

# The Web Data Commons Structured Data Extraction

Anna Primpeli, Robert Meusel, Christian Bizer, Heiner Stuckenschmidt  
Data and Web Science Group, University of Mannheim

Extracting structured data out of 3.1 billion web pages from the biggest public web corpus  
- A use case of the ViCE project -

## Structured Data on the Web

A Semantic Annotation Example

Les Misérables

itemscope itemtype="http://schema.org/Movie"

itemprop="name"

itemprop="aggregateRating" itemscope itemtype="http://schema.org/AggregateRating" itemprop="ratingValue" content="75"

itemprop="actor" itemscope itemtype="http://schema.org/Person" itemprop="name"

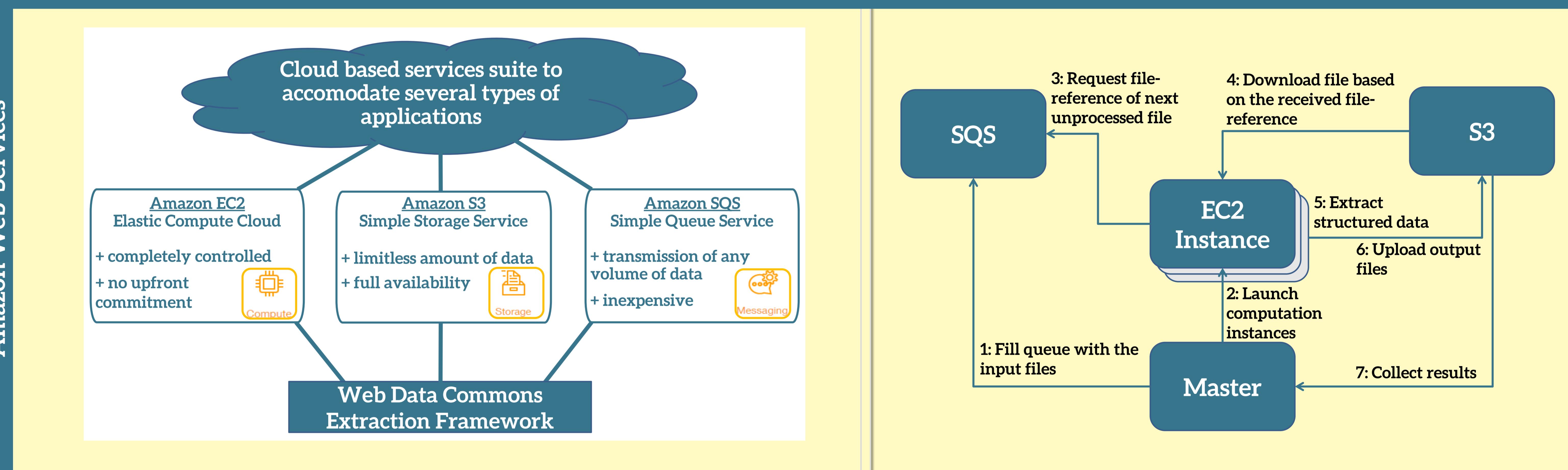
GeoCoordinates LocalBusiness Person Product Organization

Evolution of #Domains (in millions) with semantic annotations

Year	RDFa	Microdata	hCard	JSON-LD
2012	0.5	0.1	1.5	0.0
2013	0.5	0.5	1.0	0.0
2014	0.5	0.8	1.1	0.0
2015	0.5	1.0	1.2	0.5
2016	1.0	2.5	1.6	2.1

Common Topics Evolution

## Cloud Based Extraction



## Web Data Commons, Requirements and Results

The Web Data Commons project extracts structured data from the Common Crawl, the largest web corpus available to the public, and provides the extracted data for public download. Different types of data have been generated since 2012 till today including the following:

### RDFa, Microdata, and Microformat Data Sets

The extracted data as well as statistics about the deployment of different formats are published

### Web Tables Corpus

Dataset containing 147 million relational web tables

### Hyperlink Graph

Covers 3.5 billion web pages and 128 billion hyperlinks between these pages

### WebsA Database

Contains more than 400 million hypernymy relations

### Product Data Corpus

Contains over 5.6 million product records retrieved from the most visited 32 shopping websites

### Processing

- 100 computation nodes
- 8 cores, 64-bit, 2.8 GHz
- 15 GB RAM
- 2x80 GB SSD

### Data Storage

- 56 TB of input data from Common Crawl
- 1 TB of resulting data

### Time & Cost

- 50 hours
- 650 \$ total cost

### Input Data

56 TB  
34 million Domains  
3.1 billion URLs

### Extracted Data

1 TB Extracted Data  
5.6 million Domains  
1.2 billion URLs

### Extracted Data as NQuads

44 billion Triples  
9.5 billion Entities



WDC Structured Data  
<http://webdatacommons.org/structureddata>

