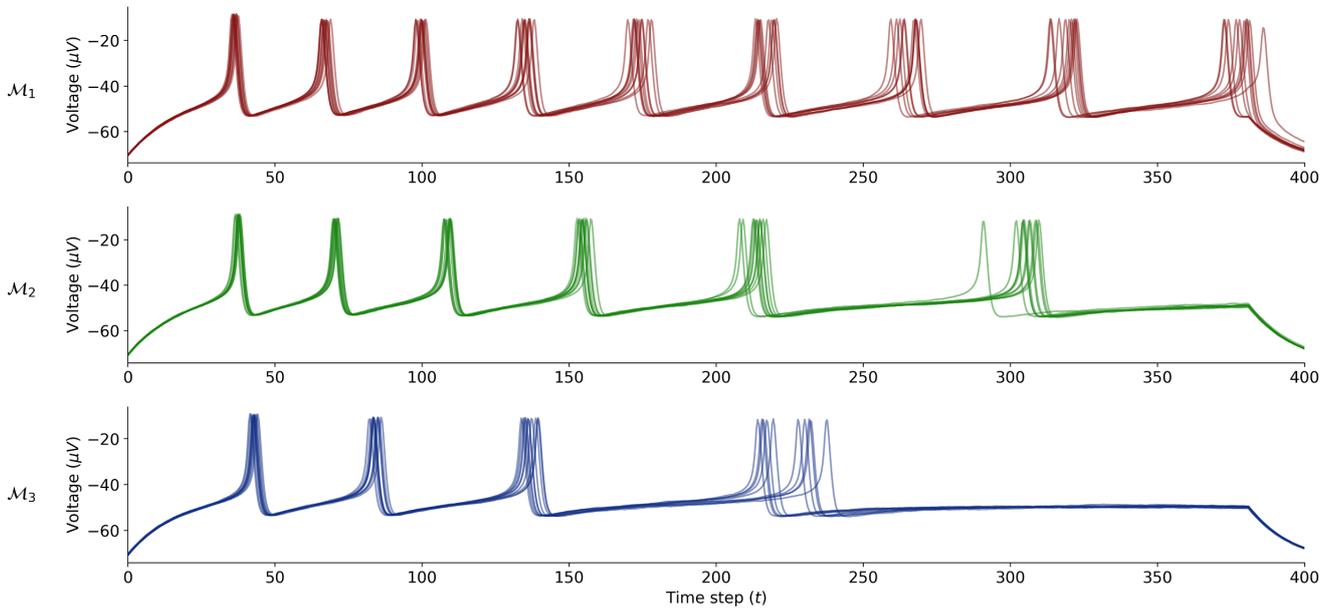
$$\bar{g}_K \sim \mathcal{U}(0.3, 15) \quad (24)$$

When not considered as free parameters, \bar{g}_M and g_L are set to default values, such that $\bar{g}_M = 0.07$ and $g_L = 0.1$, otherwise parameters are drawn from the following priors:

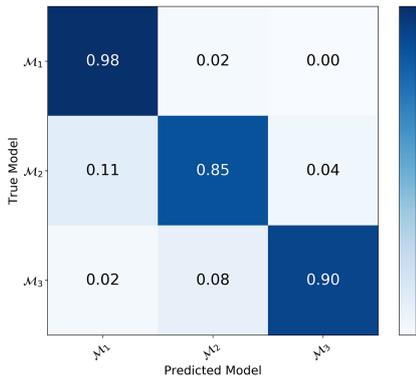
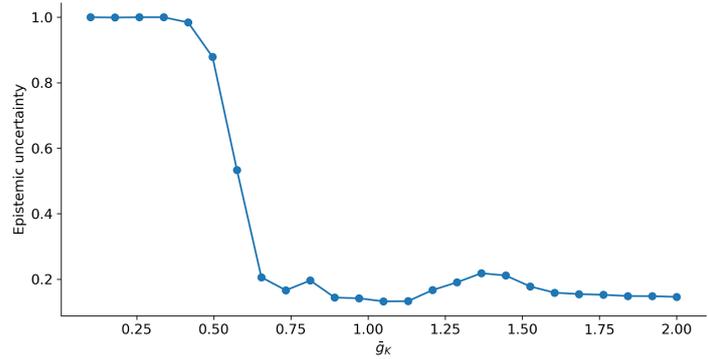
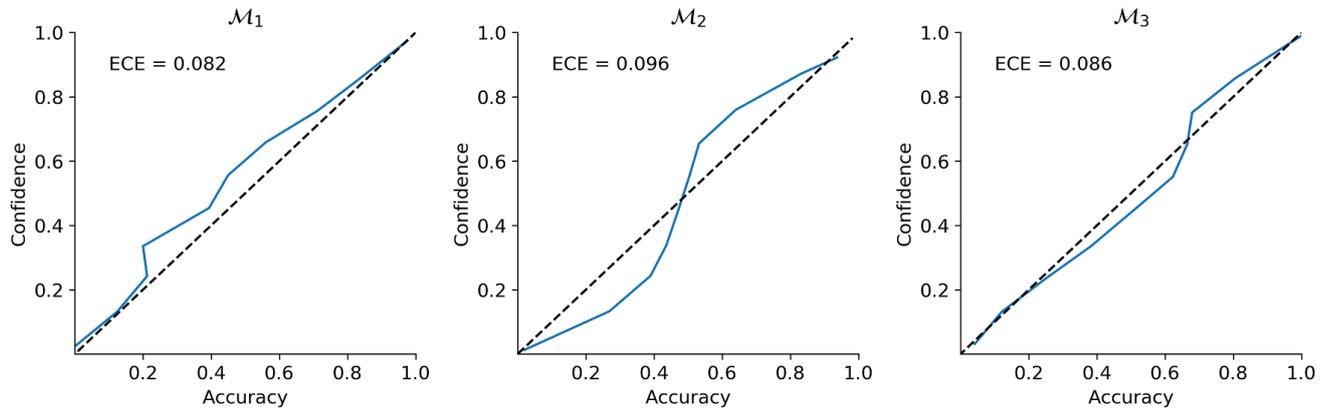
$$\bar{g}_M \sim \mathcal{U}(0.005, 0.3) \quad (25)$$

$$\bar{g}_L \sim \mathcal{U}(0.01, 0.18) \quad (26)$$

Figure S4a depicts 50 example runs from each model with respective parameters $\theta_1 = (3.0, 2.0)$, $\theta_2 = (3.0, 2.0, 0.1)$, $\theta_3 = (3.0, 2.0, 0.1, 0.11)$.



(a) Example spiking patterns from all three models

(b) Confusion matrix with $T = 100$ (c) Epistemic uncertainty as a function of \bar{g}_K (d) Calibration at $T = 400$ Fig. S4: Example Hodgkin-Huxley spiking patterns and validation results from **Experiment 5**.

APPENDIX A3 - MANUSCRIPT 3

Manuscript 3: Amortized Bayesian inference for models of cognition.

Amortized Bayesian Inference for Models of Cognition

Stefan T. Radev

Heidelberg University, Germany

Eva Marie Wieschen

Heidelberg University, Germany

Andreas Voss

Heidelberg University, Germany

Paul-Christian Bürkner

Aalto University, Finland

Abstract

As models of cognition grow in complexity and number of parameters, Bayesian inference with standard methods can become intractable, especially when the data-generating model is of unknown analytic form. Recent advances in simulation-based inference using specialized neural network architectures circumvent many previous problems of approximate Bayesian computation. Moreover, due to the properties of these special neural network estimators, the effort of training the networks via simulations amortizes over subsequent evaluations which can re-use the same network for multiple datasets and across multiple researchers. However, these methods have been largely underutilized in cognitive science and psychology so far, even though they are well suited for tackling a wide variety of modeling problems. With this work, we provide a general introduction to amortized Bayesian parameter estimation and model comparison and demonstrate the applicability of the proposed methods on a well-known class of intractable response-time models.

Keywords: Bayesian inference; Neural networks; Cognitive models; Deep learning; Simulation

Generative Models in Cognitive Science

Mathematical models formalize theories of cognition and enable the systematic investigation of cognitive processes through simulations and testable predictions. They enable a systematic joint analysis of behavioral and neural data, bridging a crucial gap between cognitive science and neuroscience (B. M. Turner, Forstmann, Steyvers, et al., 2019). Moreover, questions demanding a choice among competing cognitive theories can be resolved at the level of formal model comparison.

The *generative* property of such models arises from the fact that one can simulate the process of interest and study how it behaves under various conditions. More formally, consider a cognitive model \mathcal{M} which represents a theoretically plausible, potentially noisy, process by which observable behavior x arises from an assumed cognitive system governed by hidden parameters θ and an independent source of noise $\xi \sim p(\xi)$:

$$x = \mathcal{M}(\theta, \xi) \quad (1)$$

Generative models of this form have been developed in various domains throughout psychology and cognitive science, including decision making (Voss, Lerche, Mertens, & Voss, 2019), memory (Myung, Montenegro, & Pitt, 2007), reinforcement learning (Fontanesi, Gluth, Spektor, & Rieskamp, 2019), risky behavior (Stout, Bussemeyer, Lin, Grant, & Bonson, 2004), to name just a few. Once a model (or a set of

models) of some cognitive process of interest has been formulated, the challenge becomes to perform inference on real data. We will now briefly review the mathematical tools provided by Bayesian probability theory for parameter estimation and model comparison (Jaynes, 2003). Then, we will peruse a novel framework for performing Bayesian inference on models of cognition which are intractable with standard Bayesian methods.

Bayesian Parameter Estimation

Bayesian parameter estimation leverages prior knowledge about reasonable parameter ranges and integrates this information with the information provided by the data to arrive at a *posterior distribution* over parameters. In a Bayesian context, the posterior encodes our updated belief about plausible parameter ranges conditional on a set of N observations $X := \{x_n\}_{n=1}^N$. Bayes' rule gives us the well known analytical form of the posterior:

$$p(\theta|X) = \frac{p(X|\theta)p(\theta)}{\int p(X|\theta)p(\theta)d\theta} \quad (2)$$

where $p(X|\theta)$ represents the *likelihood* of the parameters θ and $p(\theta)$ denotes the *prior*, that is the distribution of θ before observing the data. The denominator is a normalizing constant usually referred to as the *marginal likelihood* or *evidence*. Note, that all distributions are also implicitly conditional on the particular generative model \mathcal{M} .

Based on the obtained estimate of the posterior distribution, usually in the form of random draws from the posterior, summary statistics such as posterior means or credible intervals for each parameter can be obtained. What is more, the posterior distribution can be further transformed to obtain subsequent quantities of interest, for example, the *posterior predictive distribution* which can be compared to the observed data for the purpose of model checking (Lynch & Western, 2004).

Bayesian Model Comparison

In many research domains, there is not a *single model* for a particular process, but whole *classes* of models instantiating different and often competing theories. Bayesian model comparison proceeds by assigning a plausibility value to each candidate model. These plausibility values (model weights,

model probabilities, model predictions, etc.) can be used to guide subsequent model selection.

To set the stage, consider a set of J candidate models $\mathcal{G} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_J\}$. An intuitive way to quantify plausibility is to consider the marginal likelihood of a model \mathcal{M} given by:

$$p(X|\mathcal{M}) = \int p(X|\theta, \mathcal{M}) p(\theta|\mathcal{M}) d\theta \quad (3)$$

which is also the denominator in Eq.2 (with \mathcal{M} implicit in the previous definition). This quantity is also known as *evidence*, or *prior predictive* distribution, since the likelihood is weighted by the prior (in contrast to a posterior predictive distribution where the likelihood would be weighted by the posterior). The marginal likelihood penalizes the prior complexity of a model and thus naturally embodies the principle of Occam’s razor (Jaynes, 2003). To compare two competing models, one can focus on the ratio between two marginal likelihoods, called a Bayes factor (BF):

$$\text{BF}_{ij} = \frac{p(X|\mathcal{M}_i)}{p(X|\mathcal{M}_j)} \quad (4)$$

which quantifies the relative evidence of model i over model j . Alternatively, if prior information about model plausibility is available, one can consider model posteriors $p(\mathcal{M}|X) \propto p(X|\mathcal{M}) p(\mathcal{M})$ and compute the posterior odds:

$$\frac{p(\mathcal{M}_i|X)}{p(\mathcal{M}_j|X)} = \frac{p(X|\mathcal{M}_i)}{p(X|\mathcal{M}_j)} \frac{p(\mathcal{M}_i)}{p(\mathcal{M}_j)} \quad (5)$$

which combine the relative evidence given by the BF with prior information in the form of prior odds.

Model Intractability

In order for cognitive models to be useful in practice, parameter estimation and model comparison should be feasible within reasonable time limits. As evident from their definitions, both Bayesian parameter estimation and model comparison depend on the likelihood function $p(X|\theta, \mathcal{M})$ which needs to be evaluated analytically or numerically for any triplet (\mathcal{M}, θ, X) .

When this is possible, standard Bayesian approaches for obtaining random draws from the posterior, such as Markov chain Monte Carlo (MCMC), or optimizing an approximate posterior, such as variational inference (VI), can be readily applied. However, when the likelihood function is not available in closed-form or too expensive to evaluate, standard methods no longer apply.

In fact, many interesting models from a variety of domains in cognitive science and psychology turn out to be intractable (Voss et al., 2019; B. Turner, Sederberg, & McClelland, 2016). This has precluded the wide exploration and application of these models, as researchers have often traded off complexity or neurocognitive plausibility for simplicity in order to make these models tractable. In the following, we discuss the most popular approach to inference with intractable models.

Simulation-Based Inference

Simulation-based methods leverage the generative property of mathematical models by treating a particular model as a *scientific simulator* from which synthetic data can be obtained given any configuration of the parameters. Simulation-based inference is common to many domains in science in general (Cranmer, Brehmer, & Louppe, 2019) and a variety of different approaches exist. These methods have also been dubbed *likelihood-free*, which is somewhat unfortunate, since the likelihood is implicitly defined by the generative process and sampling from the likelihood is realized through the stochastic simulator:

$$x_n \sim p(x|\theta, \mathcal{M}) \iff x_n = \mathcal{M}(\theta, \xi_n) \text{ with } \xi_n \sim p(\xi) \quad (6)$$

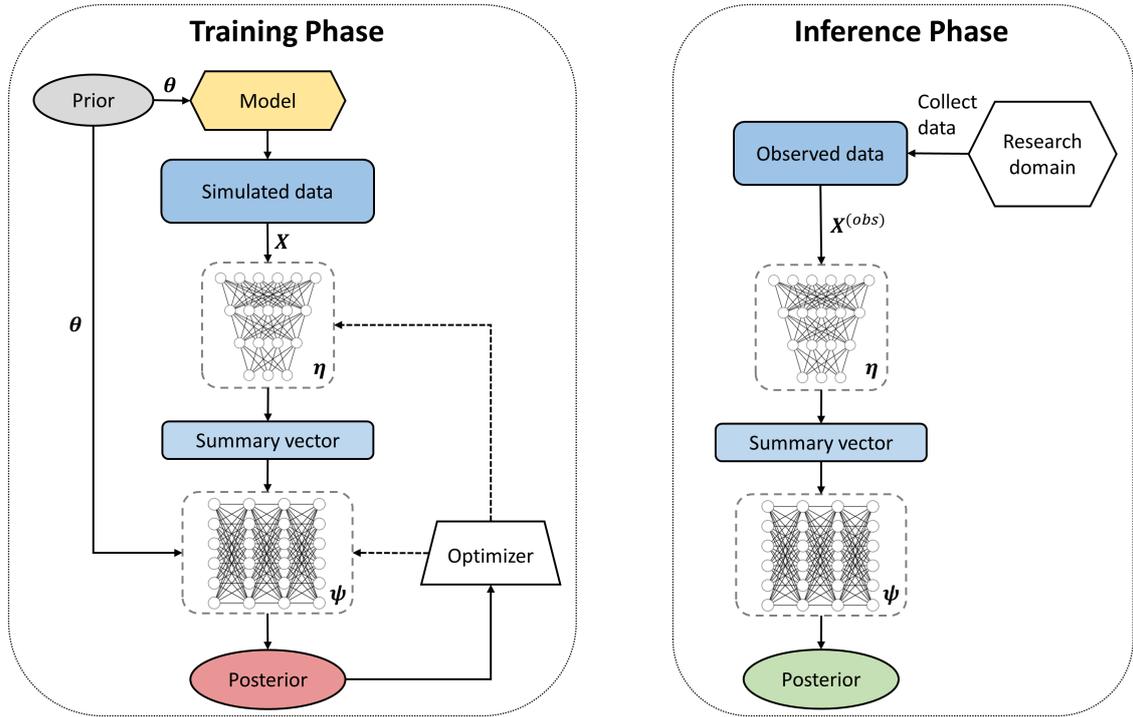
Different simulation-based methods differ mainly with respect to how they utilize the synthetic data to perform inference on real observed data (Cranmer et al., 2019). The utility of any simulation-based method depends on multiple factors, such as asymptotic guarantees, data utilization, efficiency, scalability, and software availability.

Approximate Bayesian computation (ABC) offers a standard set of theoretically sound methods for performing inference on intractable models (Cranmer et al., 2019). The core idea of ABC methods is to approximate the posterior by repeatedly sampling parameters from a proposal (prior) distribution and then generating a synthetic dataset by running the simulator with the sampled parameters. If the simulated dataset is sufficiently similar to an actually observed dataset, the corresponding parameters are retained as a sample from the desired posterior, otherwise rejected. However, in practice, ABC methods are notoriously inefficient and suffer from various problems, such as the *curse of dimensionality* or *curse of inefficiency* (Marin, Pudlo, Estoup, & Robert, 2018). More efficient methods employ various techniques to optimize sampling or correct potential biases.

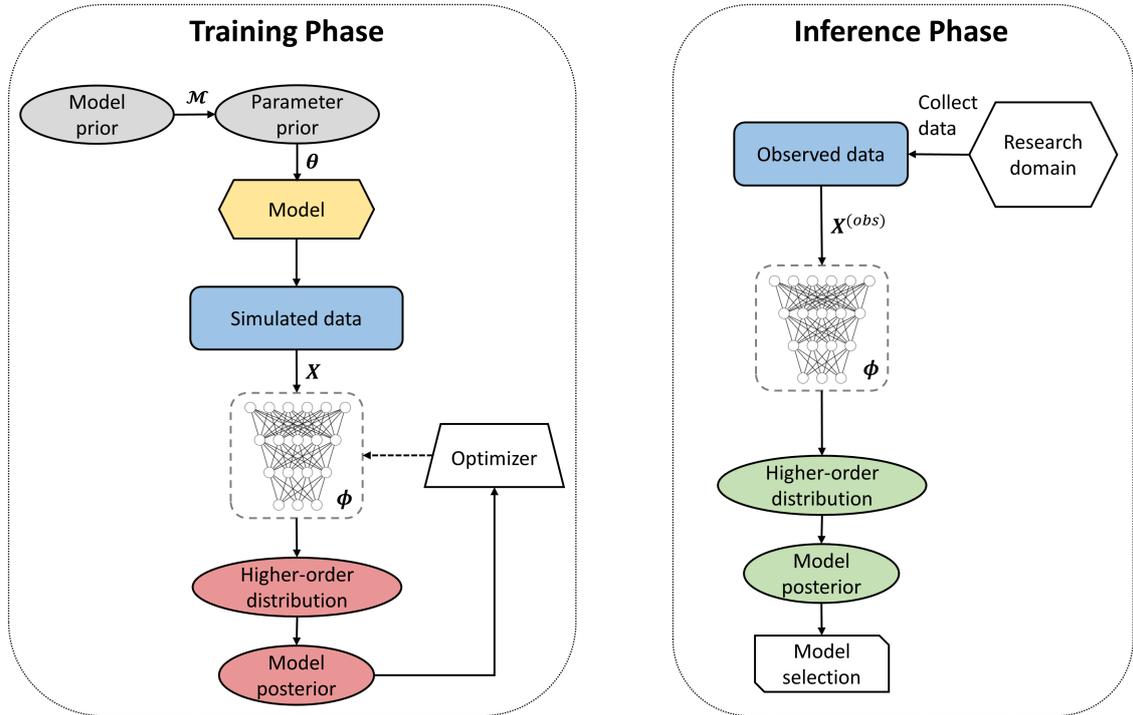
Recently, the scientific repertoire for simulation-based inference has been enhanced with ideas from deep learning and neural density estimation (NDE) in particular (Greenberg, Nonnenmacher, & Macke, 2019). These methods employ specialized neural network architectures which are trained with simulated data to perform efficient and accurate inference on previously intractable problems (Cranmer et al., 2019). NDE methods are rapidly developing and still largely underutilized in cognitive modeling, even though first applications to simulated (Radev, Mertens, Voss, Ardizzone, & Köthe, 2020; Radev, D’Alessandro, et al., 2020) as well as actual data (Wieschen, Voss, & Radev, 2020) exist.

Amortized Inference

The majority of simulation-based methods need to be applied to each dataset separately. This quickly becomes infeasible when multiple datasets are to be analysed and multiple candidate models are considered, since the expensive inference procedure needs to be repeated from scratch for each combination of dataset and model.



(a) Amortized parameter estimation



(b) Amortized model comparison

Figure 1: Graphical illustration of amortized parameter estimation and model comparison with different neural network estimators. **(a)** Amortized Bayesian parameter estimation with invertible neural networks (Radev, Mertens, et al., 2020). The left panel depicts the training phase in which the summary (f_η) and the inference network (f_ψ) are jointly optimized to approximate the true target posterior. The right panel depicts inference with already trained networks on observed data; **(b)** Amortized Bayesian model comparison with evidential neural networks (Radev, D’Alessandro, et al., 2020). The left panel depicts the training phase during which the evidential network f_ϕ is optimized to approximate the true model posteriors via a higher-order Dirichlet distribution. The right panel depicts inference with an already trained evidential network; the upfront training effort for both inference tasks is amortized over arbitrary numbers of datasets from a research domain.

In contrast, the concept of *amortized inference* refers to an approach which minimizes the cost of inference by separating the process into an expensive training (optimization) phase and a cheap inference phase which can be easily repeated for multiple datasets or models without computational overhead. Thus, the effort of training or optimization amortizes over repeated applications on multiple datasets or models. In some cases, the efficiency advantage of amortized inference becomes noticeable even for a few datasets (Radev, Mertens, et al., 2020; Radev, D’Alessandro, et al., 2020).

The field of amortized inference is rapidly growing and a variety of methods and concepts are currently being explored. For instance, *inference compilation* involves pre-training a neural network with simulations from a generative model and then using the network in combination with a probabilistic program to optimize sampling from the posterior (Le, Baydin, & Wood, 2016). The *pre-paid* estimation method (Mestdagh, Verdonck, Meers, Loossens, & Tuerlinckx, 2019) proceeds by creating a large grid of simulations which are reduced to summary statistics and stored on disk. Subsequent inference involves computing the nearest neighbors of an observed dataset in the pre-paid grid and interpolation. Sequential neural posterior estimation (SNPE) methods employ various iterative refinement schemes to transform a proposal distribution into the correct target posterior via expressive NDEs trained over multiple simulation rounds (Greenberg et al., 2019).

In line with these ideas, we recently proposed two general frameworks for amortized Bayesian parameter estimation and model comparison based on specialized neural network architectures (Radev, Mertens, et al., 2020; Radev, D’Alessandro, et al., 2020). In particular, these frameworks were designed to implement the following desirable properties:

- Fully amortized Bayesian inference for parameter estimation and model comparison of intractable models
- Asymptotic theoretical guarantees for sampling from the true parameter and model posteriors
- Learning maximally informative summary statistics directly from data instead of manual selection
- Scalability to high-dimensional problems through considerations regarding the probabilistic symmetry of the data
- Implicit preference for simpler models based purely on generative performance
- Online learning eliminating the need for storing large grids or reference tables
- Parallel computations and GPU acceleration applicable to both simulations, training/optimization, and inference

In the following, we describe our recently developed methods parameter estimation and model comparison in turn.

Amortized Parameter Estimation with Invertible Neural Networks

Recently, we proposed a novel amortization method based on invertible neural networks (Radev, Mertens, et al., 2020), which we dubbed *BayesFlow*. The method relies solely on simulations from a process model in order to learn and calibrate the full posterior over all possible parameter values and observed data patterns.

The BayesFlow method involves two separate neural networks trained jointly. A permutation invariant *summary network* is responsible for reducing an entire dataset X with a variable number N of *i.i.d.* observations¹ into a vector of *learned summary statistics*. Importantly, permutation invariant networks can deal with *i.i.d.* sequences of variable size and preserve their probabilistic symmetry. An *inference network*, implemented as an invertible neural network (Radev, Mertens, et al., 2020), is responsible for approximating the true posterior of model parameters given the output of the summary network. Invertible networks can perform asymptotically exact inference and scale well from simple low-dimensional problems to high-dimensional distributions with complex dependencies. During training, model parameters and synthetic datasets are generated on the fly and neural network parameters are adjusted via joint backpropagation (see Figure 1a, left panel, for a graphical illustration of the training phase).

Given a model and a prior over the model parameters, the goal is thus to train a conditional invertible neural network f_ψ with adjustable parameters ψ together with a summary network f_η with adjustable parameters η . These networks jointly learn an approximate posterior $p_\psi(\theta | f_\eta(X))$ over the relevant parameters for arbitrary numbers of datasets and dataset sizes N , as long as they share the same data structure. To achieve this, the networks minimize the Kullback-Leibler (KL) divergence between the true and the approximate posterior:

$$\min_{\psi, \eta} \mathbb{KL}(p(\theta | X) || p_\psi(\theta | f_\eta(X))) \quad (7)$$

Utilizing the fact that we have access to the joint distribution $p(\theta, X) = p(\theta)(X | \theta)$ via the simulator, we minimize the KL divergence in expectation over all possible datasets that can be generated given the prior and the model, resulting in the following optimization criterion:

$$\min_{\psi, \eta} \mathbb{E}_{p(\theta, x)} [-\log p_\psi(\theta | f_\eta(X))] \quad (8)$$

In practice, we approximate the criterion via its Monte Carlo (MC) estimate, since we can simulate theoretically infinite amounts of data and can easily evaluate $p_\psi(\theta | f_\eta(X))$ due to our invertible architecture. In case of perfect convergence of the networks, the summary network outputs *sufficient summary statistics* and the inference network samples from the true posterior (Radev, Mertens, et al., 2020). Importantly,

¹Note, that the *i.i.d.* assumption is not a necessary condition for the method to work, but used here only to simplify the discussion.

once the networks have been trained with sufficient amounts of simulated data, they can be stored and applied for inference on multiple datasets from a research domain (see Figure 1a, right panel).

Amortized Model Comparison with Evidential Neural Networks

In another recent work (Radev, D’Alessandro, et al., 2020), we explored a framework for Bayesian model comparison on intractable models via *evidential neural networks*. We proposed to train a permutation invariant classifier network on simulated data from multiple models. The goal of this network is to approximate posterior model probabilities as accurately as possible. To achieve this, the network is trained to output the parameters of a higher-order probability distribution (parameterized as a Dirichlet distribution) over the model probabilities themselves, which quantifies the uncertainty in model probability estimates. Thus, for a classifier network with parameters ϕ , the higher-order posterior distribution over model probabilities is given by:

$$\text{Dir}(\pi | \alpha_\phi(X)) = \frac{1}{B(\alpha_\phi(X))} \prod_{j=1}^J \pi^{\alpha_\phi(X)_j - 1} \quad (9)$$

where $\alpha_\phi(X)$ denotes the vector of *concentration parameters* obtained by the network for a dataset X and $B(\cdot)$ is the multivariate *beta* function. The mean of this Dirichlet distribution can be used as a best estimate for the posterior model probabilities:

$$p_\phi(\mathcal{M} | X) = \frac{\alpha_\phi(X)}{\sum_{j=1}^J \alpha_\phi(X)_j} \quad (10)$$

Additionally, its variance can be interpreted as the epistemic uncertainty surrounding the actual evidence which the data provide for model comparison.

For training the network, we again utilize the fact that we have access to the joint distribution $p(\mathcal{M}, \theta, X)$ via simulations (see Figure 1b, left panel). Our optimization criterion is:

$$\min_{\phi} \mathbb{E}_{p(\mathcal{M}, \theta, X)} [\mathcal{L}(p_\phi(\mathcal{M} | X), \mathcal{M})] \quad (11)$$

where $\mathcal{L}(\cdot, \cdot)$ is a *strictly proper* loss function (Gneiting & Raftery, 2007), \mathcal{M} is the true model index and the data X implicitly depend on θ . In practice, we approximate this expectation via draws from the joint distribution available via the simulator. Optimization of a strictly proper criterion, asymptotic convergence implies that the mean of the Dirichlet distribution represents the true model posteriors. Moreover, our simulation-based approach implicitly captures a preference for simpler models (Occam’s razor), since simpler models will tend to generate more similar datasets. As a consequence, when such datasets are plausible under multiple models, the comparably simpler models will be more probable.

As with parameter estimation, once the evidence network has been trained on simulated data from the candidate models, it can be applied to multiple upcoming observations from a research domain (see Figure 1b, right panel).

Example Applications

In the following, we will present two applications of amortized Bayesian parameter estimation to a recently proposed and intractable evidence accumulation model (EAM). The first illustrative application is a simulation study aimed at quantifying parameter recovery as a function of data set size. Such simulations are especially useful for planing experiments but usually too costly to perform in complex modeling scenarios. The second application is concerned with parameter estimation on real data and serves as an illustration on how researchers might utilize amortized Bayesian inference with a pre-trained density estimator in practice.

EAMs are a popular class of models in psychology and cognitive science, as they allow a model-based analysis of response time (RT) distributions. Here, we will consider a Lévy flight model (LFM) with a non-Gaussian noise assumption (Voss et al., 2019; Wieschen et al., 2020) as an example. The Lévy flight process is driven by the following stochastic ordinary differential equation (ODE):

$$dx_c = v_c dt + \xi dt^{1/\alpha} \quad (12)$$

$$\xi \sim \text{AlphaStable}(\alpha, 0, 1, 0) \quad (13)$$

where dx_c denotes accumulated cognitive evidence in condition c , v_c denotes the average speed of information accumulation (drift), and α controls how heavy the tails of the noise distribution are (i.e., smaller values increase the probability of outliers in the accumulation process). Further parameters of the model are: a decision threshold (a) which reflects the amount of information needed for selecting a response; a starting point (z_r) indicative of response biases; and a non-decision time (t_0) reflecting additive encoding and motor process. Since the relationship of the α parameter to the standard parameters of the classical diffusion model (Ratcliff, Thapar, Gomez, & McKoon, 2004) has not been previously investigated, we focus on quantifying posterior correlations in the real data application.

Simulation Example

As a first example, consider a simulated RT experiment with four conditions. How many trials are needed for accurate parameter recovery? To answer this question, we can simulate multiple experiments with varying number of trials per participant (N) and then compute some discrepancy between ground-truth parameters and their estimates. However, since the model is intractable, such a simulation scenario is not feasible with non-amortized methods, which would need weeks on standard machines (Voss et al., 2019). However, using the *BayesFlow* method (Figure 1a), we can train the networks with simulated datasets and vary the number of trials during each simulation. Such a training takes approximately one day on a standard laptop equipped with an NVIDIA[®] GTX1060 graphics card. Subsequent inference is then very cheap, as amortized parameter estimation on 500 simulated participants takes less than 2 seconds.

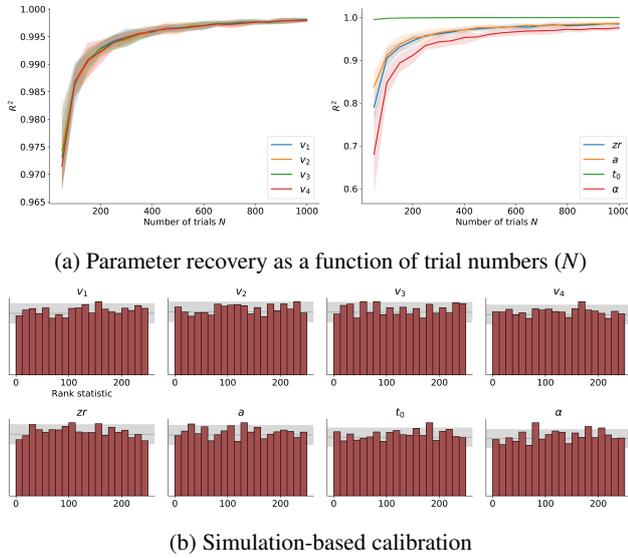


Figure 2: Simulation results. (a) The left panel depicts parameter recovery of the four drift rate parameters as a function of trial numbers per participant N . The right panel depicts recovery of the other four parameters. Posterior means are used as summaries of the full posteriors and shaded regions represent bootstrap 95% confidence intervals. (b) The panel depicts simulation-based calibration (SBC) results at $N = 800$ as a validation check for the correctness of the full posteriors.

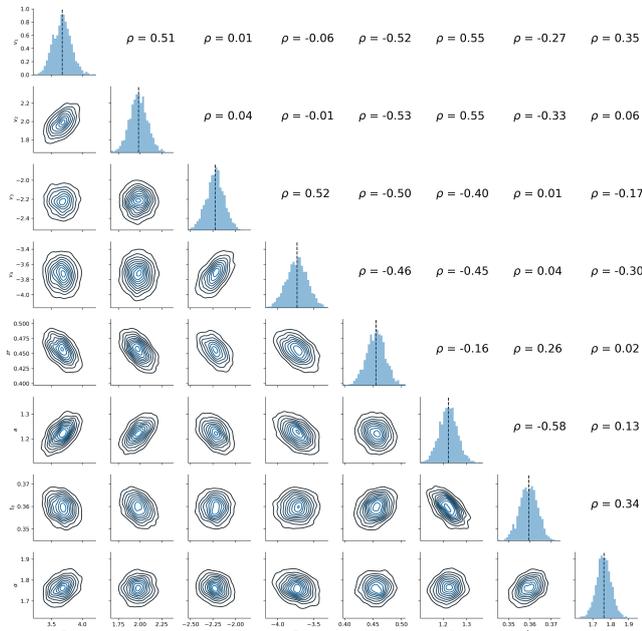


Figure 3: Example full posteriors and bivariate posterior correlations obtained from data of one participant in the long LDT via amortized Bayesian inference. Dashed lines on the main diagonal indicate posterior means.

We visualize the results by plotting the average R^2 metric obtained from fitting the LFM model to 300 simulated participants at different N between 50 and 1000 (see Figure 2a). Notably, recovery of the ground-truth parameters via posterior means is nearly perfect at higher trial numbers.

As a validation tool for visually detecting systematic biases in the approximate posteriors, we can also cheaply apply simulation-based calibration (SBC) and inspect the rank statistic of the posterior samples for uniformity (Talts, Betancourt, Simpson, Vehtari, & Gelman, 2018). Results from applying SBC to 5000 simulated participants at $N = 800$ are depicted in Figure 2b. Indeed, we confirm that no pronounced issues across all marginal posteriors are present.

Real Data Example

We can also apply the same networks from the previous simulation example for fully Bayesian inference on real data. Here, we fit the LFM model to previously unpublished data from eleven participants performing a long ($N = 800$ per condition) lexical decision task (LDT). Since the task had a 2×2 design, with a factor for *difficulty* (hard vs. easy), and a factor for *stimulus type* (word vs. non-word), we can assume a different drift rate for each design cell.

Applying the pre-trained networks, we immediately obtain samples from a full posterior over model parameters for each participant. Using the estimated posteriors, we can then test hypotheses about particular parameter values, compute individual differences, or compare means between conditions in a Bayesian way. Furthermore, we can analyze posterior correlations at an individual level and investigate task-dependent relationships between the α parameter and other parameters (see Figure 3 for results obtained from a single participant).

Across participants, α displays only small posterior correlations with drift rates as well as small posterior correlations with threshold and non-decision time parameters (mean $r < 0.5$ across all parameters of the standard diffusion model). These results provide first evidence that the α parameter can indeed be decoupled from other model parameters and possibly indicates a separate decision process.

Since the goal of this application was solely to illustrate a typical use case for amortized Bayesian inference, future research should focus on extensive external validation of the LFM model as well as proposing a neurocognitively plausible interpretation for the α parameter.

Outlook

The purpose of this work was to introduce the main ideas behind amortized Bayesian inference methods for simulation-based parameter estimation and model comparison. Although these methods come with promising theoretical guarantees and clear practical advantages, their utility for cognitive modeling is just beginning to be explored. Moreover, there are still many open questions and avenues for future research.

First, a systematic investigation of a potential *amortization gap* in certain practical application seems warranted. An amortization gap refers to a drop in estimation accuracy due

to the fact that we are relying on a single set of neural network parameters for solving an inference problem globally, instead of performing per-dataset optimization. Even though we have not observed such a scenario in our applications and simulations, this behavior might occur when the neural network estimators are not expressive enough to represent complex posterior distributions.

Second, there are still little systematic guidelines on how to best design and tune the neural network architectures so as to perform optimally across a variety of parameter estimation and model comparison tasks. Even though neural density estimation methods outperform standard ABC methods on multiple metrics and in various contexts, there is certainly room for improvement. Black-box optimization methods for hyperparameter tuning, such as Bayesian optimization or active inference (Snoek, Larochelle, & Adams, 2012), might facilitate additional performance gains and reduce potentially suboptimal architectural choices.

Finally, user-friendly software for applying Bayesian amortization methods *out-of-the-box* is still largely in its infancy. Developing and maintaining such software is a crucial future goal for increasing the applicability and usability of novel simulation-based methods.

Conclusion

We hope that the inference architectures discussed in this work will spur the interest of cognitive modelers from various domains. We believe that such architectures can greatly enhance model-based analysis in cognitive science and psychology. By leaving subsidiary tractability considerations to powerful end-to-end algorithms, researchers can focus more on the task of model development and evaluation to further improve our understanding of cognitive processes.

Acknowledgments

This research was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation; grant number GRK 2277 "Statistical Modeling in Psychology"). We thank the Technology Industries of Finland Centennial Foundation (grant 70007503; Artificial Intelligence for Research and Development) for partial support of this work.

References

- Cranmer, K., Brehmer, J., & Louppe, G. (2019). The frontier of simulation-based inference. *arXiv preprint arXiv:1911.01429*.
- Fontanesi, L., Gluth, S., Spektor, M. S., & Rieskamp, J. (2019). A reinforcement learning diffusion decision model for value-based decisions. *Psychonomic bulletin & review*, 26(4), 1099–1121.
- Gneiting, T., & Raftery, A. E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477), 359–378.
- Greenberg, D. S., Nonnenmacher, M., & Macke, J. H. (2019). Automatic posterior transformation for likelihood-free inference. *arXiv preprint arXiv:1905.07488*.
- Jaynes, E. T. (2003). *Probability theory: The logic of science*. Cambridge university press.
- Le, T. A., Baydin, A. G., & Wood, F. (2016). Inference compilation and universal probabilistic programming. *arXiv preprint arXiv:1610.09900*.
- Lynch, S. M., & Western, B. (2004). Bayesian posterior predictive checks for complex models. *Sociological methods & research*, 32(3), 301–335.
- Marin, J.-M., Pudlo, P., Estoup, A., & Robert, C. (2018). *Likelihood-free model choice*. Chapman and Hall/CRC Press Boca Raton, FL.
- Mestdagh, M., Verdonck, S., Meers, K., Loossens, T., & Tuerlinckx, F. (2019). Prepaid parameter estimation without likelihoods. *PLoS computational biology*, 15(9), e1007181.
- Myung, J. I., Montenegro, M., & Pitt, M. A. (2007). Analytic expressions for the bcdmem model of recognition memory. *Journal of Mathematical Psychology*, 51(3), 198–204.
- Radev, S. T., D'Alessandro, M., Bürkner, P.-C., Mertens, U. K., Voss, A., & Köthe, U. (2020). Amortized bayesian model comparison with evidential deep learning. *arXiv preprint arXiv:2004.10629*.
- Radev, S. T., Mertens, U. K., Voss, A., Ardizzone, L., & Köthe, U. (2020). Bayesflow: Learning complex stochastic models with invertible neural networks. *arXiv preprint arXiv:2003.06281*.
- Ratcliff, R., Thapar, A., Gomez, P., & McKoon, G. (2004). A diffusion model analysis of the effects of aging in the lexical-decision task. *Psychology and aging*, 19(2), 278.
- Snoek, J., Larochelle, H., & Adams, R. P. (2012). Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems* (pp. 2951–2959).
- Stout, J. C., Busemeyer, J. R., Lin, A., Grant, S. J., & Bonson, K. R. (2004). Cognitive modeling analysis of decision-making processes in cocaine abusers. *Psychonomic bulletin & review*, 11(4), 742–747.
- Talts, S., Betancourt, M., Simpson, D., Vehtari, A., & Gelman, A. (2018). Validating bayesian inference algorithms with simulation-based calibration. *arXiv preprint arXiv:1804.06788*.
- Turner, B., Sederberg, P., & McClelland, J. (2016). Bayesian analysis of simulation-based models. *Journal of Mathematical Psychology*, 72, 191–199.
- Turner, B. M., Forstmann, B. U., Steyvers, M., et al. (2019). *Joint models of neural and behavioral data*. Springer.
- Voss, A., Lerche, V., Mertens, U., & Voss, J. (2019). Sequential sampling models with variable boundaries and non-normal noise: A comparison of six models. *Psychonomic bulletin & review*, 26(3), 813–832.
- Wieschen, E. M., Voss, A., & Radev, S. (2020). Jumping to conclusion? a lévy flight model of decision making. *TQMP*, 16(2), 120–132.

APPENDIX A4 - MANUSCRIPT 4

Manuscript 4: A Bayesian brain model of adaptive behavior: an application to the Wisconsin Card Sorting Task.



A Bayesian brain model of adaptive behavior: an application to the Wisconsin Card Sorting Task

Marco D'Alessandro¹, Stefan T. Radev², Andreas Voss² and Luigi Lombardi¹

¹Department of Psychology and Cognitive Science, University of Trento, Rovereto, Italy

²Institute of Psychology, Heidelberg University, Heidelberg, Germany

ABSTRACT

Adaptive behavior emerges through a dynamic interaction between cognitive agents and changing environmental demands. The investigation of information processing underlying adaptive behavior relies on controlled experimental settings in which individuals are asked to accomplish demanding tasks whereby a hidden regularity or an abstract rule has to be learned dynamically. Although performance in such tasks is considered as a proxy for measuring high-level cognitive processes, the standard approach consists in summarizing observed response patterns by simple heuristic scoring measures. With this work, we propose and validate a new computational Bayesian model accounting for individual performance in the Wisconsin Card Sorting Test (WCST), a renowned clinical tool to measure set-shifting and deficient inhibitory processes on the basis of environmental feedback. We formalize the interaction between the task's structure, the received feedback, and the agent's behavior by building a model of the information processing mechanisms used to infer the hidden rules of the task environment. Furthermore, we embed the new model within the mathematical framework of the Bayesian Brain Theory (BBT), according to which beliefs about hidden environmental states are dynamically updated following the logic of Bayesian inference. Our computational model maps distinct cognitive processes into separable, neurobiologically plausible, information-theoretic constructs underlying observed response patterns. We assess model identification and expressiveness in accounting for meaningful human performance through extensive simulation studies. We then validate the model on real behavioral data in order to highlight the utility of the proposed model in recovering cognitive dynamics at an individual level. We highlight the potentials of our model in decomposing adaptive behavior in the WCST into several information-theoretic metrics revealing the trial-by-trial unfolding of information processing by focusing on two exemplary individuals whose behavior is examined in depth. Finally, we focus on the theoretical implications of our computational model by discussing the mapping between BBT constructs and functional neuroanatomical correlates of task performance. We further discuss the empirical benefit of recovering the assumed dynamics of information processing for both clinical and research practices, such as neurological assessment and model-based neuroscience.

Submitted 10 July 2020
Accepted 16 October 2020
Published 30 November 2020

Corresponding author
Marco D'Alessandro,
marco.dalessandro@unitn.it

Academic editor
Marco Tullio Liuzza

Additional Information and
Declarations can be found on
page 27

DOI 10.7717/peerj.10316

© Copyright
2020 D'Alessandro et al.

Distributed under
Creative Commons CC-BY 4.0

OPEN ACCESS

Subjects Computational Biology, Neuroscience, Psychiatry and Psychology

Keywords Adaptive behavior, Bayesian brain, Cognitive modeling, Wisconsin card sorting test, Information processing, Information theory, Belief updating

INTRODUCTION

Computational models of cognition provide a way to formally describe and empirically account for mechanistic, process-based theories of adaptive cognitive functioning (Sun, 2009; Cooper et al., 1996; Lee & Wagenmakers, 2014). A foundational theoretical framework for describing functional characteristics of neurocognitive systems has recently emerged under the hood of Bayesian brain theories (Knill & Pouget, 2004; Friston, 2010). Bayesian brain theories owe their name to their core assumption that neural computations resemble the principles of Bayesian statistical inference.

In a Bayesian theoretical framework, cognitive agents interact with an uncertain and changeable sensory environment. This requires a cognitive system to infer sensory contingencies based on an internal generative model of the environment. Such a generative model represents subjective hypotheses, or beliefs, about the causal structure of events in the environment (Friston, 2005; Knill & Pouget, 2004) and forms a basis for adaptive behavior. It is assumed that internal beliefs are constantly updated and refined to match the current state of the world as new observations become available. The core idea behind the Bayesian brain hypothesis is that computational mechanisms underlying such an internal belief updating follow the logic of Bayesian probability theory. In this respect, information about the external world provided by sensory inputs is represented as a conditional probability distribution over a set of environmental states. Consequently, the brain relies on this probabilistic representation of the world to infer the most likely environmental causes (states) which generate those inputs, and such a process follows the computational principles of Bayesian inference (Friston & Kiebel, 2009; Friston, 2010; Buckley et al., 2017).

To clarify this concept, consider a simple example of a perceptual task in which a cognitive agent is required to judge whether an item depicted on a flat plane is concave or convex. Its judgment is based solely on the basis of a set of observed perceptual features, such as, shape, orientation, texture and brightness. Here, the concave-to-convex gradient entails the set of environmental states which must be inferred. The internal generative model of the agent codifies beliefs about how different degrees of convexity might give rise to certain configurations of perceptual inputs. From a Bayesian perspective, the problem is solved by *inverting* the generative model of the environment in order to turn assumptions about how environmental states generate sensory inputs into beliefs about the most likely states (e.g., degree of convexity) given the available sensory information.

Potentially, there are no limitations regarding the complexity of environmental settings (e.g., items and rules in experimental tasks) and cognitive processes to be described in light of the Bayesian brain framework. Indeed, the latter has proven to be a consistent computational modeling paradigm for the investigation of a variety of neurocognitive mechanisms, such as motor control (Friston et al., 2010), oculomotor dynamics (Friston et al., 2012), object recognition (Kersten, Mamassian & Yuille, 2004), attention (Feldman & Friston, 2010), perceptual inference (Petzschner, Glasauer & Stephan, 2015; Knill & Pouget, 2004), multisensory integration (Körding et al., 2007), as well as for providing a foundational theoretical account of general neural systems' functioning (Lee & Mumford, 2003; Friston, 2005; Friston, 2003) and complex clinical scenarios such as Schizophrenia

(Stephan, Baldeweg & Friston, 2006), and Autistic Spectrum Disorder (Haker, Schneebeli & Stephan, 2016; Lawson, Rees & Friston, 2014). For this reason, such a modeling approach might provide a comprehensive and unified framework under which several cognitive impairments can be measured and understood in the light of a general process-based theory of neural functioning.

In this work, we address the challenging problem of modeling adaptive behavior in a dynamic environment. The empirical assessment of adaptive functioning often relies on dynamic reinforcement learning scenarios which require participants to adapt their behavior during the unfolding of a (possibly) demanding task. Typically, these tasks are designed with the aim to figure out how adaptive behavior unfolds through multiple trials as participants observe certain environmental contingencies, take actions, and receive feedback based on their actions. From a Bayesian theoretical perspective, optimal performance in such adaptive experimental paradigms require that agents infer the probabilistic model underlying the hidden environmental states. Since these models usually change as the task progresses, agents, in turn, need to adapt their inferred model, in order to take optimal actions.

Here, we propose and validate a computational Bayesian model which accounts for the dynamic behavior of cognitive agents in the Wisconsin Card Sorting Test (WCST; Berg, 1948; Heaton, 1981), which is perhaps the most widely adopted neuropsychological setting employed to investigate adaptive functioning. Due to its structure, the WCST can account for executive components underlying observed behavior, such as set-shifting, cognitive flexibility and impulsive response modulation (Bishara et al., 2010; Alvarez & Emory, 2006). For this reason, we consider the WCST as a fundamental paradigm for investigating adaptive behavior from a Bayesian perspective.

The environment of the WCST consists of a target and a set of stimulus cards with geometric figures which vary according to three perceptual features. The WCST requires participants to infer the correct classification rule by trial and error using the examiner's feedback. The feedback is thought to carry a positive or negative information signaling the agent whether the immediate action was appropriate or not. Modeling adaptive behavior in the WCST from a Bayesian perspective is straightforward, since observable actions emerge from the interaction between the internal probabilistic model of the agent and a set of discrete environmental states.

Performance in WCST is usually measured via a rough summary metric such as the number of correct/incorrect responses or pre-defined psychological scoring criteria (see for instance Heaton, 1981). These metrics are then used to infer the underlying cognitive processes involved in the task. A major shortcoming of this approach is that it simply assumes the cognitive processes to be inferred without specifying an explicit *process model*. Moreover, summary measures do not utilize the full information present in the data, such as trial-by-trial fluctuations or various interesting agent-environment interactions. For this reason, crude scoring measures are often insufficient to disentangle the dynamics of the relevant cognitive (sub)processes involved in solving the task. Consequently, an entanglement between processes at the metric level can prevent us from answering interesting research questions about aspects of adaptive behavior.

In our view, a sound computational account for adaptive behavior in the WCST needs to provide at least a quantitative measure of effective belief updating about the environmental states at each trial. This measure should be complemented by a measure of how feedback-related information influences behavior. The first measure should account for the integration of meaningful information. In other words, it should describe how prior beliefs about the current environmental state change after an observation has been made. The second measure should account for signaling the (im)probability of observing a certain environmental configuration (e.g., an (un)expected feedback given a response) ([Schwartenbeck, FitzGerald & Dolan, 2016](#)).

Indeed, recent studies suggest that the meaningful information content and the pure unexpectedness of an observation are processed differently at the neural level. Moreover, such disentanglement appears to be of crucial importance to the understanding of how new information influences adaptive behavior ([Nour et al., 2018](#); [Schwartenbeck, FitzGerald & Dolan, 2016](#); [O'Reilly et al., 2013](#)). Inspired by these results and previous computational proposals ([Koechlin & Summerfield, 2007](#)), we integrate these different information processing aspects into the current model from an information-theoretic perspective.

Our computational cognitive model draws heavily on the mathematical frameworks of Bayesian probability theory and information theory ([Sayood, 2018](#)). First, it provides a parsimonious description of observed data in the WCST via two neurocognitively meaningful parameters, namely, *flexibility* and *information loss* (to be motivated and explained in the next section). Moreover, it captures the main response patterns obtainable in the WCST via different parameter configurations. Second, we formulate a functional connection between cognitive parameters and underlying information processing mechanisms related to belief updating and prediction formation. We formalize and distinguish between *Bayesian surprise* and *Shannon surprise* as the main mechanisms for adaptive belief updating. Moreover, we introduce a third quantity, which we named predictive *Entropy* and which quantifies an agent's subjective uncertainty about the current internal model. Finally, we propose to measure these quantities on a trial-by-trial basis and use them as a proxy for formally representing the dynamic interplay between agents and environments.

The rest of the paper is organized as follows. First, the WCST is described in more detail and a mathematical representation of the new Bayesian computational model is provided. Afterwards, we explore the model's characteristics through simulations and perform parameter recovery on simulated data using a powerful Bayesian deep neural network method ([Radev et al., 2020](#)). We then apply the model to real behavioral data from an already published dataset. Finally, we discuss the results as well as the main strengths and limitations of the proposed model.

THE WISCONSIN CARD SORTING TEST

In a typical WCST ([Heaton, 1981](#); [Berg, 1948](#)), participants learn to pay attention and respond to relevant stimulus features, while ignoring irrelevant ones, as a function of

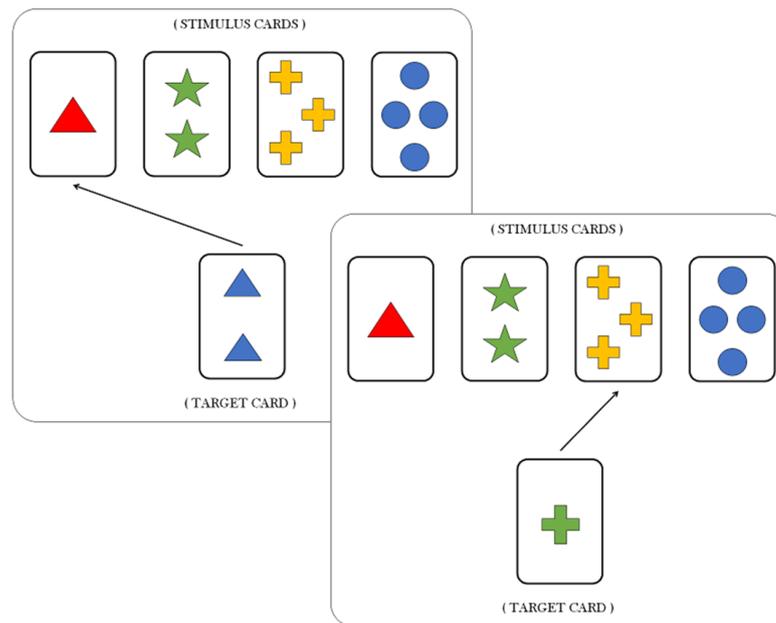


Figure 1 Suppose that the current sorting rule is the feature shape. The target card in the first trial (left box) contains two blue triangles. A correct response requires that the agent matches the target card with the stimulus card containing the single triangle (arrow represents the correct choice), regardless of the features color and number. The same applies for the second trial (right box) in which matching the target card with the stimulus card containing three yellow crosses is the correct response.

Full-size DOI: [10.7717/peerj.10316/fig-1](https://doi.org/10.7717/peerj.10316/fig-1)

experimental feedback. In particular, Individuals are asked to match a target card with one of four stimulus cards according to a proper sorting principle, or sorting rule. Each card depicts geometric figures that vary in terms of three features, namely, color (red, green, blue, yellow), shape (triangle, star, cross, circle) and number of objects (1, 2, 3 and 4). For each trial, the participant is required to identify the sorting rule which is valid for that trial, that is, which of the three feature has to be considered as a criterion to matching the target card with the right stimulus card (see Fig. 1). Notice that both features and sorting rules refer to the same concept. However, the feature still codifies a property of the card, whilst the sorting rule refers to the particular feature which is valid for the current trial.

Each response in the WCST is followed by a feedback informing the participant if his/her response is correct or incorrect. After some fixed number of consecutive responses, the sorting rule is changed by the experimenter without warning, and participants are required to infer the new sorting rule. Clearly, the most adaptive response would be to explore the remaining possible rules. However, participants sometimes would persist responding according to the old rule and produce what is called a *perseverative response*.

METHODS

The model

The core idea behind our computational framework is to encode the concept of *belief* into a generative probabilistic model of the environment. Belief updating then corresponds to recursive Bayesian updating of the internal model based on current and past interactions between the agent and its environment. Optimal or sub-optimal actions are selected according to a well specified or a misspecified internal model and, in turn, cause perceptible changes in the environment.

We assume that the cognitive agent aims to infer the *true hidden state* of the environment by processing and integrating sensory information from the environment. Within the context of the WCST, the hidden environmental states might change as a function of both the structure of the task and the (often sub-optimal) behavioral dynamics, so the agent constantly needs to rely on environmental feedback and own actions to infer the current state. We assume that the agent maintains an internal probability distribution over the states at each individual trial of the WCST. The agent then updates this distribution upon making new observations. In particular, the hidden environmental states to be inferred are the three features, $s_t \in \{1, 2, 3\}$, which refer the three possible sorting rules in the task environment such that 1: color, 2: shape and 3: number of objects. The posterior probability of the states depends on an observation vector $x_t = (a_t, f_t)$, which consists of the pair of agent's response $a_t \in \{1, 2, 3, 4\}$, codifying the action of choosing deck 1, 2, 3 or 4, and received feedback $f_t \in \{0, 1\}$, referring to the fact that a given response results in a failure (0) or in a success (1), in a given trial $t = 0, \dots, T$. The discrete response a_t represents the stimulus card indicator being matched with a target card at trial t . We denote a sequence of observations as $x_{0:t} = (x_0, x_1, \dots, x_t) = ((a_0, f_0), (a_1, f_1), (a_2, f_2), \dots, (a_t, f_t))$ and set $x_0 = \emptyset$ in order to indicate that there are no observations at the onset of the task. Thus, trial-by-trial belief updating is recursively computed according to Bayes' rule:

$$p(s_t | x_{0:t}) = \frac{p(x_t | s_t, x_{0:t-1})p(s_t | x_{0:t-1})}{p(x_t | x_{0:t-1})}. \quad (1)$$

Accordingly, the agent's posterior belief about the task-relevant features s_t after observing a sequence of response-feedback pairs $x_{0:t}$ is proportional to the product of the likelihood of observing a particular response-feedback pair and the agent's prior belief about the task-relevant feature in the current trial. The likelihood of an observation is computed as follows:

$$p(x_t | s_t, x_{0:t-1}) = \frac{f_t p(a_t | s_t = i) + (1 - f_t)(1 - p(a_t | s_t = i))}{f_t \sum_j p(a_t | s_t = j) + (1 - f_t) \sum_j (1 - p(a_t | s_t = j))} \quad (2)$$

where $j = 1, 2, 3$ and $p(a_t | s_t = i)$ indicates the probability of a matching between the target and the stimulus card assumed that the current feature is i . Here, we assume the likelihood of a current observation to be independent from previous observations without loss of generality, that is:

$$p(x_t | s_t, x_{0:t-1}) = p(x_t | s_t).$$

The prior belief for a given trial t is computed based on the posterior belief generated in the previous trial, $p(s_{t-1}|x_{0:t-1})$, and the agent's belief about the probability of transitions between the hidden states, $p(s_t|s_{t-1})$. The prior belief can also be considered as a predictive probability over the hidden states. The predictive distribution for an upcoming trial t is computed according to the Chapman–Kolmogorov equation:

$$p(s_{t+1} = k|x_{0:t}) = \sum_{i=1}^3 p(s_{t+1} = k|s_t = i, \Gamma(t))p(s_t = i|x_{0:t}) \quad (3)$$

where $\Gamma(t)$ represents a stability matrix describing transitions between the states (to be explained shortly). Thus, the agent combines information from the updated belief (posterior distribution) and the belief about the transition properties of the environmental states to predict the most probable future state. The predictive distribution represents the internal model of the cognitive agent according to which actions are generated.

The stability matrix $\Gamma(t)$ encodes the agent's belief about the probability of states being stable or likely to change in the next trial. In other words, the stability matrix reflects the cognitive agent's internal representation of the dynamic probabilistic model of the task environment. It is computed on each trial based on the response-feedback pair, x_t , and a matching signal, m_t , which are observed.

The matching signal m_t is a vector informing the cognitive agent which features are currently relevant (meaningful), such that $m_t^{(i)} = 1$ when a positive feedback is associated with a response implying feature $s_t = i$, and $m_t^{(i)} = 0$ otherwise. Note, that the matching signal is not a free parameter of the model, but is completely determined by the task contingencies. The matching signal vector allows the agent to compute the *state activation level* $\omega_t^{(i)} \in [0, 1]$ for the hidden state $s_t = i$, which provides an internal measure of the (accumulated) evidence for each hidden state at trial t . Thus, the activation levels of the hidden states are represented by a vector ω_t . The stability matrix is a square and asymmetric matrix related to hidden state activation levels such that:

$$\Gamma(t) = \begin{bmatrix} \omega_t^{(1)} & \frac{1}{2}(1 - \omega_t^{(1)}) & \frac{1}{2}(1 - \omega_t^{(1)}) \\ \frac{1}{2}(1 - \omega_t^{(2)}) & \omega_t^{(2)} & \frac{1}{2}(1 - \omega_t^{(2)}) \\ \frac{1}{2}(1 - \omega_t^{(3)}) & \frac{1}{2}(1 - \omega_t^{(3)}) & \omega_t^{(3)} \end{bmatrix} \quad (4)$$

where the entries $\Gamma_{ii}(t)$ in the main diagonal represent the elements of the activation vector ω_t , and the non-diagonal elements are computed so as to ensure that rows sum to 1. The state activation vector is computed in each trial as follows:

$$\begin{bmatrix} \omega_t^{(1)} \\ \omega_t^{(2)} \\ \omega_t^{(3)} \end{bmatrix} = f_t \omega_{t-1}^\delta \begin{bmatrix} m_t^{(1)} \\ m_t^{(2)} \\ m_t^{(3)} \end{bmatrix} + \lambda \left[(1 - f_t) \omega_{t-1}^\delta \begin{bmatrix} 1 - m_t^{(1)} \\ 1 - m_t^{(2)} \\ 1 - m_t^{(3)} \end{bmatrix} \right] \begin{bmatrix} \omega_{t-1}^{(1)} \\ \omega_{t-1}^{(2)} \\ \omega_{t-1}^{(3)} \end{bmatrix}. \quad (5)$$

This equation reflects the idea that state activations are simultaneously affected by the observed feedback, f_t , and the matching signal vector, m_t . However, the matching signal vector conveys different information based on the current feedback. Matching a target card

with a stimulus card makes a feature (or a subset of features) informative for a specific state. The vector m_t contributes to increase the activation level of a state if the feature is informative for that state when a positive feedback is received, as well as to decrease the activation level when a negative feedback is received.

The parameter $\lambda \in [0, 1]$ modulates the efficiency to disengage attention to a given state-activation configuration when a negative feedback is processed. We therefore term this parameter *flexibility*. We also assume that information from the matching signal vector can degrade by slowing down the rate of evidence accumulation for the hidden states. This means that the matching signal vector can be re-scaled based on the current state activation level. The parameter $\delta \in [0, 1]$ is introduced to achieve this re-scaling. When $\delta = 0$, there is no re-scaling and updating of the state activation levels relies on the entire information conveyed by m_t . On the other extreme, when $\delta = 1$, several trials have to be accomplished before converging to a given configuration of the state activation levels. Equivalently, higher values of δ affect the entropy of the distribution over hidden states by decreasing the probability of sampling of the correct feature. We therefore refer to δ as *information loss*.

The free parameters λ and δ are central to our computational model, since they regulate the rate at which the internal model converges to the true task environmental model. can be expressed in compact notation as follows:

$$\omega_t = f_t \omega_{t-1}^\delta m_t + \lambda [(1 - f_t) \omega_{t-1}^\delta (1 - m_t)] \omega_{t-1}. \quad (6)$$

Note that the information loss parameter δ affects the amount of information that a cognitive agent acquires from environmental contingencies, irrespective of the type of feedback received. Global information loss thus affects the rate at which the divergence between the agent's internal model and the true model is minimized. [Figure 2](#) illustrates these ideas.

The probabilistic representation of adaptive behaviour provided by our Bayesian agent model allows us to quantify latent cognitive dynamics by means of meaningful information-theoretic measures. Information theory has proven to be an effective and natural mathematical language to account for functional integration of structured cognitive processes and to relate them to brain activity ([Koechlin & Summerfield, 2007](#); [Friston et al., 2017](#); [Collell & Fauquet, 2015](#); [Strange et al., 2005](#); [Friston, 2003](#)). In particular, we are interested in three key measures, namely, *Bayesian surprise*, \mathcal{B}_t , *Shannon surprise*, \mathcal{I}_t , and *entropy*, \mathcal{H}_t . The subscript t indicates that we can compute each quantity on a trial-by-trial basis. Each quantity is amenable to a specific interpretation in terms of separate neurocognitive processes. Bayesian surprise \mathcal{B}_t quantifies the magnitude of the update from prior belief to posterior belief. Shannon surprise \mathcal{I}_t quantifies the improbability of an observation given an agent's prior expectation. Finally, entropy \mathcal{H}_t measures the degree of epistemic uncertainty regarding the true environmental states. Such measures are thought to account for the ability of the agent to manage uncertainty as emerging as a function of competing behavioral affordances ([Hirsh, Mar & Peterson, 2012](#)). We expect an adaptive system to attenuate uncertainty over environmental states (current features) by reducing the entropy of its internal probabilistic model.

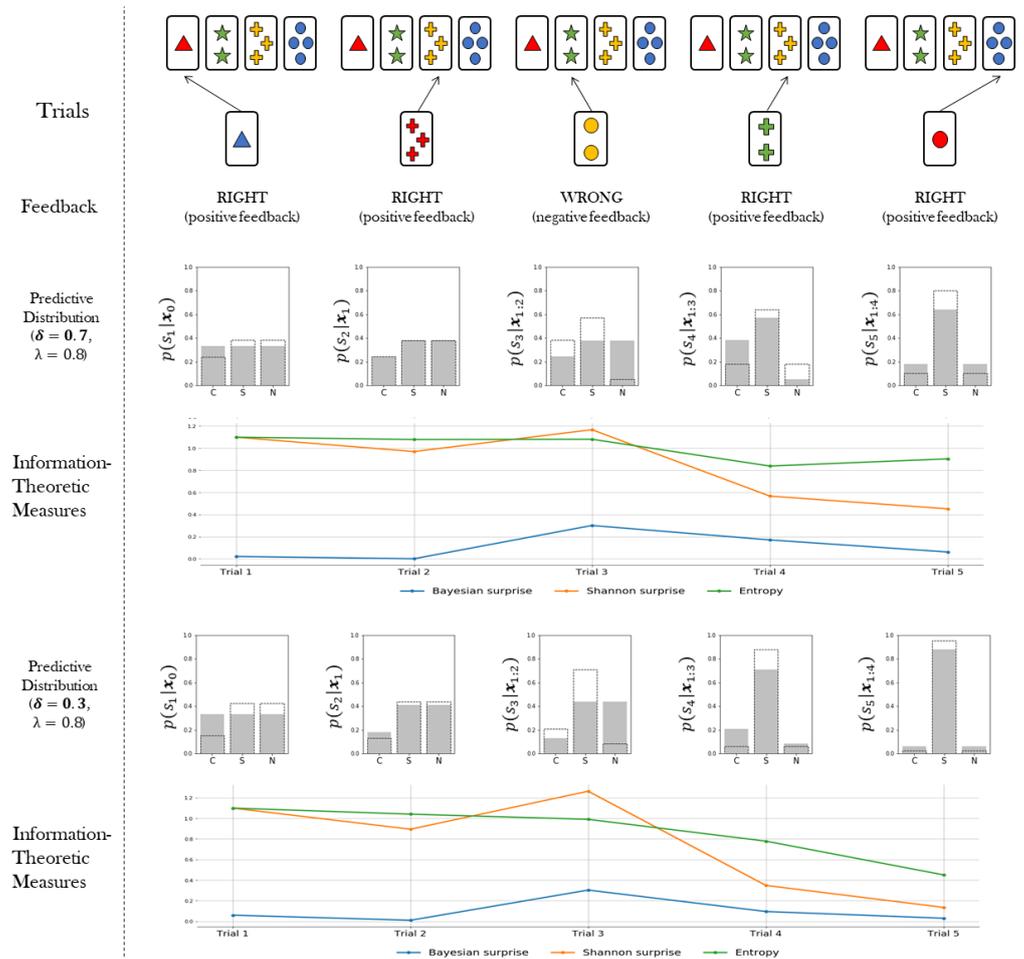


Figure 2 The figure shows the rate of convergence of the predictive distributions to the true task environmental model. The predictive distributions at trial $t + 1$ depends on the sorting action a_t (first row) and the received feedback f_t (second row). Two examples of updating a predictive distribution are shown: one in which information loss is high ($\delta = 0.7$, third row), and one in which information loss is low ($\delta = 0.3$, fifth row). High information loss slows down the convergence of the internal model to the true environmental model. The gray bar plots represent the predictive probability distribution over the rules from which an action is sampled at each trial. Dotted bars represent the updated predictive distribution after the feedback observation. For each scenario, trial-by-trial information-theoretic measures are shown.

Full-size [DOI: 10.7717/peerj.10316/fig-2](https://doi.org/10.7717/peerj.10316/fig-2)

Bayesian surprise can be computed as the Kullback–Leibler (\mathbb{KL}) divergence between prior and posterior beliefs about the environmental states. Thus, Bayesian surprise accounts for the divergence between the predictive model for the current trial and the updated predictive model for the upcoming trial. It is computed as follows:

$$\begin{aligned} \mathcal{B}_t &= \mathbb{KL}[p(s_{t+1}|x_{0:t})||p(s_t|x_{0:t-1})] \\ &= \sum_{i=1}^3 \left[p(s_{t+1} = i|x_{0:t}) \log \left(\frac{p(s_{t+1} = i|x_{0:t})}{p(s_t = i|x_{0:t-1})} \right) \right] \end{aligned} \quad (7)$$

The Shannon surprise of a current observation given a previous one is computed as the conditional information content of the observation:

$$\begin{aligned} \mathcal{I}_t &= -\log p(x_t | x_{0:t-1}) \\ &= -\log \sum_{i=1}^3 [p(x_t | s_t = i) p(s_t = i | x_{0:t-1})] \end{aligned} \quad (8)$$

Finally, the entropy is computed over the predictive distribution in order to account for the uncertainty in the internal model of the agent in trial t as follows:

$$\begin{aligned} \mathcal{H}_t &= \mathbb{E}[-\log p(s_t | x_{0:t-1})] \\ &= -\sum_{i=1}^3 p(s_t = i | x_{0:t-1}) \log p(s_t = i | x_{0:t-1}) \end{aligned} \quad (9)$$

Once the flexibility (λ) and information loss (δ) parameters are estimated from data, the information-theoretic quantities can be easily computed and visualized for each trial of the WCST (see Fig. 2). This allows to rephrase standard neurocognitive constructs in terms of measurable information-theoretic quantities. Moreover, the dynamics of these quantities, as well as their interactions, can be used for formulating and testing hypotheses about the neurocognitive underpinnings of adaptive behavior in a principled way, as discussed later in the paper. A summary of all quantities relevant for our computational model is provided in Table 1.

Simulations

In this section we evaluate the expressiveness of the model by assessing its ability to reproduce meaningful behavioral patterns as a function of its two free parameters. We study how the generative model behaves when performing the WCST in a 2-factorial simulated Monte Carlo design where flexibility (λ) and information loss (δ) are systematically varied.

In this simulation, the Heaton version of the task (Heaton, 1981) is administered to the Bayesian cognitive agent. In this particular version, the sorting rule (true environmental state) changes after a fixed number of consecutive correct responses. In particular, when the agent correctly matches the target card in 10 consecutive trials, the sorting rule is automatically changed. The task ends after completing a maximum of 128 trials.

Generative model

The cognitive agent's responses are generated at each time step (trial) by processing the experimental feedback. Its performance depends on the parameters governing the computation of the relevant quantities. The generative algorithm is outlined in Algorithm 1.

Algorithm 1 Bayesian cognitive agent

-
- 1: Set parameters $\theta = (\lambda, \delta)$.
 - 2: Set initial activation levels $\omega_0 = (0.5, 0.5, 0.5)$.
 - 3: Set initial observation $\mathbf{x}_0 = \emptyset$ and $p(s_1|\mathbf{x}_0) = p(s_1)$.
 - 4: **for** $t = 1, \dots, T$ **do**
 - 5: Sample feature from prior/predictive internal model $s_t \sim p(s_t|\mathbf{x}_{0:t-1})$.
 - 6: Obtain a new observation $\mathbf{x}_t = (a_t, f_t)$.
 - 7: Compute state posterior $p(s_t|\mathbf{x}_{0:t})$.
 - 8: Compute new activation levels ω_t .
 - 9: Compute stability matrix $\Gamma(t)$.
 - 10: Update prior/predictive internal model to $p(s_{t+1}|\mathbf{x}_{0:t})$.
 - 11: **end for**
-

Simulation 1: clinical assessment of the Bayesian agent

Ideally, the qualitative performance of the Bayesian cognitive agent will resemble human performance. To this aim, we adopt a metric which is usually employed in clinical assessment of test results in neurological and psychiatric patients (Braff *et al.*, 1991; Zakzanis, 1998; Bechara & Damasio, 2002; Landry & Al-Taie, 2016). Thus, agent performance is codified according to a neuropsychological criterion (Heaton, 1981; Flashman, Homer & Freides, 1991) which allows to classify responses into several response types. These response types provide the scoring measures for the test.

Here, we are interested in: (1) non-perseverative errors (E); (2) perseverative errors (PE); (3) number of trials to complete the first category (TFC); and (4) number of failures to maintain set (FMS). Perseverative errors occur when the agent applies a sorting rule which was valid before the rule has been changed. Usually, detecting a perseveration error is far from trivial, since several response configurations could be observed when individuals are required to shift a sorting rule after completing a category (see Flashman, Homer & Freides (1991) for details). On the other hand, non-perseverative errors refer to all errors which do not fit the above description, or in other words, do not occur as a function of changing the sorting rule, such as casual errors.

The number of trials to complete the first category tells us how many trials the agent needs in order to achieve the first sorting principle, and can be seen as an index of conceptual ability (Anderson, 2008; Singh, Aich & Bhattarai, 2017). Finally, a failure to maintain a set occurs when the agent fails to match cards according to the sorting rule after it can be determined that the agent has acquired the rule. A given sorting rule is assumed to be acquired when the individual correctly sorts at least five cards in a row (Heaton, 1981; Figueroa & Youmans, 2013). Thus, a failure to maintain a set arises whenever a participant suddenly changes the sorting strategy in the absence of negative feedback. Failures to maintain a set are mostly attributed to distractibility. We compute this measure by counting the occurrences of first errors after the acquisition of a rule.

We run the generative model by varying flexibility across four levels, $\lambda \in \{0.3, 0.5, 0.7, 0.9\}$, and information loss across three levels, $\delta \in \{0.4, 0.7, 0.9\}$. We generate data from 150 synthetic cognitive agents per parameter combination and compute standard

Table 1 Descriptive summary of all quantities involved in our model representation.

Expression	Name	Description
$s_t \in \{1, 2, 3\}$	Sorting rule	Card feature relevant for the sorting criterion in trial t .
$a_t \in \{1, 2, 3, 4\}$	Choice action	Action of choosing one of the four stimulus cards in trial t .
$f_t \in \{0, 1\}$	Feedback	Indicates whether the action of matching a stimulus to a target card is correct or not in trial t .
$x_t = (a_t, f_t)$	Observation	Pair of action and feedback which constitutes the agent's observation in trial t .
$\Gamma(t)$	Stability matrix	Matrix encoding the agent's beliefs about state transitions from trial t to the next trial $t + 1$.
$\lambda \in [0, 1]$	Flexibility	Parameter encoding the efficiency to disengage attention from a currently attended hidden state when signaled by the environment.
$\delta \in [0, 1]$	Information loss	Parameter encoding how efficiently the agent's internal model converges to the true environmental model based on experience.
$m_t^{(i)} \in \{0, 1\}$	Matching signal	Signal indicating whether feature i is relevant in trial t based on the feedback received.
$\omega_t^{(i)} \in [0, 1]$	State activation level	Agent's internal measure of the accrued evidence for the hidden environmental state i in trial t .
$\mathcal{B}_t \in \mathbb{R}^+$	Bayesian surprise	Kullback–Leibler divergence between prior and posterior beliefs about hidden environmental states in trial t .
$\mathcal{I}_t \in \mathbb{R}^+$	Shannon surprise	Information-theoretic surprise encoding the improbability or unexpectedness of an observation in trial t .
$\mathcal{H}_t \in \mathbb{R}^+$	Entropy	Degree of epistemic uncertainty in the internal model of the environment in trial t .

scoring measures for each of the agents simulated responses. Results from the simulation runs are depicted in [Table 2](#) and a graphical representation is provided in [Fig. 3](#).

The simulated performance of our Bayesian cognitive agents demonstrates that different parameter combinations capture different meaningful behavioral patterns. In other words, flexibility and information loss seem to interact in a theoretically meaningful way.

First, overall errors increase when flexibility (λ) decreases, which is reflected by the inverse relation between the number of casual, as well as perseverative, errors and the values of parameter λ . Moreover, this pattern is consistent across all the levels of parameter δ . More precisely, information loss (δ) seems to contribute to the characterization of the casual and the perseverative components of the error in a different way. Perseverative errors are likely to occur after a sorting rule has changed and reflect the inability of the agent to use feedback to disengage attention from the currently attended feature. They therefore result from local cognitive dynamics conditioned on a particular stage of the task (e.g., after completing a series of correct responses).

Second, information loss does not interact with flexibility when perseverative errors are considered. This is due to the fact that high information loss affects general performance by yielding a dysfunctional response strategy which increases the probability of making an error at any stage of the task. The lack of such interaction provides evidence that our

Table 2 Mean clinical scoring measures as functions of flexibility (λ) and information loss (δ). Cells show the average scores across simulated agents (standard deviation is shown in parenthesis).

Scoring measure	Info. Loss (δ)	Flexibility (λ)			
		$\lambda = 0.3$	$\lambda = 0.5$	$\lambda = 0.7$	$\lambda = 0.9$
Casual Errors (E)	$\delta = 0.4$	9.07 (2.68)	7.95 (2.07)	7.50 (2.13)	6.85 (1.75)
	$\delta = 0.7$	10.84 (2.35)	9.60 (2.2)	8.25 (2.23)	7.37 (1.74)
	$\delta = 0.9$	12.75 (2.96)	11.25 (2.43)	9.12 (2.09)	7.79 (1.73)
Perseverative Errors (PE)	$\delta = 0.4$	20.81 (2.27)	18.18 (1.88)	14.99 (1.88)	12.37 (1.12)
	$\delta = 0.7$	19.77 (2.55)	17.65 (2.26)	15.42 (1.94)	12.39 (1.47)
	$\delta = 0.9$	18.56 (2.76)	16.58 (2.53)	14.49 (2.03)	12.33 (1.44)
Trials to First Category (TFC)	$\delta = 0.4$	12.20 (1.46)	11.91 (1.35)	11.83 (1.24)	11.67 (1.04)
	$\delta = 0.7$	13.82 (2.76)	13.32 (2.52)	12.97 (2.13)	12.29 (1.53)
	$\delta = 0.9$	17.27 (4.21)	16.63 (4.04)	14.39 (3.58)	12.91 (1.91)
Failures to Maintain Set (FMS)	$\delta = 0.4$	0.11 (0.31)	0.09 (0.31)	0.05 (0.32)	0.02 (0.14)
	$\delta = 0.7$	1.65 (1.4)	1.41 (1.3)	0.84 (0.91)	0.35 (0.69)
	$\delta = 0.9$	4.44 (1.96)	3.88 (1.86)	2.79 (1.56)	1.54 (1.25)

computational model can disentangle between error patterns due to perseveration and those due to general distractibility, according to neuropsychological scoring criteria.

However, in our framework, flexibility (λ) is allowed to yield more general and non-local cognitive dynamics as well. Indeed, λ plays a role whenever belief updating is demanded as a function of negative feedback. An error classified as non-perseverative (e.g., casual error) by the scoring criteria might still be processed as a feedback-related evidence for belief updating. Consistently, the interaction between λ and δ in accounting for causal errors shows that performance worsens when both flexibility and information loss become less optimal, and that such pattern becomes more pronounced for lower values of δ .

On the other hand, a specific effect of information loss (δ) can be observed for the scoring measures related to slow information processing and distractibility. The number of trials to achieve the first category reflects the efficiency of the agent in arriving at the first true environmental model. Flexibility does not contribute meaningfully to the accumulation of errors before completing the first category for some levels of information loss. This is reflected by the fact that the mean number of trials increases as a function of δ , and do not change across levels of λ for low and mid values of δ . A similar pattern applies for failures to maintain a set. Both scoring measures index a deceleration of the process of evidence accumulation for a specific environmental configuration, although the latter is a more exhaustive measures of dysfunctional adaptation.

Therefore, an interaction between parameters can be observed when information loss is high. A slow internal model convergence process increases the amount of errors due to improper rule sampling from the internal environmental model. However, internal model convergence also plays a role when a new category has to be accomplished after completing an older one. On the one hand, compromised flexibility increases the amount of errors due to inefficient feedback processing. This leads to longer trial windows needed to achieve the first category. On the other hand, when information loss is high, belief updating upon

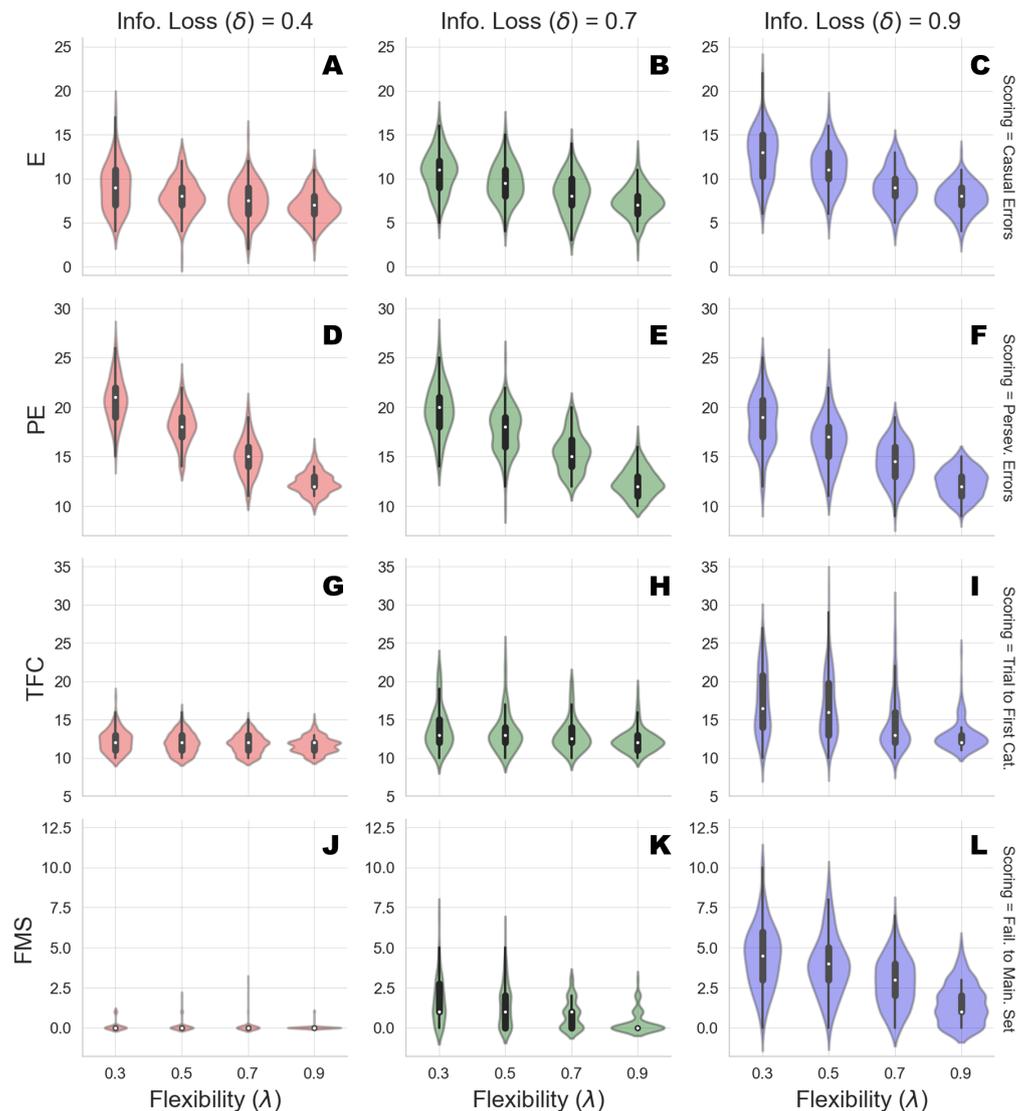


Figure 3 Clinical scoring measures as functions of flexibility (λ) and information loss (δ) - simulated scenarios. The different cells show the violin plots for the estimated distribution densities of the scoring measures obtained from the group of synthetic individuals, for the levels of λ across different levels of δ . In particular, they show the distribution of non-perseverative errors (E: A–C), perseverative errors (PE: D–F), number of trials to complete the first category (TFC: G–I), number of failures to maintain set (FMS: J–L) obtained from 150 synthetic agent's response simulations for each cell of the factorial design.

Full-size [DOI: 10.7717/peerj.10316/fig-3](https://doi.org/10.7717/peerj.10316/fig-3)

negative feedback is compromised due to high internal model uncertainty. At this point, the probability to err due to distractibility increases, as accounted by the failures to maintain a set measures.

Finally, the joint effect of δ and λ for high levels of information loss suggests that the roles played by the two cognitive parameters in accounting for adaptive functioning can be entangled when neuropsychological scoring criteria are considered.

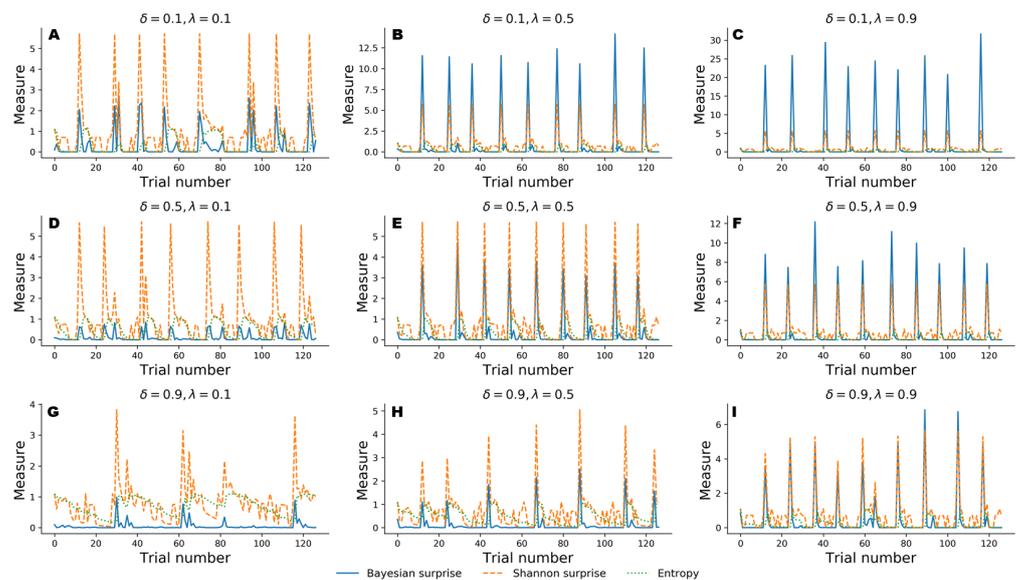


Figure 4 Information-theoretic measures varying as a function of flexibility λ and information loss δ across 128 trials of the WCST. Trajectories depicted in A, D, and G show cognitive dynamics across the levels of information loss when flexibility is low. B, E, and H show the unfolding of information-theoretic quantities when flexibility is mildly impaired, whilst C, F, and I refer to an optimal flexibility value. Optimal belief updating and uncertainty reduction are achieved with low information loss and high flexibility (C).

Full-size DOI: [10.7717/peerj.10316/fig-4](https://doi.org/10.7717/peerj.10316/fig-4)

Simulation 2: Information-theoretic analysis of the Bayesian agent

In the following, we explore a different simulation scenario in which information-theoretic measures are derived to assess performance of the Bayesian cognitive agent. In particular, we explore the functional relationship between cognitive parameters and the dynamics of the recovered information-theoretic measures by simulating observed responses by varying flexibility across three levels, $\lambda \in \{0.1, 0.5, 0.9\}$, and information loss across three levels, $\delta \in \{0.1, 0.5, 0.9\}$.

For this simulation scenario, we make no prior assumptions about sub-types of error classification. Instead, we investigate the dynamic interplay between Bayesian surprise, \mathcal{B}_t , Shannon surprise, \mathcal{I}_t , and entropy, \mathcal{H}_t over the entire course of 128 trials in the WCST.

Figure 4 depicts results from the nine simulation scenarios. Although an exhaustive discussion on cognitive dynamics should couple information-theoretic measures with patterns of correct and error responses, we focus solely on the information-theoretic time series for illustrative purposes. We refer to the ‘Application’ for a more detailed description of the relation between observed responses and estimated information-theoretic measures in the context of data from a real experiment.

Again, simulated performance of the Bayesian cognitive agent shows that different parameter combinations yield different patterns of cognitive dynamics. Observed spikes and their related magnitudes signal informative task events (e.g., unexpected negative feedback), as accounted by Shannon surprise, or belief updating, as accounted by Bayesian

surprise. Finally, entropy encodes the epistemic uncertainty about the environmental model on a trial-by-trial basis.

In general, low information loss (δ) ensures optimal behavior by speeding up internal model convergence by decreasing the number of trials needed to minimize uncertainty about the environmental states. Low uncertainty reflects two main aspects of adaptive behavior. On the one hand, the probability that a response occurs due to sampling of improper rules decreases, allowing the agent to prevent random responses due to distractibility. On the other hand, model convergence entails a peaked Shannon surprise when a negative feedback occurs, due to the divergence between predicted and actual observations.

Flexibility (λ) plays a crucial role in integrating feedback information in order to enable belief updating. The first row depicted in Fig. 4 shows cognitive dynamics related to low information loss, across the levels of flexibility. As can be noticed, there is a positive relation between the magnitude of the Bayesian surprise and the level of flexibility, although unexpectedness yields approximately the same amount of signaling, as accounted by peaked Shannon surprise. From this perspective, surprise and belief updating can be considered functionally separable, where the first depends on the particular internal model probability configuration related to δ , whilst the second depends on flexibility λ .

However, more interesting patterns can be observed when information loss increases. In particular, model convergence slows down and several trials are needed to minimize predictive model entropy. Casual errors might occur within trial windows characterized by high uncertainty, and interactions between entropy and Shannon surprise can be observed in such cases. In particular, Shannon surprise magnitude increases when model's entropy decreases, that is, during task phases in which the internal model has already converged. As a consequence, negative feedback could be classified as informative or uninformative, based on the uncertainty in the current internal model. This is reflected by the negative relation between entropy and Shannon surprise, as can be noticed by inspecting the graphs depicted in the third row of Fig. 4. Therefore, the magnitude of belief updating depends on the interplay between entropy and Shannon surprise, and can differ based on the values of the two measures in a particular task phase.

In sum, both simulation scenarios suggest that the simulated behavior of our generative model is in accord with theoretical expectations. Moreover, the flexibility and information loss parameters can account for a wide range of observed response patterns and inferred dynamics of information processing.

Parameter estimation

In this section, we discuss the computational framework for estimating the parameters of our model from observed behavioral data. Parameter estimation is essential to inferring the cognitive dynamics underlying observed behavior in real-world applications of the model. This section is slightly more technical and can be skipped without significantly affecting the flow of the text.

Computational framework

Rendering our cognitive model suitable for application in real-world contexts also entails accounting for uncertainty about parameter estimates. Indeed, uncertainty quantification turns out to be a fundamental and challenging goal when first-level quantities, that is, cognitive parameter estimates, are used to recover (second-level) information-theoretic measures of cognitive dynamics. The main difficulties arise when model complexity makes estimation and uncertainty quantification intractable at both analytical and numerical levels. For instance, in our case, probability distributions for the hidden model are generated at each trial, and the mapping between hidden states and responses changes depending on the structure of the task environment.

Identifying such a dynamic mapping is relatively easy from a generative perspective, but it becomes challenging, and almost impossible, when inverse modeling is required. Generally, this problem arises when the likelihood function relating model parameters to the data is not available in closed-form or too complex to be practically evaluated (Sisson & Fan, 2011). To overcome these limitations, we apply the first version of the recently developed *BayesFlow* method (see Radev et al., 2020 for mathematical details). At a high-level, *BayesFlow* is a simulation-based method that estimates parameters and quantifies estimation uncertainty in a unified Bayesian probabilistic framework when inverting the generative model is intractable. The method is based on recent advances in deep generative modeling and makes no assumptions about the shape of the true parameter posteriors. Thus, our ultimate goal becomes to approximate and analyze the joint posterior distribution over the model parameters. The parameter posterior is given via an application of Bayes' rule:

$$p(\theta|x_{0:T}, m_{0:T}) = \frac{p(x_{0:T}, m_{0:T}|\theta)p(\theta)}{\int p(x_{0:T}, m_{0:T}|\theta)p(\theta)d\theta} \quad (10)$$

where we set $\theta = (\lambda, \delta)$ and stack all observations and matching signals into the vectors $x_{0:T} = (x_0, x_1, \dots, x_T)$ and $m_{0:T} = (m_0, m_1, \dots, m_T)$, respectively. The *BayesFlow* method uses simulations from the generative model to optimize a neural density estimator which learns a probabilistic mapping between raw data and parameters. It relies on the fact that data can easily be simulated by repeatedly running the generative model with different parameter configurations θ sampled from the prior. During training, the neural network estimator iteratively minimizes the divergence between the true posterior and an approximate posterior. Once the network has been trained, we can efficiently obtain samples from the approximate joint posterior distribution of the cognitive parameters of interest, which can be further processed in order to extract meaningful summary statistics (e.g., posterior means, medians, modes, etc.). Importantly, we can apply the same pre-trained inference network to an arbitrary number of real or simulated data sets (i.e., the training effort *amortizes* over multiple evaluations of the network).

For our purposes of validation and application, we train the network for 50 epochs which amount to 50000 forward simulations. As a prior, we use a bivariate continuous uniform distribution $p(\theta) \sim \mathcal{U}([0, 0], [1, 1])$. We then validate performance on a separate validation set of 1000 simulated data sets with known *ground-truth* parameter values. Training the

networks took less than a day on a single machine with an NVIDIA[®] GTX1060 graphics card (CUDA version 10.0) using TensorFlow (version 1.13.1) (Abadi et al., 2016). In contrast, obtaining full parameter posteriors from the entire validation set took approximately 1.78 s. In what follows, we describe and report all performance validation metrics.

Performance metrics and validation results

To assess the accuracy of point estimates, we compute the root mean squared error (RMSE) and the coefficient of determination (R^2) between posterior means and true parameter values. To assess the quality of the approximate posteriors, we compute a calibration error (Radev et al., 2020) of the empirical coverage of each marginal posterior. Finally, we implement simulation-based calibration (SBC, Talts et al., 2018) for visually detecting systematic biases in the approximate posteriors.

Point Estimates. Point estimates obtained by posterior means as well as corresponding RMSE and R^2 metrics are depicted in Figs. 5A–5B. Note, that point estimates do not have any special status in Bayesian inference, as they could be misleading depending on the shape of the posteriors. However, they are simple to interpret and useful for ease-of-comparison. We observe that pointwise recovery of λ is better than that of δ . This is mainly due to suboptimal pointwise recovery in the lower (0, 0.1) range of δ . This pattern is evident in Figs. 5A–5B and is due to the fact that δ values in this range produce almost indistinguishable data patterns. Bootstrap estimates yielded an average RMSE of 0.155 ($SD = 0.004$) and an average R^2 of 0.708 ($SD = 0.015$) for the δ parameter. An average RMSE of 0.094 ($SD = 0.002$) and an average R^2 of 0.895 ($SD = 0.007$) were obtained for the λ parameter. These results suggest good global pointwise recovery but also warrant the inspection of full posteriors, especially in the low ranges of δ .

Full Posteriors. Average bootstrap calibration error was 0.011 ($SD = 0.005$) for the marginal posterior of δ and 0.014 ($SD = 0.007$) for the marginal posterior of λ . Calibration error is perhaps the most important metric here, as it measures potential under- or overconfidence across all confidence intervals of the approximate posterior (i.e., an α -confidence interval should contain the true posterior with a probability of α , for all $\alpha \in (0, 1)$). Thus, low calibration error indicates a faithful uncertainty representation of the approximate posteriors. Additionally, SBC-histograms are depicted in Figs. 5C–5D. As shown by Talts et al. (2018), deviations from the uniformity of the rank statistic (also known as a PIT histogram) indicate systematic biases in the posterior estimates. A visual inspection of the histograms reveals that the posterior means slightly overestimate the true values of δ . This corroborates the pattern seen in Figs. 5A–5B for the lower range of δ .

Finally, Figs. 5E–5H depicts the full marginal posteriors on two example validation sets. Even on these two data sets, we observe strikingly different posterior shapes. The marginal posterior of δ obtained from the first data set is slightly left-skewed and has its density concentrated over the (0.8, 1.0) range. On the other hand, the marginal posterior of δ from the second data set is noticeably right-skewed and peaked across the lower range of the parameter. The marginal posteriors of λ appear more symmetric and warrant the use of the posterior mean as a useful summary of the distribution. These two examples underline the importance of investigating full posterior distributions as a means to encode

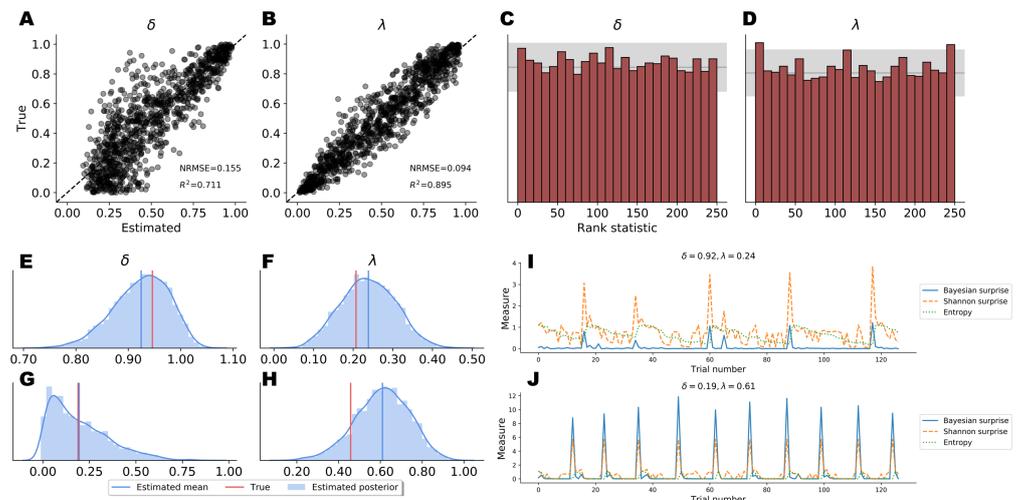


Figure 5 Parameter recovery results on validation data; (A and B) posterior means vs. true parameter values; (C and D) histograms of the rank statistic used for simulation-based calibration; (E–H) example full posteriors for two validation data sets; (I and J) example information-theoretic dynamics recovered from the parameter posteriors.

Full-size [DOI: 10.7717/peerj.10316/fig-5](https://doi.org/10.7717/peerj.10316/fig-5)

epistemic uncertainty about parameter values. Moreover, they demonstrate the advantage of imposing no distributional assumptions on the resulting posteriors, as their form and sharpness can vary widely depending on the concrete data set.

APPLICATION

In this section we fit the Bayesian cognitive model to real clinical data. The aim of this application is to evaluate the ability of our computational framework to account for dysfunctional cognitive dynamics of information processing in substance dependent individuals (SDI) as compared to healthy controls.

Rationale

The advantage of modeling cognitive dynamics in individuals from a clinical population is that model predictions can be examined in light of available evidence about individual performance. For instance, SDIs are known to demonstrate inefficient conceptualization of the task and dysfunctional, error-prone response strategies. This has been attributed to defective error monitoring and behavior modulation systems, which depend on cingulate and frontal brain regions functionality (Kübler, Murphy & Garavan, 2005; Willuhn, Sun & Steiner, 2003). On the other hand, the WCST should be a rather easy and straightforward task for healthy participants to obtain excellent performance. Therefore, we expect our model to consistently capture such characteristics. To test these expectations, we estimate the two relevant parameters λ and δ from both clinical patients and healthy controls from an already published dataset (Bechara & Damasio, 2002).

The data

The dataset used in this application consists of responses collected by administering the standard Heaton version of the WCST ([Heaton, 1981](#)) to healthy participants and SDIs. In this version of the task, the sorting rule changes when a participant collects a series of 10 consecutive correct responses, and the task ends when this happens for 6 times. Participants in the study consisted of 39 SDIs and 49 healthy individuals. All participants were adults (>18 years old) and gave their informed consent for inclusion which was approved by the appropriate human subject committee at the University of Iowa. SDIs were diagnosed as substance dependent based on the Structured Clinical Interview for DSM-IV criteria ([First, 1997](#)).

Model fitting

We fit the Bayesian cognitive agent separately to data from each participant in order to obtain individual-level posterior distributions. We apply the same *BayesFlow* network trained for the previous simulation studies, so obtaining posterior samples for each participant is almost instant (due to amortized inference).

Results

The means of the joint posterior distributions are depicted for each individual in [Fig. 6](#), and provide a complete overview of the heterogeneity in cognitive sub-components at both individual and group levels (individual-level full joint posterior distributions can be found in the [Appendix A1](#)).

The estimates reveal a rather interesting pattern across both healthy and SDI participants. In particular, in both clinical and control groups, individuals with a poor flexibility (e.g., low values of λ) can be detected. However, the group parameter space appears to be partitioned into two main clusters consisting of individuals with high and low flexibility, respectively. As can be noticed, the majority of SDIs belongs to the latter cluster, which suggests that the model is able to capture error-related defective behavior in the clinical population and attribute it specifically to the flexibility parameter. On the other hand, individual performance seems hardly separable along the information loss parameter dimension.

As a further validation, we compare the classification performance of two logistic regression models. The first uses the estimated parameter means as inputs and the participants' binary group assignment (patient vs. control) as an outcome. The second uses the four standard clinical measures (non-perseverative errors (E), perseverative errors (PE), number of trials to complete the first category (TFC), number of failures to maintain set (FMS) computed from the sample as inputs and the same outcome. Since we are interested solely in classification performance and want to mitigate potential overfitting due to small sample size, we compute leave-one-out cross-validated (LOO-CV) performance for both models. Interestingly, both logistic regression models achieve the same accuracy of 0.70, with a sensitivity of 0.71 and specificity of 0.70. Thus, it appears that our model is able to differentiate between SDIs and healthy individuals as good as the standard clinical measures.

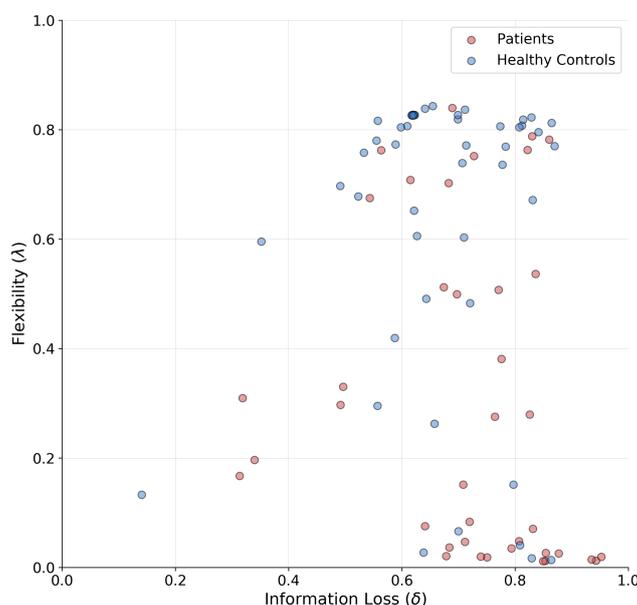


Figure 6 Joint posterior mean coordinates of the cognitive parameters, flexibility (λ) and information loss (δ), estimated for each individual. We observe a great heterogeneity in the distribution of posterior means, most pronouncedly for the flexibility parameter. However, a moderate between-subject variability in information loss can still be observed in both groups.

Full-size  DOI: [10.7717/peerj.10316/fig-6](https://doi.org/10.7717/peerj.10316/fig-6)

However, as pointed out in the previous sections, estimated parameters serve merely as a basis to reconstruct cognitive dynamics by means of the trial-by-trial unfolding of information-theoretic measures. Moreover, cognitive dynamics can only be analysed and interpreted by relying on the joint contribution of both estimated parameters and individual-specific observed response patterns.

To further clarify this concept, we investigate the reconstructed time series of information-theoretic quantities based on the response patterns of two exemplary individuals (Fig. 7). In particular, Fig. 7A depicts the behavioral outcomes of a SDI with sub-optimal performance where the information-theoretic trajectories are reconstructed by taking the corresponding posterior means ($[\bar{\lambda} = 0.07, \bar{\delta} = 0.82]$), thus representing compromised flexibility and high information loss. Differently, Fig. 7B shows the information-theoretic path related to response dynamics of an optimal control participant, according to the parameter set $[\bar{\lambda} = 0.60, \bar{\delta} = 0.35]$, representing relatively high flexibility, and low information loss. Note, that in both cases, the reconstructed information-theoretic measures are based on the estimated posterior means for ease of comparison (see Appendix A1 for the full joint posterior densities of the two exemplary individuals and the rest of the sample).

Results in Fig. 7A account for a typical sub-optimal behavior observed in the SDI group, where several errors are produced in different phases of the task. The error patterns produced by such an individual might be induced by a non-trivial interaction between cognitive sub-components. Lower values of flexibility imply that errors are likely to be

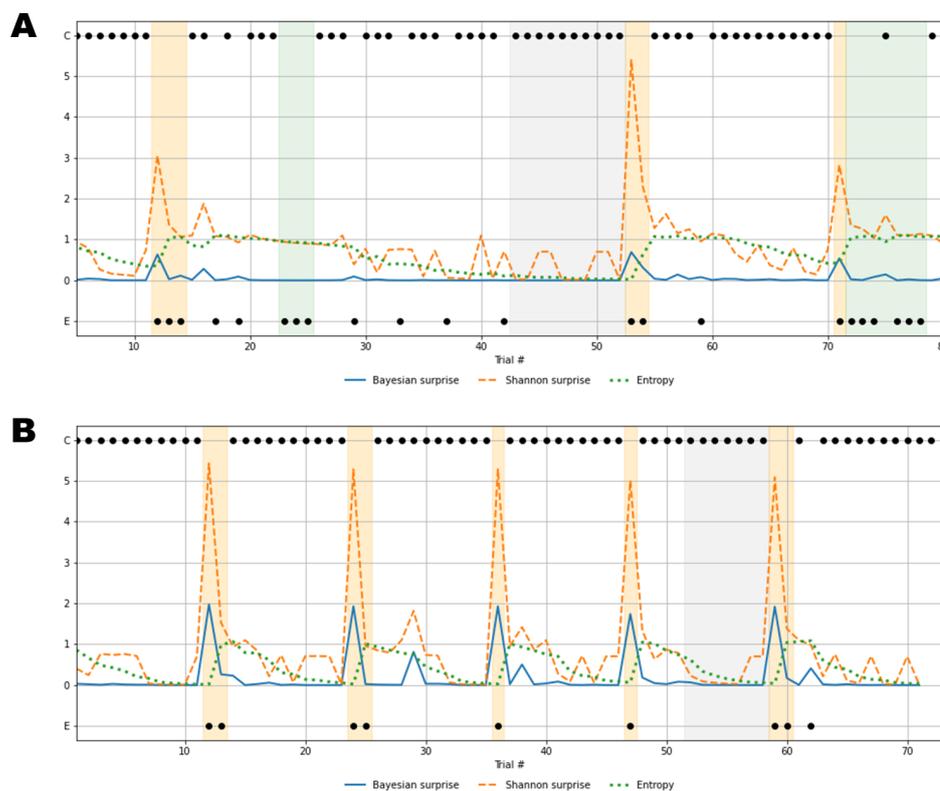


Figure 7 Recovered cognitive dynamics of two exemplary individuals. (A) Trial-by-trial information-theoretic measures of a SDI characterized by very low flexibility and very high information loss; (B) trial-by-trial information-theoretic measures of a healthy individual characterized by relatively high flexibility and low information loss. Labels C and E in the y-axis indicate correct and error responses.

Full-size DOI: [10.7717/peerj.10316/fig-7](https://doi.org/10.7717/peerj.10316/fig-7)

produced by generating responses from an internal environmental model which is no longer valid. In other words, the agent is unable to rely on local feedback-related information in order to update beliefs about hidden states. On the other hand, higher values of information loss reflect a general inefficiency of belief updating processes due to slow convergence to the optimal probabilistic environmental model. From this perspective, Bayesian surprise \mathcal{B}_t and Shannon surprise \mathcal{I}_t might play different roles in regulating behavior based on different internal model probability configurations. In addition, errors might be processed differently based on the status of the internal environmental representation, as reflected by the entropy of the predictive model, \mathcal{H}_t . Thus, information-theoretic measures allow to describe cognitive dynamics on a trial-by-trial basis and, further, to disentangle the effect that different feedback-related information processing dynamics exert on adaptive behavior.

Processing unexpected observations is accounted by the quantification of surprise upon observing a response-feedback pair which is inconsistent with the current internal model of the task environment. Negative feedback is maximally informative when errors occur after the internal model has converged to the true task model (grey area, Fig. 7A), or

the entropy approaches zero (grey line, Fig. 7A). The Shannon surprise (orange line) is maximal when errors occur within trial windows in which the agent's uncertainty about environmental states is minimal (orange areas, Fig. 7A). However, internal model updates following an informative feedback are not optimally performed, which is reflected by very small Bayesian surprise (blue line, Fig. 7A). This can be attributed to impaired flexibility and reflects the fact that after internal model convergence, informative feedback is not processed adequately and the internal model becomes impervious to change.

Conversely, errors occurring when the agent is uncertain about the true environmental state carry no useful information for belief updating, since the system fails to conceive such errors as unexpected and informative. The information loss parameter plays a crucial role in characterizing this cognitive behavior. The slow convergence to the true environmental model, accompanied by the slow reduction of entropy in the predictive model, leads to a large number of trials required to achieve a good representation of the current task environment (white areas, Fig. 7A). Errors occurring within trial windows with large predictive model entropy (green area, Fig. 7A) do not affect subsequent behavior, and feedback is maximally uninformative.

Rather different cognitive dynamics can be observed in Fig. 7B, accounting for a typical optimal behavior where the errors produced fall within the trial windows which follow a rule completion (e.g., when the individual completes a sequence of 10 consecutive correct responses), and, thus, the environmental model becomes obsolete. However, the high flexibility, λ , allows to rely on local feedback-related information to suddenly update beliefs about the hidden states, that is, the most appropriate sorting rule. In this case, negative feedback become maximally informative after model convergence (grey area, Fig. 7B) and the process of entropy reduction (green line, Fig. 7B) is faster (e.g., less trials are needed) compared to the sub-optimal behavior scenario. Since uncertainty about the environmental states decreases faster, the Shannon surprise is always highly peaked when errors occur (orange line, Fig. 7B), thus ensuring an efficient employment of the local feedback-related information. Accordingly, higher values of Bayesian surprise are observed (blue line, Fig. 7B), revealing optimal internal model updating.

In general, the role that predictive (internal) model uncertainty plays in characterizing the way the agent processes feedback allows to disentangle sub-types of errors based on the information they convey for subsequent belief updating. From this perspective, error classification is entirely dependent on the status of the internal environmental model across task phases. Identifying such a dynamic latent process is therefore fundamental, since the error codification criterion evolves with respect to the internal information processing dynamics. Otherwise, the problem of inferring which errors are due to perseverance in maintaining an older (converged) internal model and which due to uncertainty about the true environmental state becomes intractable, or even impossible.

DISCUSSION

Investigating information processing related to changing environmental contingencies is fundamental to understanding adaptive behavior. For this purpose, cognitive scientists

mostly rely on controlled settings in which individuals are asked to accomplish (possibly) highly demanding tasks whose demands are assumed to resemble those of natural environments. Even in the most trivial cases, such as the WCST, optimal performance requires integrated and distributed neurocognitive processes. Moreover, these processes are unlikely to be isolated by simple scoring or aggregate performance measures.

In the current work, we developed and validated a new computational Bayesian model which maps distinct cognitive processes into separable information-theoretic constructs underlying observed adaptive behavior. We argue that these constructs could help describe and investigate the neurocognitive processes underlying adaptive behavior in a principled way.

Furthermore, we couple our computational model with a novel neural density estimation method for simulation-based Bayesian inference ([Radev et al., 2020](#)). Accordingly, we can quantify the entire information contained in the data about the assumed cognitive parameters via a full joint posterior over plausible parameter values. Based on the joint posterior, a representative summary statistic can be computed to simulate the most plausible unfolding of information-theoretic quantities on a trial-by-trial basis.

Several computational models have been proposed to describe and explain performance in the WCST, ranging from behavioral ([Bishara et al., 2010](#); [Gläscher, Adolphs & Tranel, 2019](#); [Steinke et al., 2020](#)) to neural network models ([Dehaene & Changeux, 1991](#); [Amos, 2000](#); [Levine, Parks & Prueitt, 1993](#); [Monchi, Taylor & Dagher, 2000](#)). These models aim to provide psychologically interpretable parameters or biologically inspired network structures, respectively, accounting for specific qualitative patterns of observed data. Behavioral models, in particular, abstract the main cognitive features underlying individual performance in the WCST according to different theoretical frameworks (e.g., attentional updating ([Bishara et al., 2010](#))), or reinforcement learning ([Steinke et al., 2020](#)) and disentangle psychological sub-processes explaining observed task performance. However, the main advantage of our Bayesian model is that it provides both a cognitive and a measurement model which coexist within the overarching theoretical framework of Bayesian brain theories. More precisely, the presented model is specifically designed to capture trial-by-trial fluctuations in information processing as described by second-order information-theoretic quantities. The latter can be seen as a multivariate quantitative account of the interaction between the agent and its environment. Moreover, it is worth noting that such a model representation might not be applicable outside a Bayesian theoretical framework.

Even though our computational model is not a neural model, it might provide a suitable description of cognitive dynamics at a representational and/or a computational level ([Marr, 1982](#)). This description can then be related to neural functioning underlying adaptive behavioral. Indeed, there is some evidence to suggest that neural processes related to belief maintenance/updates and unexpectedness are crucial for performance in the WCST. In particular, brain circuits associated with cognitive control and belief formation, such as the parietal cortex and prefrontal regions, seem to share a functional basis with neural substrates involved in adaptive tasks ([Nour et al., 2018](#)). Prefrontal regions appear to mediate the relation between feedback and belief updates ([Lie et al., 2006](#)) and efficient

functioning in such brain structures seems to be heavily dependent on dopaminergic neuromodulation (Ott & Nieder, 2019). Moreover, the dopaminergic system plays a role in the processing of salient and unexpected environmental stimuli, in learning based on error-related information, and in evaluating candidate actions (Nour et al., 2018; Daw et al., 2011; Gershman, 2018). Accordingly, dopaminergic system functioning has been put in relation with performance in the WCST (Hsieh et al., 2010; Rybakowski et al., 2005) and shown to be critical for the main executive components involved in the task, that is, cognitive flexibility and set-shifting (Bestmann et al., 2014; Stelzel et al., 2010). Further, neural activity in the anterior cingulate cortex (ACC) is increased when a negative feedback occurs in the context of the WCST (Lie et al., 2006). This finding corroborates the view that the ACC is part of an error-detection network which allocates attentional resources to prevent future errors. The ACC might play a crucial role in adaptive functioning by encoding error-related or, more generally, feedback-related information. Thus, it could facilitate the updating of internal environmental models (Rushworth & Behrens, 2008).

The neurobiological evidence suggests that brain networks involved in the WCST might endow adaptive behavior by accounting for maintaining/updating of an internal model of the environment and efficient processing of unexpected information. Is it noteworthy, that these processing aspects are incorporated into our computational framework. At this point, we briefly outline the empirical and theoretical potentials of the proposed computational framework for investigating adaptive functioning and discuss future research vistas.

Model-Based Neuroscience. Recent studies have pointed out the advantage of simultaneously modeling and analyzing neural and behavioral data within a joint modeling framework. In this way, the latter can be used to provide information for the former, as well as the other way around (Turner et al., 2017; Turner et al., 2013; Forstmann et al., 2011). This involves the development of joint models which encode assumptions about the probabilistic relationships between neural and cognitive parameters.

Within our framework, the reconstruction of information-theoretic discrete time series yields a quantitative account of the agent's internal processing of environmental information. Event-related cognitive measures of belief updating, epistemic uncertainty and surprise can be put in relation with neural measurements by explicitly providing a formal account of the statistical dependencies between neural and cognitive (information-theoretic) quantities. In this way, latent cognitive dynamics can be directly related to neural event-related measures (e.g., fMRI, EEG). Applications in which information-theoretic measures are treated as dependent variables in standard statistical analysis are also possible.

Neurological Assessment. Although neuroscientists have considered performance in the WCST as a proxy for measuring high-level cognitive processes, the usual approach to the analysis of human adaptive behavior consists in summarizing response patterns by simple heuristic scoring measures (e.g., occurrences of correct responses and sub-types of errors produced) and classification rules (Flashman, Homer & Freides, 1991). However, the theoretical utility of such a summary approach remains questionable. Indeed, adaptive behavior appears to depend on a complex and intricate interplay between multiple network structures (Barcelo et al., 2006; Monchi et al., 2001; Lie et al., 2006; Barceló & Rubia, 1998; Buchsbaum et al., 2005). This posits a great challenge for disentangling high-level cognitive

constructs at a model level and further investigating their relationship with neurobiological substrates. It appears that standard scoring measures might not be able to fulfil these tasks. Moreover, there is a pronounced lack of anatomical specificity in previous research concerning the neural and functional substrates of the WCST (*Nyhus & Barceló, 2009*).

Thus, there is a need for more sophisticated modeling approaches. For instance, disentangling errors due to perseverative processing of previously relevant environmental models from those due to uncertainty about task environmental states, is important and nontrivial. Sparse and distributed error patterns might depend on several internal model probability configurations. Such internal models are latent, and can only be uncovered through cognitive modeling. Therefore, information-based criteria to response (error) classification can enrich clinical evaluation beyond heuristically motivated criteria.

Generalizability. Another important advantage of the proposed computational framework is that it is not solely confined to the WCST. In fact, one can argue that the seventy-year old WCST does not provide the only or even the most suitable setting for extracting information about cognitive dynamics from general populations or maladaptive behavior in clinical populations. One can envision tasks which embody probabilistic (uncertain) or even chaotic environments (for instance with partially observable or unreliable feedback or partially observable states) and demand integrating information from different modalities (*O'Reilly et al., 2013; Nour et al., 2018*). These settings might prove more suitable for investigating changes in uncertainty-related processing or cross-modal integration than deterministic and fully observable WCST-like settings.

Despite these advantages, our proposed computational framework has certain limitations. A first limitation might concern the fact that the new Bayesian cognitive model accounts for the main dynamics in adaptive tasks by relying on only two parameters. Although such a parsimonious proposal suffices to disentangle latent data-generating processes, a more exhaustive formal description of cognitive sub-components might be envisioned. However, parameter estimation can become challenging in such a scenario, especially when one-dimensional response data is used as a basis for parameter recovery. Second, the information loss parameter appears to be more challenging to estimate than the flexibility parameter in some datasets. There are at least two possible remedies for this problem. On the one hand, global estimation of information loss might be hampered due to the model's current functional (algorithmic) formulation and can therefore be optimized via an alternative formulation/parameterization. On the other hand, it might be the case that the data obtainable in the simple WCST environment is not particularly informative about this parameter and, in general, not suitable for modeling more complex and non-linear cognitive dynamics in general. Future works should therefore focus on designing and exploring more data-rich controlled environments which can provide a better starting point for investigating complex latent cognitive dynamics in a principled way. Additionally, the information loss parameter seems to be less effective in differentiating between substance abusers and healthy controls in the particular sample used in this work. Thus, further model-based analyses on individuals from different clinical populations are needed to fully understand the potential of our 2-parameter model as a clinical neuropsychological tool. Finally, in this work, we did not perform formal model comparison, as this would

require an extensive consideration of various nested and non-nested model within the same theoretical framework and between different theoretical frameworks. We therefore leave this important endeavor for future research.

CONCLUSIONS

In conclusion, the proposed model can be considered as the basis for a (bio)psychometric tool for measuring the dynamics of cognitive processes under changing environmental demands. Furthermore, it can be seen as a step towards a theory-based framework for investigating the relation between such cognitive measures and their neural underpinnings. Further investigations are needed to refine the proposed computational model and systematically explore the advantages of the Bayesian brain theoretical framework for empirical research on high-level cognition.

ACKNOWLEDGEMENTS

We thank Karin Prillinger and Luca D'Alessandro for reading the manuscript and providing useful suggestions which significantly improved the original text.

ADDITIONAL INFORMATION AND DECLARATIONS

Funding

The authors received no funding for this work. Stefan T. Radev was supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation; grant number GRK 2277 “Statistical Modeling in Psychology”).

Grant Disclosures

The following grant information was disclosed by the authors:
Deutsche Forschungsgemeinschaft: GRK 2277.

Competing Interests

The authors declare there are no competing interests.

Author Contributions

- Marco D'Alessandro and Stefan Radev conceived and designed the experiments, performed the experiments, analyzed the data, prepared figures and/or tables, authored or reviewed drafts of the paper, and approved the final draft.
- Andreas Voss and Luigi Lombardi conceived and designed the experiments, authored or reviewed drafts of the paper, and approved the final draft.

Data Availability

The following information was supplied regarding data availability:

The codes for generating data and estimate parameters are available at GitHub: <https://github.com/stefanradev93/DBN>.

The codes of the generative model and parameter estimation are available in “WCST_INN_Final.ipynb”. The data used for parameter estimation are in

“Data128.ipynb” and “MathingMat.ipynb”. The codes for producing the information-theoretic measures from estimated parameters are shown in “ITmeasures.ipynb”. The codes require the library Bayes Flow (<https://github.com/stefanradev93/BayesFlow>).

Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj.10316#supplemental-information>.

REFERENCES

- Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, Devin M, Ghemawat S, Irving G, Isard M, Kudlur M, Levenberg J, Monga R, Moore S, Murray DG, Steiner B, Tucker P, Vasudevan V, Warden P, Wicke M, Yu Y, Zheng X. 2016. Tensorflow: a system for large-scale machine learning. In: *12th USENIX symposium on operating systems design and implementation (OSDI 16)*. 265–283.
- Alvarez JA, Emory E. 2006. Executive function and the frontal lobes: a meta-analytic review. *Neuropsychology Review* **16**(1):17–42 DOI [10.1007/s11065-006-9002-x](https://doi.org/10.1007/s11065-006-9002-x).
- Amos A. 2000. A computational model of information processing in the frontal cortex and basal ganglia. *Journal of Cognitive Neuroscience* **12**(3):505–519 DOI [10.1162/089892900562174](https://doi.org/10.1162/089892900562174).
- Anderson PJ. 2008. Towards a developmental model of executive function. In: Anderson V, Jacobs R, Anderson PJ, eds. *Executive functions and the frontal lobes*. New York: Psychology Press, 3–21.
- Barcelo F, Escera C, Corral MJ, Periañez JA. 2006. Task switching and novelty processing activate a common neural network for cognitive control. *Journal of Cognitive Neuroscience* **18**(10):1734–1748 DOI [10.1162/jocn.2006.18.10.1734](https://doi.org/10.1162/jocn.2006.18.10.1734).
- Barceló F, Rubia FJ. 1998. Non-frontal P3b-like activity evoked by the Wisconsin Card Sorting Test. *Neuroreport* **9**(4):747–751 DOI [10.1097/00001756-199803090-00034](https://doi.org/10.1097/00001756-199803090-00034).
- Bechara A, Damasio H. 2002. Decision-making and addiction (part I): impaired activation of somatic states in substance dependent individuals when pondering decisions with negative future consequences. *Neuropsychologia* **40**(10):1675–1689 DOI [10.1016/S0028-3932\(02\)00015-5](https://doi.org/10.1016/S0028-3932(02)00015-5).
- Berg EA. 1948. A simple objective technique for measuring flexibility in thinking. *The Journal of General Psychology* **39**(1):15–22 DOI [10.1080/00221309.1948.9918159](https://doi.org/10.1080/00221309.1948.9918159).
- Bestmann S, Ruge D, Rothwell J, Galea JM. 2014. The role of dopamine in motor flexibility. *Journal of Cognitive Neuroscience* **27**(2):365–376.
- Bishara AJ, Kruschke JK, Stout JC, Bechara A, McCabe DP, Bussemeyer JR. 2010. Sequential learning models for the Wisconsin card sort task: assessing processes in substance dependent individuals. *Journal of Mathematical Psychology* **54**(1):5–13 DOI [10.1016/j.jmp.2008.10.002](https://doi.org/10.1016/j.jmp.2008.10.002).
- Braff DL, Heaton R, Kuck J, Cullum M, Moranville J, Grant I, Zisook S. 1991. The generalized pattern of neuropsychological deficits in outpatients with chronic schizophrenia with heterogeneous Wisconsin Card Sorting Test results. *Archives of General Psychiatry* **48**(10):891–898 DOI [10.1001/archpsyc.1991.01810340023003](https://doi.org/10.1001/archpsyc.1991.01810340023003).

- Buchsbaum BR, Greer S, Chang W-L, Berman KF. 2005.** Meta-analysis of neuroimaging studies of the Wisconsin Card-Sorting task and component processes. *Human Brain Mapping* **25**(1):35–45 DOI [10.1002/hbm.20128](https://doi.org/10.1002/hbm.20128).
- Buckley CL, Kim CS, McGregor S, Seth AK. 2017.** The free energy principle for action and perception: a mathematical review. *Journal of Mathematical Psychology* **81**:55–79 DOI [10.1016/j.jmp.2017.09.004](https://doi.org/10.1016/j.jmp.2017.09.004).
- Collell G, Fauquet J. 2015.** Brain activity and cognition: a connection from thermodynamics and information theory. *Frontiers in Psychology* **6**:818.
- Cooper R, Fox J, Farringdon J, Shallice T. 1996.** A systematic methodology for cognitive modelling. *Artificial Intelligence* **85**(1–2):3–44 DOI [10.1016/0004-3702\(95\)00112-3](https://doi.org/10.1016/0004-3702(95)00112-3).
- Daw ND, Gershman SJ, Seymour B, Dayan P, Dolan RJ. 2011.** Model-based influences on humans' choices and striatal prediction errors. *Neuron* **69**(6):1204–1215 DOI [10.1016/j.neuron.2011.02.027](https://doi.org/10.1016/j.neuron.2011.02.027).
- Dehaene S, Changeux J-P. 1991.** The Wisconsin Card Sorting Test: theoretical analysis and modeling in a neuronal network. *Cerebral Cortex* **1**(1):62–79 DOI [10.1093/cercor/1.1.62](https://doi.org/10.1093/cercor/1.1.62).
- Feldman H, Friston K. 2010.** Attention, uncertainty, and free-energy. *Frontiers in Human Neuroscience* **4**:215.
- Figueroa IJ, Youmans RJ. 2013.** Failure to maintain set: a measure of distractibility or cognitive flexibility? In: *Proceedings of the human factors and ergonomics society annual meeting*. Los Angeles: Sage Publications, 828–832 DOI [10.1177/1541931213571180](https://doi.org/10.1177/1541931213571180).
- First MB. 1997.** Structured clinical interview for DSM-IV axis I disorders. New York: Biometrics Research, New York State Psychiatric Institute.
- Flashman LA, Homer MD, Freides D. 1991.** Note on scoring perseveration on the Wisconsin Card Sorting Test. *The Clinical Neuropsychologist* **5**(2):190–194 DOI [10.1080/13854049108403303](https://doi.org/10.1080/13854049108403303).
- Forstmann BU, Wagenmakers E-J, Eichele T, Brown S, Serences JT. 2011.** Reciprocal relations between cognitive neuroscience and formal cognitive models: opposites attract? *Trends in Cognitive Sciences* **15**(6):272–279 DOI [10.1016/j.tics.2011.04.002](https://doi.org/10.1016/j.tics.2011.04.002).
- Friston K. 2003.** Learning and inference in the brain. *Neural Networks* **16**(9):1325–1352 DOI [10.1016/j.neunet.2003.06.005](https://doi.org/10.1016/j.neunet.2003.06.005).
- Friston K. 2005.** A theory of cortical responses. *Philosophical Transactions of the Royal Society B: Biological Sciences* **360**(1456):815–836 DOI [10.1098/rstb.2005.1622](https://doi.org/10.1098/rstb.2005.1622).
- Friston K. 2010.** The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience* **11**(2):127–138 DOI [10.1038/nrn2787](https://doi.org/10.1038/nrn2787).
- Friston K, Adams R, Perrinet L, Breakspear M. 2012.** Perceptions as hypotheses: saccades as experiments. *Frontiers in Psychology* **3**:151.
- Friston K, FitzGerald T, Rigoli F, Schwartenbeck P, Pezzulo G. 2017.** Active inference: a process theory. *Neural Computation* **29**(1):1–49 DOI [10.1162/NECO_a_00912](https://doi.org/10.1162/NECO_a_00912).
- Friston K, Kiebel S. 2009.** Predictive coding under the free-energy principle. *Philosophical Transactions of the Royal Society B: Biological Sciences* **364**(1521):1211–1221 DOI [10.1098/rstb.2008.0300](https://doi.org/10.1098/rstb.2008.0300).

- Friston KJ, Daunizeau J, Kilner J, Kiebel SJ. 2010.** Action and behavior: a free-energy formulation. *Biological Cybernetics* **102**(3):227–260
[DOI 10.1007/s00422-010-0364-z](https://doi.org/10.1007/s00422-010-0364-z).
- Gershman SJ. 2018.** The successor representation: its computational logic and neural substrates. *Journal of Neuroscience* **38**(33):7193–7200
[DOI 10.1523/JNEUROSCI.0151-18.2018](https://doi.org/10.1523/JNEUROSCI.0151-18.2018).
- Gläscher J, Adolphs R, Tranel D. 2019.** Model-based lesion mapping of cognitive control using the Wisconsin Card Sorting Test. *Nature Communications* **10**(1):1–12
[DOI 10.1038/s41467-018-07882-8](https://doi.org/10.1038/s41467-018-07882-8).
- Haker H, Schneebeli M, Stephan KE. 2016.** Can Bayesian theories of autism spectrum disorder help improve clinical practice? *Frontiers in Psychiatry* **7**:107.
- Heaton R. 1981.** Wisconsin card sorting test manual; revised and expanded. Odessa, FL: Psychological Assessment Resources.
- Hirsh JB, Mar RA, Peterson JB. 2012.** Psychological entropy: a framework for understanding uncertainty-related anxiety. *Psychological Review* **119**(2):304–320
[DOI 10.1037/a0026767](https://doi.org/10.1037/a0026767).
- Hsieh PC, Yeh TL, Lee IH, Huang HC, Chen PS, Yang YK, Chiu NT, Lu RB, Liao M-H. 2010.** Correlation between errors on the Wisconsin Card Sorting Test and the availability of striatal dopamine transporters in healthy volunteers. *Journal of Psychiatry & Neuroscience* **35**(2):90–94
[DOI 10.1503/jpn.090007](https://doi.org/10.1503/jpn.090007).
- Kersten D, Mamassian P, Yuille A. 2004.** Object perception as Bayesian inference. *Annual Review of Psychology* **55**:271–304
[DOI 10.1146/annurev.psych.55.090902.142005](https://doi.org/10.1146/annurev.psych.55.090902.142005).
- Knill DC, Pouget A. 2004.** The Bayesian brain: the role of uncertainty in neural coding and computation. *TRENDS in Neurosciences* **27**(12):712–719
[DOI 10.1016/j.tins.2004.10.007](https://doi.org/10.1016/j.tins.2004.10.007).
- Koechlin E, Summerfield C. 2007.** An information theoretical approach to prefrontal executive function. *Trends in Cognitive Sciences* **11**(6):229–235
[DOI 10.1016/j.tics.2007.04.005](https://doi.org/10.1016/j.tics.2007.04.005).
- Körding KP, Beierholm U, Ma WJ, Quartz S, Tenenbaum JB, Shams L. 2007.** Causal inference in multisensory perception. *PLOS ONE* **2**(9):e943
[DOI 10.1371/journal.pone.0000943](https://doi.org/10.1371/journal.pone.0000943).
- Kübler A, Murphy K, Garavan H. 2005.** Cocaine dependence and attention switching within and between verbal and visuospatial working memory. *European Journal of Neuroscience* **21**(7):1984–1992
[DOI 10.1111/j.1460-9568.2005.04027.x](https://doi.org/10.1111/j.1460-9568.2005.04027.x).
- Landry O, Al-Taie S. 2016.** A meta-analysis of the Wisconsin Card Sort Task in autism. *Journal of autism and developmental disorders* **46**(4):1220–1235
[DOI 10.1007/s10803-015-2659-3](https://doi.org/10.1007/s10803-015-2659-3).
- Lawson RP, Rees G, Friston KJ. 2014.** An aberrant precision account of autism. *Frontiers in Human Neuroscience* **8**:302.
- Lee TS, Mumford D. 2003.** Hierarchical Bayesian inference in the visual cortex. *JOSA A* **20**(7):1434–1448
[DOI 10.1364/JOSAA.20.001434](https://doi.org/10.1364/JOSAA.20.001434).
- Lee MD, Wagenmakers E-J. 2014.** Bayesian cognitive modeling: a practical course. Cambridge: Cambridge University Press.

- Levine DS, Parks RW, Prueitt PS. 1993.** Methodological and theoretical issues in neural network models of frontal cognitive functions. *International Journal of Neuroscience* 72(3-4):209–233 DOI [10.3109/00207459309024110](https://doi.org/10.3109/00207459309024110).
- Lie C-H, Specht K, Marshall JC, Fink GR. 2006.** Using fMRI to decompose the neural processes underlying the Wisconsin Card Sorting Test. *NeuroImage* 30(3):1038–1049 DOI [10.1016/j.neuroimage.2005.10.031](https://doi.org/10.1016/j.neuroimage.2005.10.031).
- Marr D. 1982.** *Vision: a computational investigation into the human representation and processing of visual information*. San Francisco: W.H. Freeman.
- Monchi O, Petrides M, Petre V, Worsley K, Dagher A. 2001.** Wisconsin Card Sorting revisited: distinct neural circuits participating in different stages of the task identified by event-related functional magnetic resonance imaging. *Journal of Neuroscience* 21(19):7733–7741 DOI [10.1523/JNEUROSCI.21-19-07733.2001](https://doi.org/10.1523/JNEUROSCI.21-19-07733.2001).
- Monchi O, Taylor JG, Dagher A. 2000.** A neural model of working memory processes in normal subjects, Parkinson's disease and schizophrenia for fMRI design and predictions. *Neural Networks* 13(8–9):953–973 DOI [10.1016/S0893-6080\(00\)00058-7](https://doi.org/10.1016/S0893-6080(00)00058-7).
- Nour MM, Dahoun T, Schwartenbeck P, Adams RA, FitzGerald TH, Coello C, Wall MB, Dolan RJ, Howes OD. 2018.** Dopaminergic basis for signaling belief updates, but not surprise, and the link to paranoia. *Proceedings of the National Academy of Sciences of the United States of America* 115(43):E10167–E10176 DOI [10.1073/pnas.1809298115](https://doi.org/10.1073/pnas.1809298115).
- Nyhus E, Barceló F. 2009.** The Wisconsin Card Sorting Test and the cognitive assessment of prefrontal executive functions: a critical update. *Brain and Cognition* 71(3):437–451 DOI [10.1016/j.bandc.2009.03.005](https://doi.org/10.1016/j.bandc.2009.03.005).
- Ott T, Nieder A. 2019.** Dopamine and cognitive control in prefrontal cortex. *Trends in Cognitive Sciences* 23(3):213–234.
- O'Reilly JX, Schüffelgen U, Cuell SF, Behrens TE, Mars RB, Rushworth MF. 2013.** Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proceedings of the National Academy of Sciences* 110(38):E3660–E3669 DOI [10.1073/pnas.1305373110](https://doi.org/10.1073/pnas.1305373110).
- Petzschner FH, Glasauer S, Stephan KE. 2015.** A Bayesian perspective on magnitude estimation. *Trends in Cognitive Sciences* 19(5):285–293 DOI [10.1016/j.tics.2015.03.002](https://doi.org/10.1016/j.tics.2015.03.002).
- Radev ST, Mertens UK, Voss A, Ardizzone L, Kthe U. 2020.** BayesFlow: Learning complex stochastic models with invertible neural networks. ArXiv preprint. [arXiv:2003.06281](https://arxiv.org/abs/2003.06281).
- Rushworth MF, Behrens TE. 2008.** Choice, uncertainty and value in prefrontal and cingulate cortex. *Nature Neuroscience* 11(4):389–397 DOI [10.1038/nn2066](https://doi.org/10.1038/nn2066).
- Rybakowski J, Borkowska A, Czernski P, Kapelski P, Dmitrzak-Weglarczyk M, Hauser J. 2005.** An association study of dopamine receptors polymorphisms and the Wisconsin Card Sorting Test in schizophrenia. *Journal of Neural Transmission* 112(11):1575–1582 DOI [10.1007/s00702-005-0292-6](https://doi.org/10.1007/s00702-005-0292-6).
- Sayood K. 2018.** Information theory and cognition: a review. *Entropy* 20(9):706 DOI [10.3390/e20090706](https://doi.org/10.3390/e20090706).

- Schwartenbeck P, FitzGerald TH, Dolan R. 2016.** Neural signals encoding shifts in beliefs. *NeuroImage* 125:578–586 DOI [10.1016/j.neuroimage.2015.10.067](https://doi.org/10.1016/j.neuroimage.2015.10.067).
- Singh S, Aich TK, Bhattarai R. 2017.** Wisconsin Card Sorting Test performance impairment in schizophrenia: an Indian study report. *Indian journal of psychiatry* 59(1):88–93 DOI [10.4103/0019-5545.204440](https://doi.org/10.4103/0019-5545.204440).
- Sisson SA, Fan Y. 2011.** *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839].
- Steinke A, Lange F, Seer C, Hendel MK, Kopp B. 2020.** Computational modeling for neuropsychological assessment of bradyphrenia in Parkinsons disease. *Journal of Clinical Medicine* 9(4):1158 DOI [10.3390/jcm9041158](https://doi.org/10.3390/jcm9041158).
- Stelzel C, Basten U, Montag C, Reuter M, Fiebach CJ. 2010.** Frontostriatal involvement in task switching depends on genetic differences in d2 receptor density. *Journal of Neuroscience* 30(42):14205–14212 DOI [10.1523/JNEUROSCI.1062-10.2010](https://doi.org/10.1523/JNEUROSCI.1062-10.2010).
- Stephan KE, Baldeweg T, Friston KJ. 2006.** Synaptic plasticity and dysconnection in schizophrenia. *Biological Psychiatry* 59(10):929–939 DOI [10.1016/j.biopsych.2005.10.005](https://doi.org/10.1016/j.biopsych.2005.10.005).
- Strange BA, Duggins A, Penny W, Dolan RJ, Friston KJ. 2005.** Information theory, novelty and hippocampal responses: unpredicted or unpredictable? *Neural Networks* 18(3):225–230 DOI [10.1016/j.neunet.2004.12.004](https://doi.org/10.1016/j.neunet.2004.12.004).
- Sun R. 2009.** Theoretical status of computational cognitive modeling. *Cognitive Systems Research* 10(2):124–140 DOI [10.1016/j.cogsys.2008.07.002](https://doi.org/10.1016/j.cogsys.2008.07.002).
- Talts S, Betancourt M, Simpson D, Vehtari A, Gelman A. 2018.** Validating Bayesian inference algorithms with simulation-based calibration. ArXiv preprint. [arXiv:1804.06788](https://arxiv.org/abs/1804.06788).
- Turner BM, Forstmann BU, Love BC, Palmeri TJ, Van Maanen L. 2017.** Approaches to analysis in model-based cognitive neuroscience. *Journal of Mathematical Psychology* 76:65–79 DOI [10.1016/j.jmp.2016.01.001](https://doi.org/10.1016/j.jmp.2016.01.001).
- Turner BM, Forstmann BU, Wagenmakers E-J, Brown SD, Sederberg PB, Steyvers M. 2013.** A Bayesian framework for simultaneously modeling neural and behavioral data. *NeuroImage* 72:193–206 DOI [10.1016/j.neuroimage.2013.01.048](https://doi.org/10.1016/j.neuroimage.2013.01.048).
- Willuhn I, Sun W, Steiner H. 2003.** Topography of cocaine-induced gene regulation in the rat striatum: relationship to cortical inputs and role of behavioural context. *European Journal of Neuroscience* 17(5):1053–1066 DOI [10.1046/j.1460-9568.2003.02525.x](https://doi.org/10.1046/j.1460-9568.2003.02525.x).
- Zakzanis KK. 1998.** The subcortical dementia of Huntington’s disease. *Journal of Clinical and Experimental Neuropsychology* 20(4):565–578 DOI [10.1076/j.jcen.20.4.565.1468](https://doi.org/10.1076/j.jcen.20.4.565.1468).

DECLARATION IN ACCORDANCE TO § 8 (1) C)
AND (D) OF THE DOCTORAL DEGREE
REGULATION OF THE FACULTY



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

FAKULTÄT FÜR VERHALTENS-
UND EMPIRISCHE KULTURWISSENSCHAFTEN

**Promotionsausschuss der Fakultät für Verhaltens- und Empirische Kulturwissenschaften
der Ruprecht-Karls-Universität Heidelberg**
Doctoral Committee of the Faculty of Behavioural and Cultural Studies of Heidelberg University

**Erklärung gemäß § 8 (1) c) der Promotionsordnung der Universität Heidelberg
für die Fakultät für Verhaltens- und Empirische Kulturwissenschaften**
Declaration in accordance to § 8 (1) c) of the doctoral degree regulation of Heidelberg University,
Faculty of Behavioural and Cultural Studies

Ich erkläre, dass ich die vorgelegte Dissertation selbstständig angefertigt, nur die angegebenen Hilfsmittel benutzt und die Zitate gekennzeichnet habe.

I declare that I have made the submitted dissertation independently, using only the specified tools and have correctly marked all quotations.

**Erklärung gemäß § 8 (1) d) der Promotionsordnung der Universität Heidelberg
für die Fakultät für Verhaltens- und Empirische Kulturwissenschaften**
Declaration in accordance to § 8 (1) d) of the doctoral degree regulation of Heidelberg University,
Faculty of Behavioural and Cultural Studies

Ich erkläre, dass ich die vorgelegte Dissertation in dieser oder einer anderen Form nicht anderweitig als Prüfungsarbeit verwendet oder einer anderen Fakultät als Dissertation vorgelegt habe.

I declare that I did not use the submitted dissertation in this or any other form as an examination paper until now and that I did not submit it in another faculty.

Vorname Nachname

First name Family name Stefan Radev

Datum, Unterschrift

Date, Signature 15.03.2021, *Stefan Radev*