

# Web Archive Analytics

---

Martin Potthast  
Universität Leipzig  
webis.de

ZIH Colloquium · Dresden · April 28, 2022

# Outline

- ① The Global Datasphere
- ② The Internet Archive
- ③ Web Archive Analytics @ Webis
- ④ Web Archive Processing
- ⑤ Webis Archive Research



# The Global Datasphere



# The Global Datasphere

*“A measure of all new data captured, created, and replicated in a single year.”*

[IDC, 2018]



*“... images and videos on mobile phones uploaded to YouTube, digital movies populating the pixels of our high-definition TVs, security footage at airports and major events such as the Olympic Games, subatomic collisions recorded by the Large Hadron Collider at CERN, banking data swiped in an ATM, transponders recording highway tolls, voice calls zipping through digital phone lines, texting as a widespread means of communications, ...”*

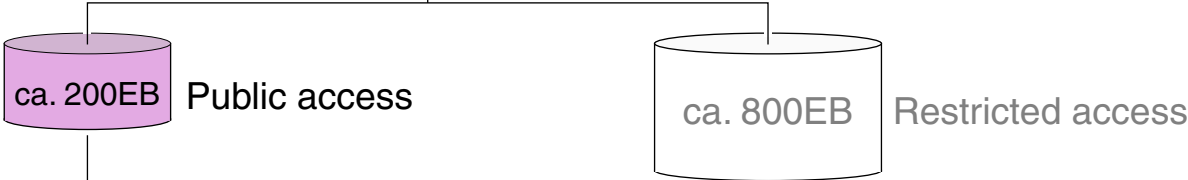
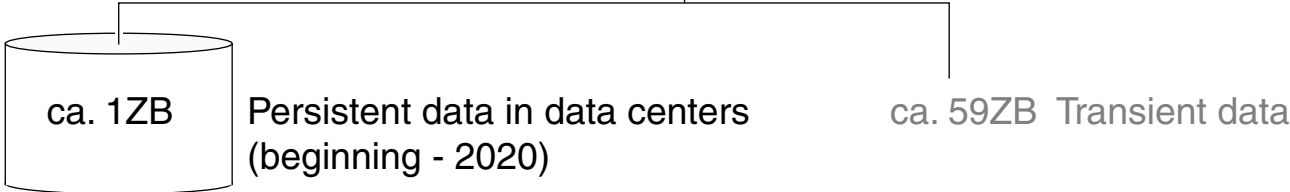
[IDC, 2012]

© WEBIS 2022

# The Global Datasphere in 2020



ca. 59ZB Entire data generated in 2020



- Web pages (< 1EB)
- Books and texts
- Audio recordings
- Videos
- Images
- Software programs

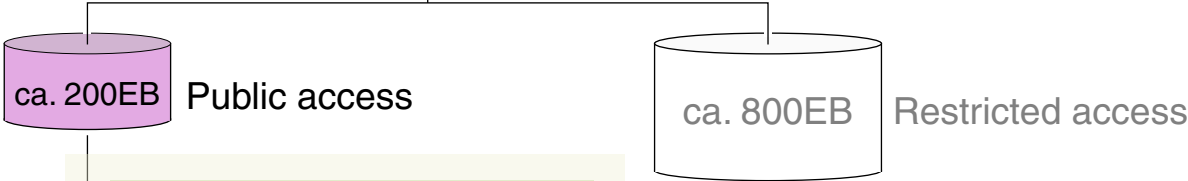
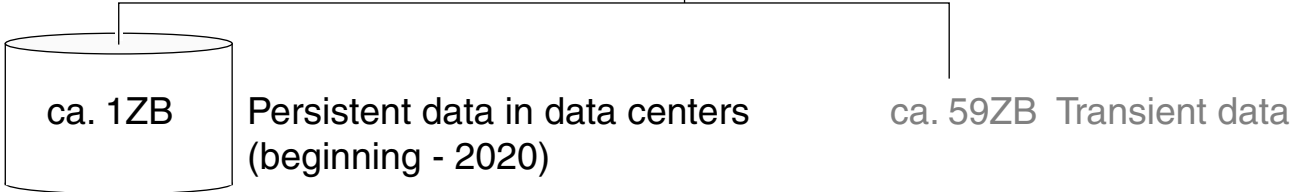
- Data of individuals
- Data in enterprises
- Data of public bodies

|     |   |           |       |
|-----|---|-----------|-------|
| 1GB | = | $10^9$    | Bytes |
| 1TB | = | $10^{12}$ | Bytes |
| 1PB | = | $10^{15}$ | Bytes |
| 1EB | = | $10^{18}$ | Bytes |
| 1ZB | = | $10^{21}$ | Bytes |

# The Global Datasphere in 2020



ca. 59ZB Entire data generated in 2020



**Public access**

- Web pages (< 1EB) 
- Books and texts
- Audio recordings
- Videos
- Images
- Software programs

INTERNET ARCHIVE 

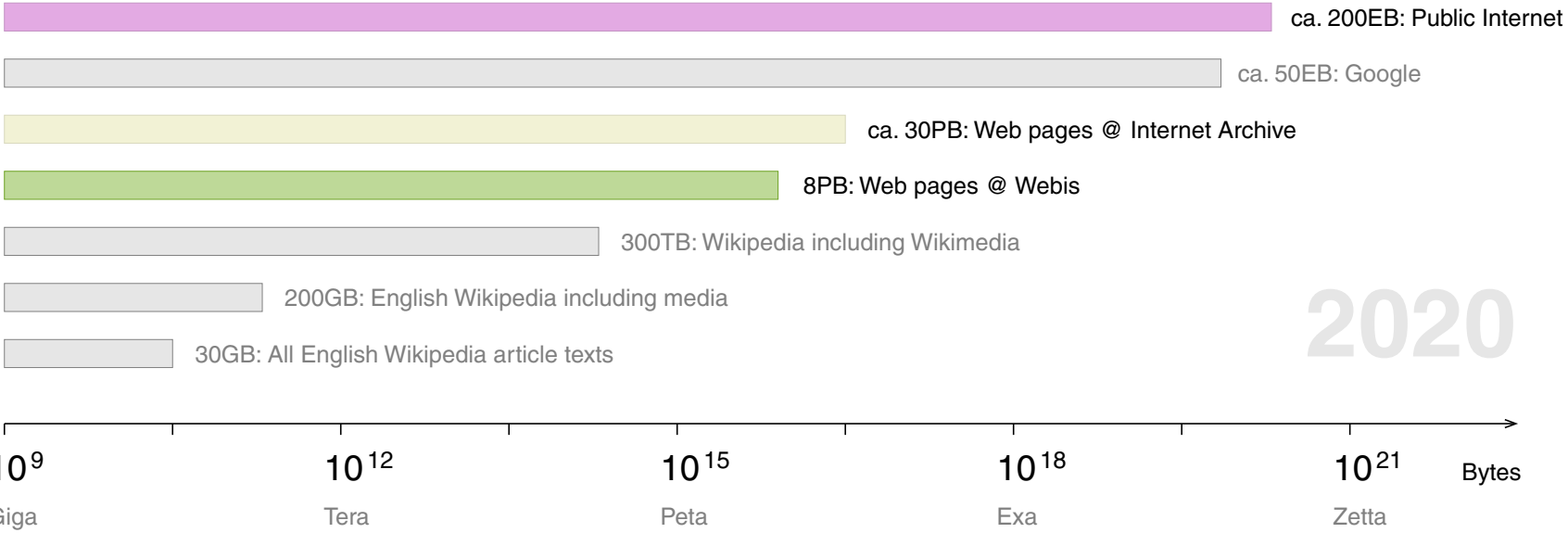
- Restricted access**
- Data of individuals
  - Data in enterprises
  - Data of public bodies

|     |   |                  |       |
|-----|---|------------------|-------|
| 1GB | = | 10 <sup>9</sup>  | Bytes |
| 1TB | = | 10 <sup>12</sup> | Bytes |
| 1PB | = | 10 <sup>15</sup> | Bytes |
| 1EB | = | 10 <sup>18</sup> | Bytes |
| 1ZB | = | 10 <sup>21</sup> | Bytes |

Basis: IDC (2014-20) • Seagate (2018-20) • Cisco Systems (2018) • Statista (2020) • Domo Inc. (2018-20)

# The Global Datasphere in 2020

## Relating Data Source Sizes

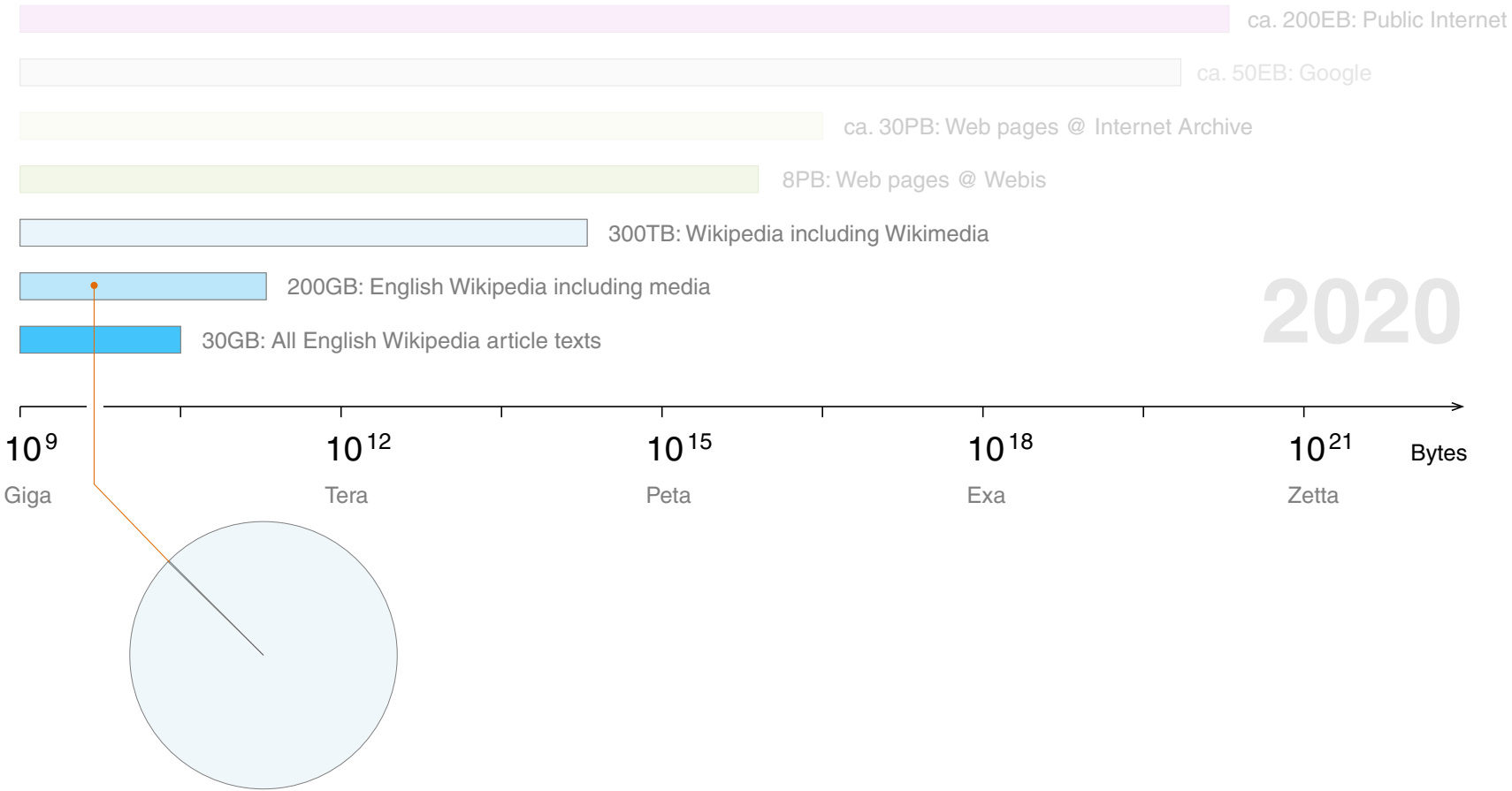


2020



# The Global Datasphere in 2020

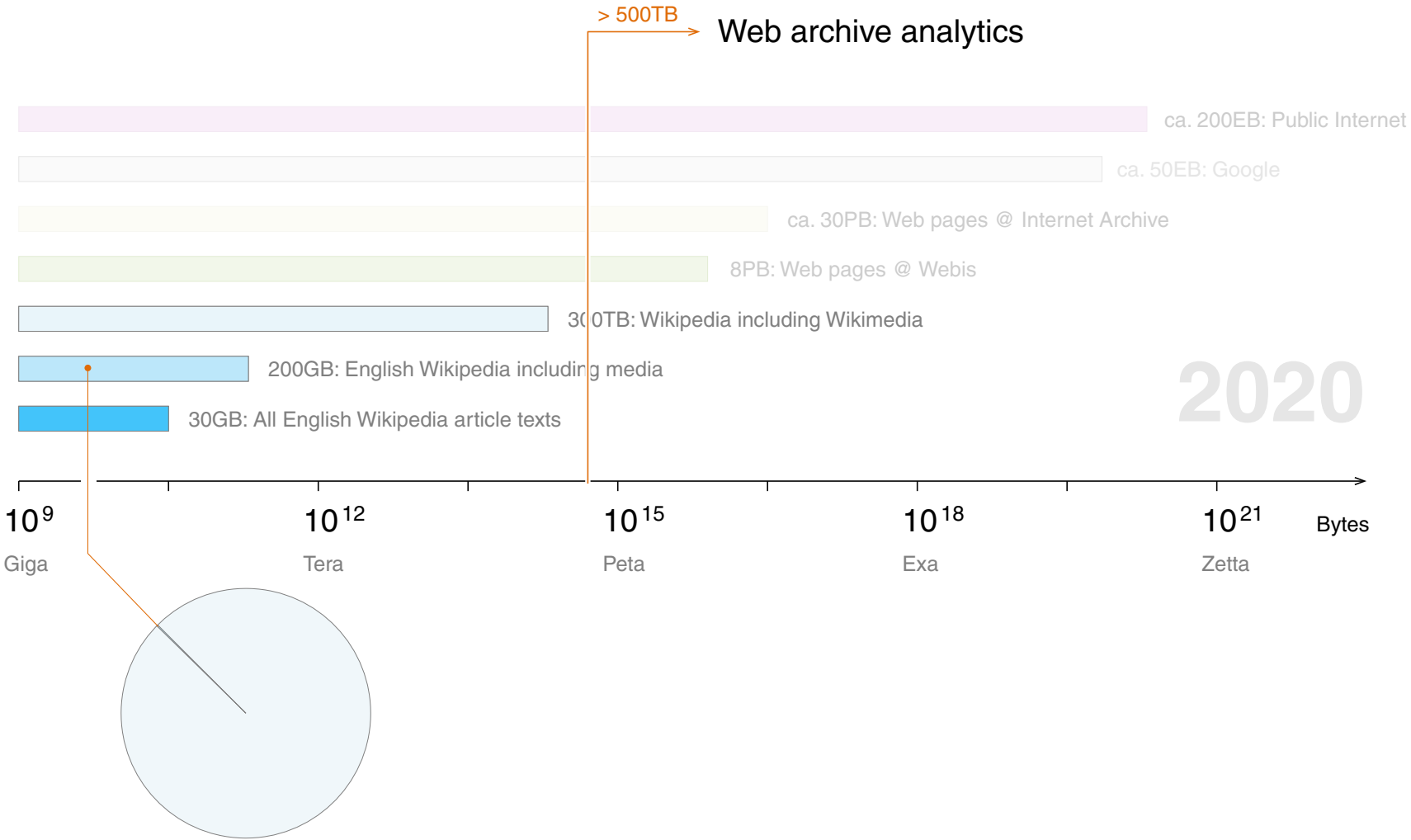
## Relating Data Source Sizes



2020

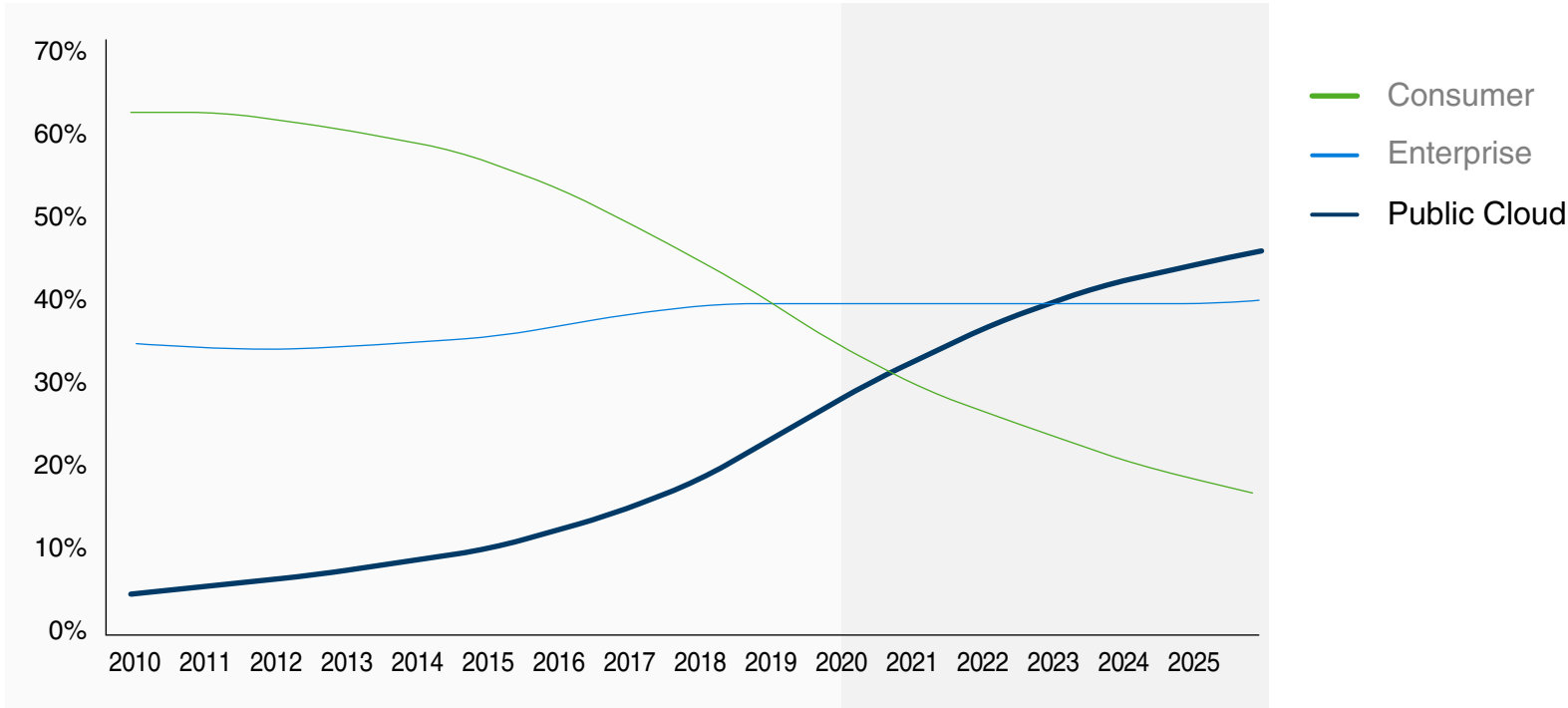
# The Global Datasphere in 2020

## Relating Data Source Sizes



# The Global Datasphere in 2020

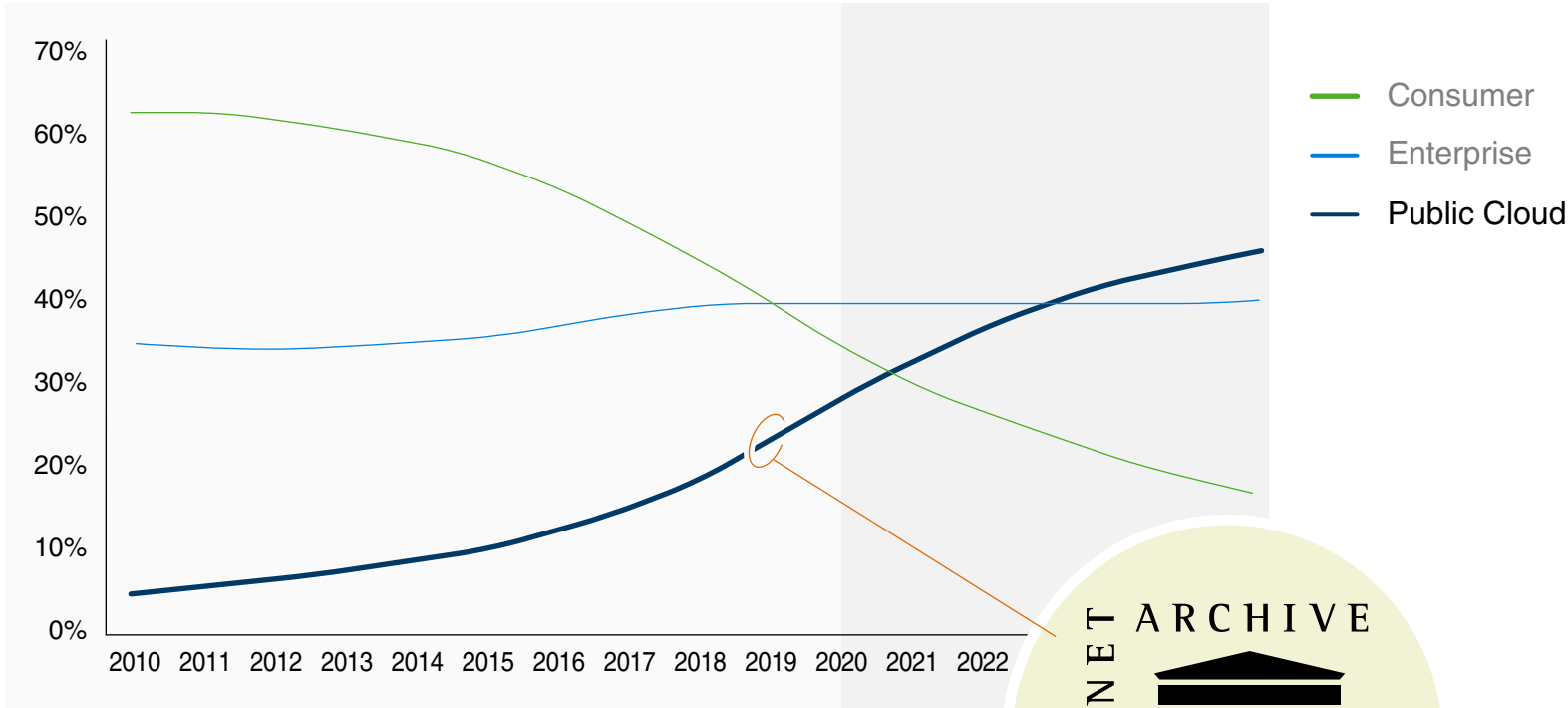
## Where is the Data Stored?



**Basis:** Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, May 2020.

# The Global Datasphere in 2020

## Where is the Data Stored?



Among others:



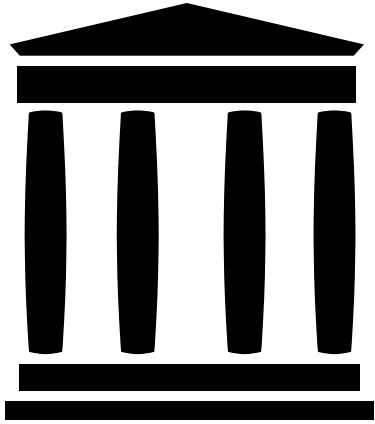
Basis: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, May 2020.




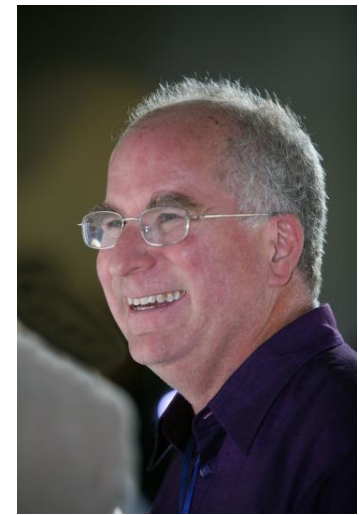
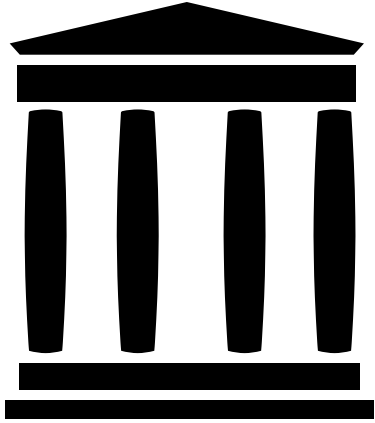
The Internet Archive



- ❑ Founded 1996 by Brewster Kahle
- ❑ For all things digital:
  - 477 billion web pages (ca. 30PB) – accessible via the INTERNET ARCHIVE WayBackMachine
  - 20 million books and texts
  - 4.5 million audio recordings (including 180,000 live concerts)
  - 4 million videos (including 1.6 million Television News programs)
  - 3 million images
  - 200,000 software programs



- ❑ Founded 1996 by Brewster Kahle
- ❑ For all things digital:
  - 477 billion web pages (ca. 30PB) – accessible via the 
  - 20 million books and texts
  - 4.5 million audio recordings (including 180,000 live concerts)
  - 4 million videos (including 1.6 million Television News programs)
  - 3 million images
  - 200,000 software programs



Mission: “Universal access to all knowlege.”

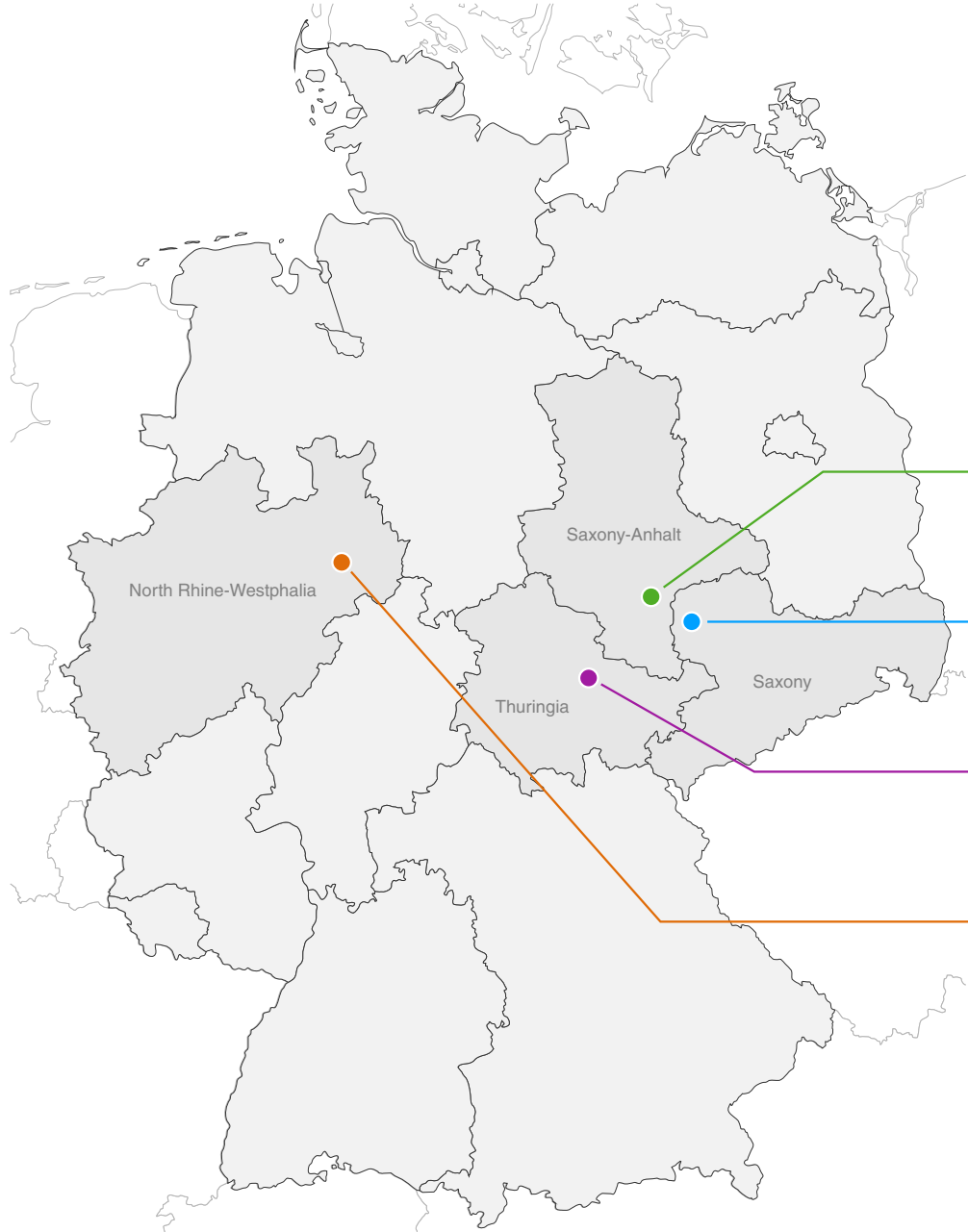
- ❑ One full copy in San Francisco
- ❑ Part at the new Library of Alexandria
- ❑ Part in Amsterdam
- ❑ Copy representative portion (8PB) to the Digital Bauhaus Lab / Webis group:

[[archive.webis.de](https://archive.webis.de)]





Web Archive Analytics @ Webis



MLU Halle-Wittenberg

Prof. Dr. Matthias Hagen

Leipzig University

Prof. Dr. Martin Potthast

Bauhaus-Universität Weimar

Prof. Dr. Benno Stein


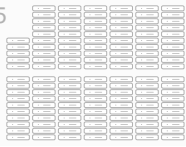









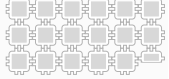
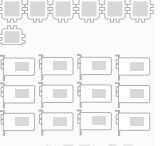
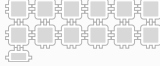


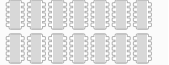



Paderborn University

Prof. Dr. Henning Wachsmuth

# Webis Data Center (Digital Bauhaus Lab)



# Webis Data Center (Digital Bauhaus Lab)

|           | $\alpha$ -web [2009]                                                                                            | $\beta$ -web [2015]                                                                                                | $\gamma$ -web [2016 + 2021]                                                                                              | $\delta$ -web [2018]                                                                                                  | $\epsilon$ -web [2020]                                                                                             |
|-----------|-----------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------|
| Nodes     | 44                             | 135                               | 9                                      | 78                                 | 55                              |
| Disk [PB] | 0.2                            | 4.1                               | 0.08                                    | 12                                 | 0.1                             |
| Cores     | 176 <br><br>$\cong 3.2$ TFLOPs | 1,740 <br><br>$\cong 67.4$ TFLOPs | 672 + 227,328 <br><br>$\cong 8$ PFLOPs | 1,248 <br><br>$\cong 119.8$ TFLOPs | 1,100 <br><br>$\cong 44$ TFLOPs |
| RAM [TB]  | 0.8                            | 28                                | 7.5                                    | 10                                 | 7                               |

## Typical research:

$\alpha$ -Web. Teaching, Staging environment


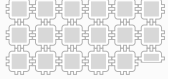
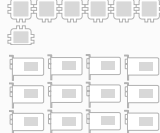


$\beta$ -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

$\gamma$ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

$\delta$ -Web. Web archiving, Virtualization (storage)

$\epsilon$ -Web. Search index construction, Argument search

# Webis Data Center (Digital Bauhaus Lab)

|           | $\alpha$ -web [2009]                                                                                            | $\beta$ -web [2015]                                                                                                | $\gamma$ -web [2016 + 2021]                                                                                              | $\delta$ -web [2018]                                                                                                  | $\epsilon$ -web [2020]                                                                                             |
|-----------|-----------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------|
| Nodes     | 44                             | 135                               | 9                                      | 78                                 | 55                              |
| Disk [PB] | 0.2                            | 4.1                               | 0.08                                    | 12                                 | 0.1                             |
| Cores     | 176 <br><br>$\cong 3.2$ TFLOPs | 1,740 <br><br>$\cong 67.4$ TFLOPs | 672 + 227,328 <br><br>$\cong 8$ PFLOPs | 1,248 <br><br>$\cong 119.8$ TFLOPs | 1,100 <br><br>$\cong 44$ TFLOPs |
| RAM [TB]  | 0.8                            | 28                                | 7.5                                    | 10                                 | 7                               |

## Typical research:

$\alpha$ -Web. Teaching, Staging environment


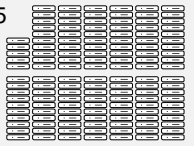

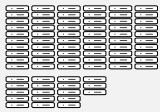







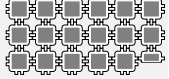
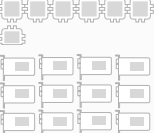



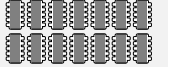



$\beta$ -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

$\gamma$ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

$\delta$ -Web. Web archiving, Virtualization (storage)

$\epsilon$ -Web. Search index construction, Argument search

# Webis Data Center (Digital Bauhaus Lab)

|           | $\alpha$ -web [2009]                                                                                              | $\beta$ -web [2015]                                                                                                  | $\gamma$ -web [2016 + 2021]                                                                                                | $\delta$ -web [2018]                                                                                                    | $\epsilon$ -web [2020]                                                                                               |
|-----------|-------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------|
| Nodes     | 44                               | 135                                 | 9                                        | 78                                   | 55                                |
| Disk [PB] | 0.2                              | 4.1                                 | 0.08                                      | 12                                   | 0.1                               |
| Cores     | 176 <br><br>$\approx 3.2$ TFLOPs | 1,740 <br><br>$\approx 67.4$ TFLOPs | 672 + 227,328 <br><br>$\approx 8$ PFLOPs | 1,248 <br><br>$\approx 119.8$ TFLOPs | 1,100 <br><br>$\approx 44$ TFLOPs |
| RAM [TB]  | 0.8                              | 28                                  | 7.5                                      | 10                                   | 7                                 |

## Typical research:

$\alpha$ -Web. Teaching, Staging environment

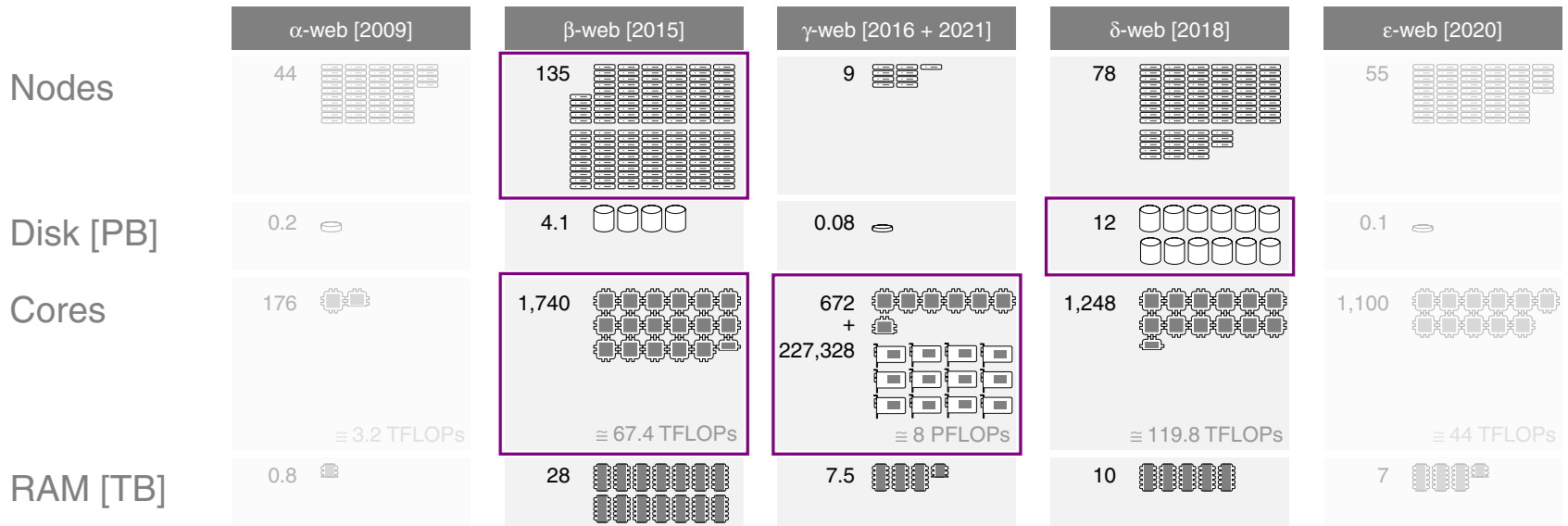
$\beta$ -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

$\gamma$ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

$\delta$ -Web. Web archiving, Virtualization (storage)

$\epsilon$ -Web. Search index construction, Argument search

# Webis Data Center (Digital Bauhaus Lab)



## Typical research:

$\alpha$ -Web. Teaching, Staging environment

$\beta$ -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

$\gamma$ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

$\delta$ -Web. Web archiving, Virtualization (storage)

$\epsilon$ -Web. Search index construction, Argument search

# Webis Analytics Stack

Data  
Consumption  
Layer

Data  
Analytics  
Layer

Data  
Management  
Layer

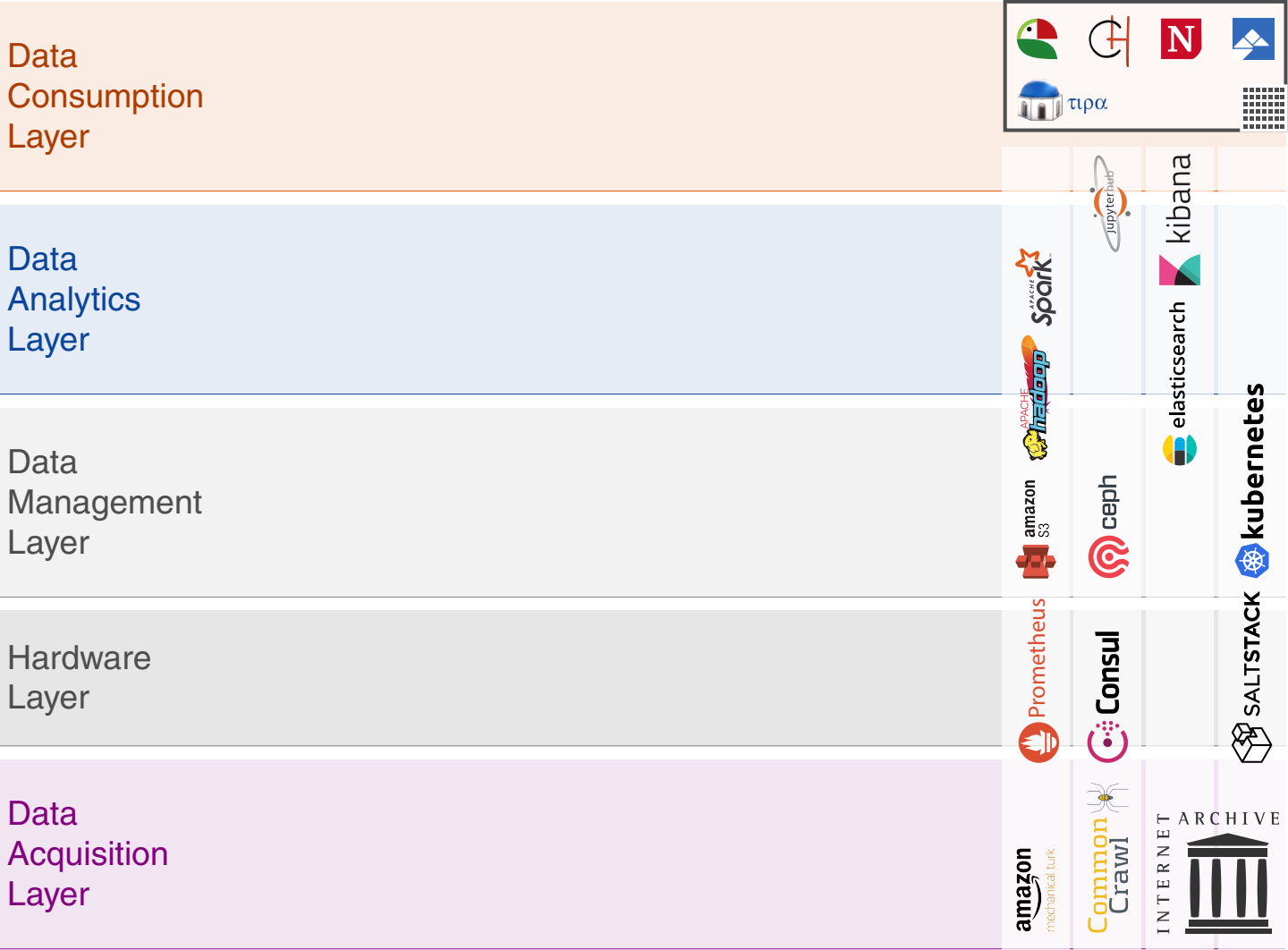
Hardware  
Layer

Data  
Acquisition  
Layer



# Webis Analytics Stack

Vendor stack



# Webis Analytics Stack

|                                      | Technology stack                                                                                                                               | Vendor stack |
|--------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|--------------|
| <p><b>Data Consumption Layer</b></p> | <ul style="list-style-type: none"> <li>- Visual analytics</li> <li>- Immersive technologies</li> <li>- Intelligent agents</li> </ul>           |              |
| <p><b>Data Analytics Layer</b></p>   | <ul style="list-style-type: none"> <li>- Distributed learning</li> <li>- State-space search</li> <li>- Symbolic inference</li> </ul>           |              |
| <p><b>Data Management Layer</b></p>  | <ul style="list-style-type: none"> <li>- Key-value store</li> <li>- RDF triple store</li> <li>- Graph store</li> <li>- Object store</li> </ul> |              |
| <p><b>Hardware Layer</b></p>         | <ul style="list-style-type: none"> <li>- Orchestration</li> <li>- Parallelization</li> <li>- Virtualization</li> </ul>                         |              |
| <p><b>Data Acquisition Layer</b></p> | <ul style="list-style-type: none"> <li>- Distant supervision</li> <li>- Crowdsourcing</li> <li>- Crawling and archiving</li> </ul>             |              |

# Webis Analytics Stack

## Task Stack

## Technology stack

## Vendor stack

|                               | Task Stack                                                                                                                                    | Technology stack                                                                                                                               | Vendor stack |
|-------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|--------------|
| <b>Data Consumption Layer</b> | <ul style="list-style-type: none"> <li>- Query and explore</li> <li>- Visualize and interact</li> <li>- Explain and justify</li> </ul>        | <ul style="list-style-type: none"> <li>- Visual analytics</li> <li>- Immersive technologies</li> <li>- Intelligent agents</li> </ul>           |              |
| <b>Data Analytics Layer</b>   | <ul style="list-style-type: none"> <li>- Diagnose and reason</li> <li>- Structure identification</li> <li>- Structure verification</li> </ul> | <ul style="list-style-type: none"> <li>- Distributed learning</li> <li>- State-space search</li> <li>- Symbolic inference</li> </ul>           |              |
| <b>Data Management Layer</b>  | <ul style="list-style-type: none"> <li>- Provenance tracking</li> <li>- Normalization</li> <li>- Cleansing</li> </ul>                         | <ul style="list-style-type: none"> <li>- Key-value store</li> <li>- RDF triple store</li> <li>- Graph store</li> <li>- Object store</li> </ul> |              |
| <b>Hardware Layer</b>         | <ul style="list-style-type: none"> <li>- Monitoring</li> <li>- Replication</li> </ul>                                                         | <ul style="list-style-type: none"> <li>- Orchestration</li> <li>- Parallelization</li> <li>- Virtualization</li> </ul>                         |              |
| <b>Data Acquisition Layer</b> | <ul style="list-style-type: none"> <li>- Replay</li> <li>- Collect</li> <li>- Log</li> </ul>                                                  | <ul style="list-style-type: none"> <li>- Distant supervision</li> <li>- Crowdsourcing</li> <li>- Crawling and archiving</li> </ul>             |              |

# Webis Analytics Stack

## Task Stack

## Technology stack

## Vendor stack

## Roles

|                               | Task Stack                                                                                                                                    | Technology stack                                                                                                                               | Vendor stack | Roles                                                                                                              |
|-------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|--------------|--------------------------------------------------------------------------------------------------------------------|
| <b>Data Consumption Layer</b> | <ul style="list-style-type: none"> <li>- Query and explore</li> <li>- Visualize and interact</li> <li>- Explain and justify</li> </ul>        | <ul style="list-style-type: none"> <li>- Visual analytics</li> <li>- Immersive technologies</li> <li>- Intelligent agents</li> </ul>           |              | <p>Experts:</p> <ul style="list-style-type: none"> <li>- IR</li> <li>- NLP</li> <li>- CSS</li> <li>- VA</li> </ul> |
| <b>Data Analytics Layer</b>   | <ul style="list-style-type: none"> <li>- Diagnose and reason</li> <li>- Structure identification</li> <li>- Structure verification</li> </ul> | <ul style="list-style-type: none"> <li>- Distributed learning</li> <li>- State-space search</li> <li>- Symbolic inference</li> </ul>           |              | Data scientist                                                                                                     |
| <b>Data Management Layer</b>  | <ul style="list-style-type: none"> <li>- Provenance tracking</li> <li>- Normalization</li> <li>- Cleansing</li> </ul>                         | <ul style="list-style-type: none"> <li>- Key-value store</li> <li>- RDF triple store</li> <li>- Graph store</li> <li>- Object store</li> </ul> |              | Data engineer                                                                                                      |
| <b>Hardware Layer</b>         | <ul style="list-style-type: none"> <li>- Monitoring</li> <li>- Replication</li> </ul>                                                         | <ul style="list-style-type: none"> <li>- Orchestration</li> <li>- Parallelization</li> <li>- Virtualization</li> </ul>                         |              |                                                                                                                    |
| <b>Data Acquisition Layer</b> | <ul style="list-style-type: none"> <li>- Replay</li> <li>- Collect</li> <li>- Log</li> </ul>                                                  | <ul style="list-style-type: none"> <li>- Distant supervision</li> <li>- Crowdsourcing</li> <li>- Crawling and archiving</li> </ul>             |              | Data scientist                                                                                                     |

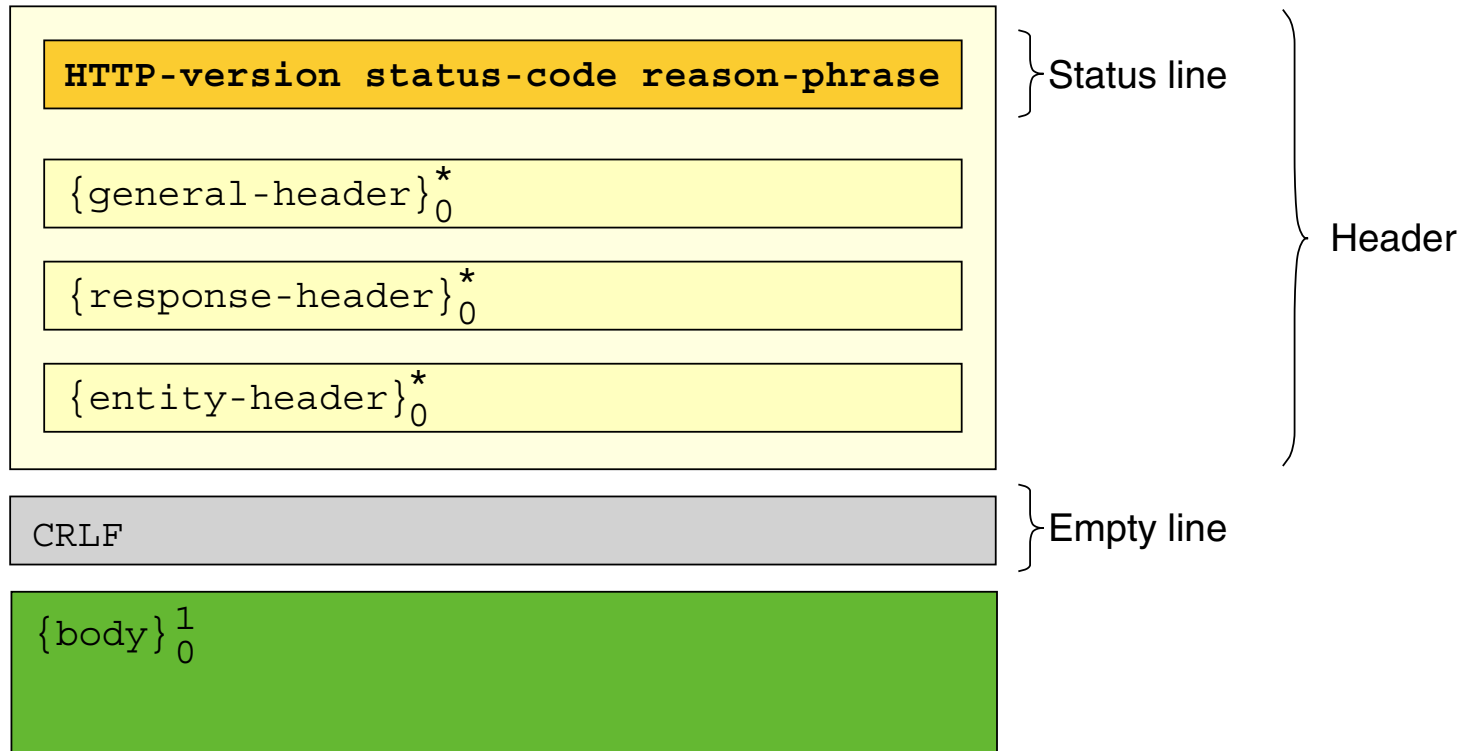


## Web Archive Processing

# Web Archive Data

## WARC Standard

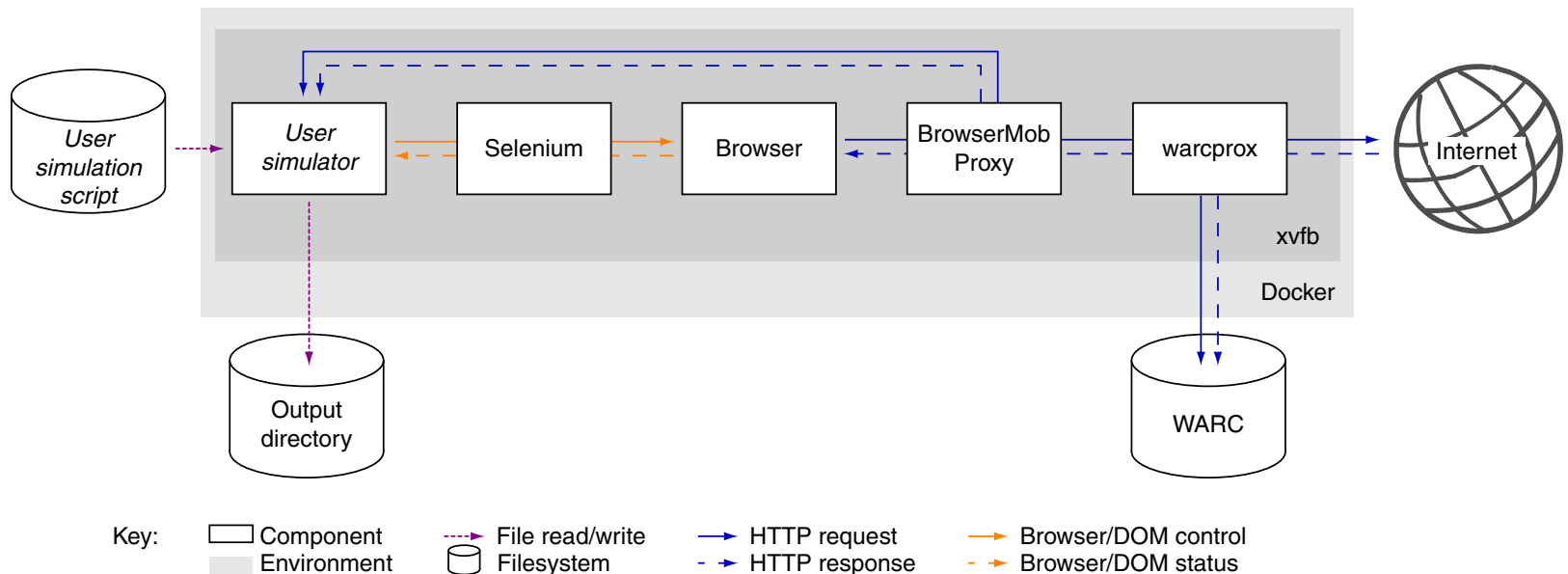
- ❑ WARC is a standard format for web archives.
- ❑ A WARC file consists of a zipped sequence of WARC records. (~1 GiB / file)
- ❑ A WARC record corresponds to one HTTP request/response for a given URI:



# Web Archive Data

## Web Archiving

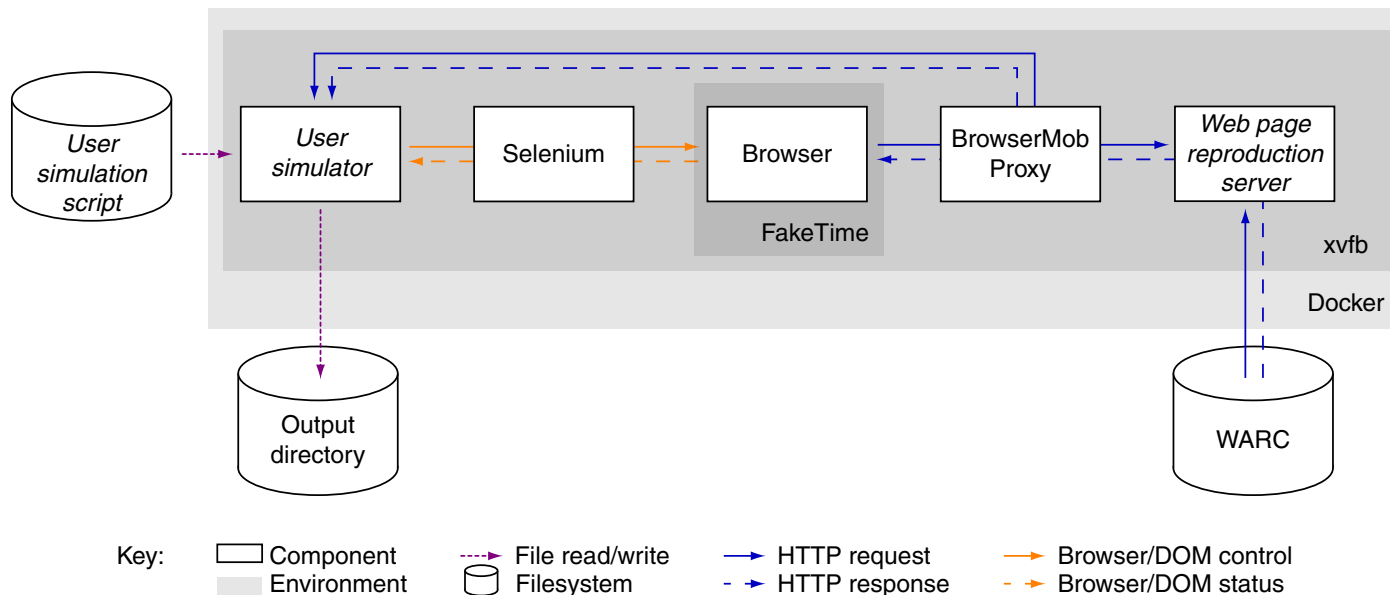
- ❑ A web page: Record all HTTP communication between browser and server.
- ❑ A browser is simulated to ensure the human-readable version is obtained.
- ❑ During web crawling, a web archiver “browses” every crawled page.



# Web Archive Data

## Web Archiving

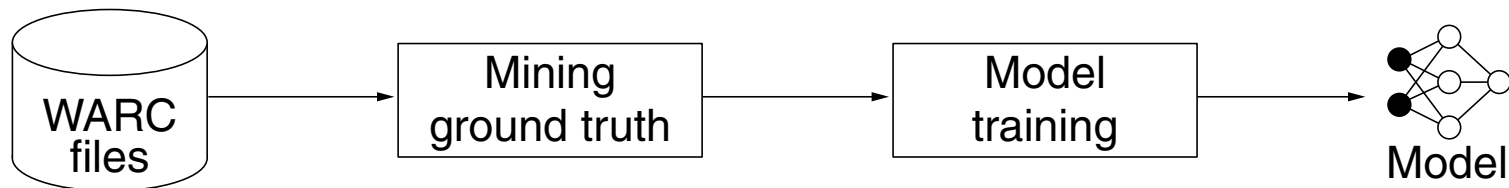
- ❑ A web page: Record all HTTP communication between browser and server.
- ❑ A browser is simulated to ensure the human-readable version is obtained.
- ❑ During web crawling, a web archiver “browses” every crawled page.





# Web Archive Processing

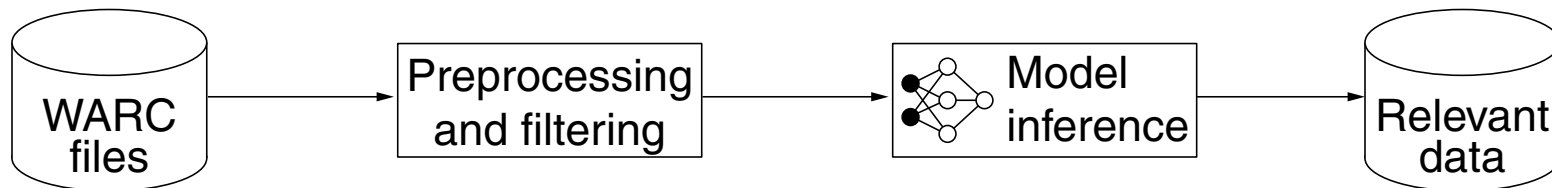
## Streamed Model Training Pipeline



- ❑ Given a learning task and ground truth within WARC files, train a model.  
Only a fraction of the records within the WARC files are ground truth.
- ❑ Goal: Training at web scale (billions of WARC files)

# Web Archive Processing

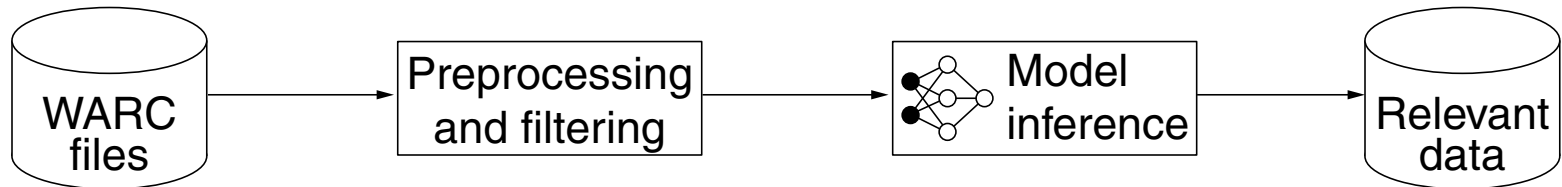
## Streamed Model Training Pipeline



- ❑ Given a mining task and a trained (classification) model, collect relevant data.  
Only a fraction of the records within the WARC files are relevant.
- ❑ Goal: Mining at web scale (billions of WARC files)

# Web Archive Processing

## Streamed Model Training Pipeline



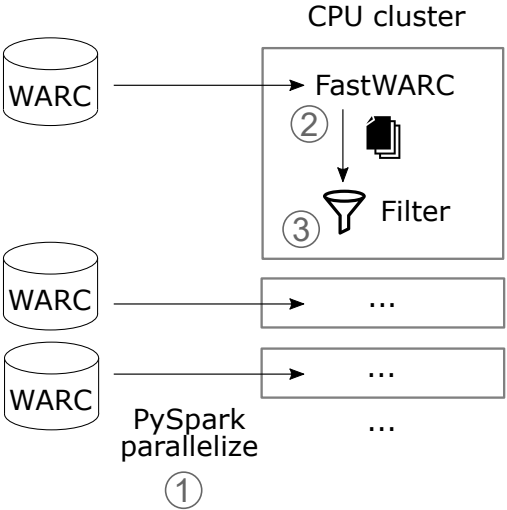
- ❑ Given a mining task and a trained (classification) model, collect relevant data. Only a fraction of the records within the WARC files are relevant.
- ❑ Goal: Mining at web scale (billions of WARC files)

### Observations:

- ❑ Mining / filtering WARC files is “embarrassingly parallel”.
- ❑ Decompressing WARC files, and processing WARC records are CPU bound.
- ❑ The mining / preprocessing step results in a variational data flow.
- ❑ Training of neural networks is GPU bound and presumes constant data flow.
- ❑ WARC storage, parallel processing, and GPU bound processing are on separate clusters.

# Web Archive Processing

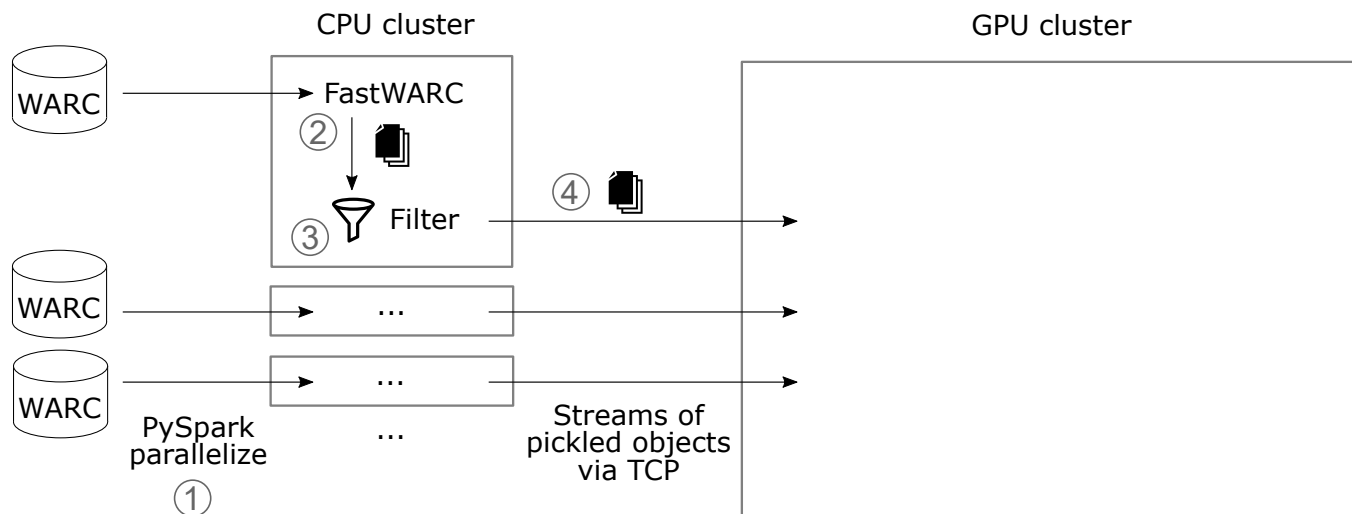
## Streamed Model Training Pipeline



1. PySpark distributes WARCs among workers
2. FastWARC decompresses and iterates records
3. First filtering step of records

# Web Archive Processing

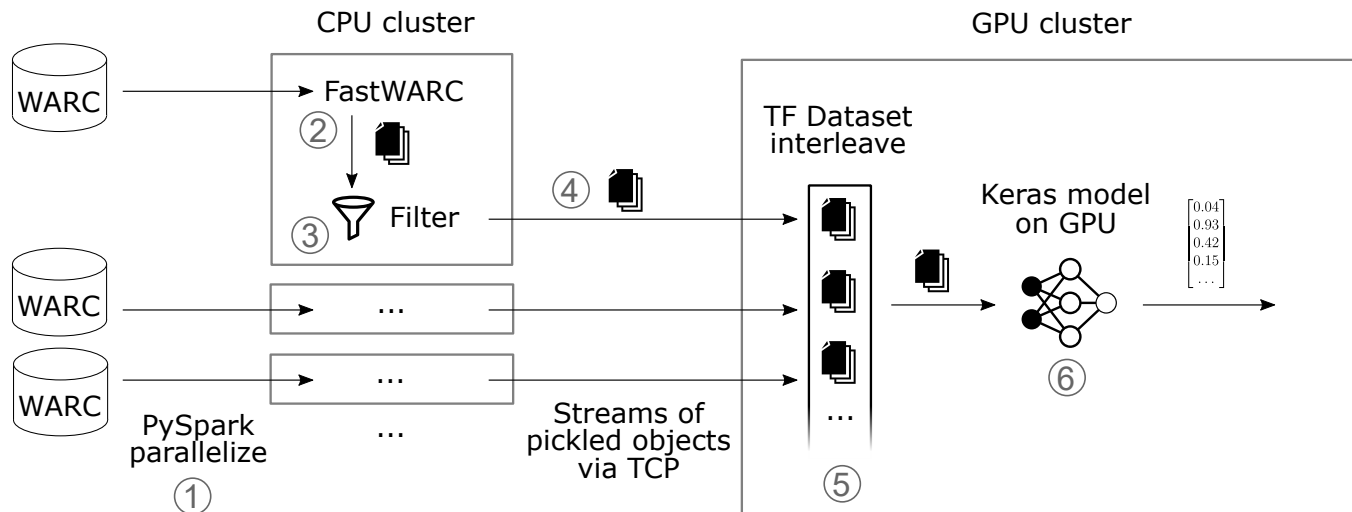
## Streamed Model Training Pipeline



1. PySpark distributes WARC files among workers
2. FastWARC decompresses and iterates records
3. First filtering step of records
4. Pickled record streams

# Web Archive Processing

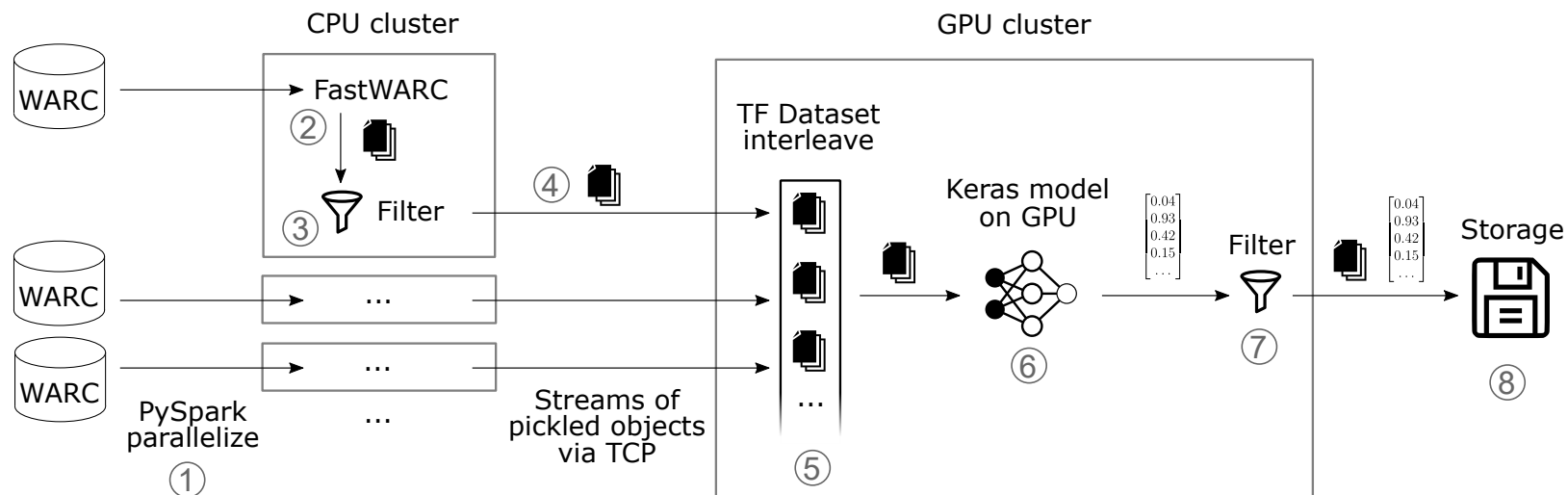
## Streamed Model Training Pipeline



1. PySpark distributes WARC files among workers
2. FastWARC decompresses and iterates records
3. First filtering step of records
4. Pickled record streams
5. Conversion to Tensorflow datasets and source interleaving
6. Batched processing by a Keras model

# Web Archive Processing

## Streamed Model Training Pipeline



1. PySpark distributes WARC files among workers
2. FastWARC decompresses and iterates records
3. First filtering step of records
4. Pickled record streams
5. Conversion to Tensorflow datasets and source interleaving
6. Batched processing by a Keras model
7. Second filtering based on classification results
8. Storage of relevant data



Webis Archive Research



Archival support

Argumentation

Language models

Search engines

Social sciences

Text reuse

Text synthesis

## ❑ Web Page Segmentation

Goal: Improve reliability of semantic web page segmentation.

## ❑ Web Crawling Quality Analysis

Goals: (1) Detect incomplete crawls.

(2) Improve the web page reconstructability from crawls.

## ❑ Personal Web Archival

Goal: Technology for individual web archive creation and search.

- ❑ **Learn Discussion Strategies**

Approach: Harvesting talk pages, email repositories, Reddit threads.

- ❑ **Acquire Justification and Reasoning Knowledge**

Approach: Construction of a causality graph from causal statements.

- ❑ **Compute Ranking Functions for Arguments**

Approach: Analysis of the hyperlink graph of web pages.

## □ Truths and Myths of the Mnemonic Password Advice

Approach: Construction of a position-dependent, higher-order language model, based on word initials of two billion sentences of verified casual language.

Example:

“The quick brown fox jumps over the lazy dog.”

~> Is “**Tqbfjotld**” a strong password?



**args.me**

The first (2017) search engine for arguments on the web.



**ChatNoir**

Search engine with rank explanation, indexing the ClueWeb and the CommonCrawl.



**Netspeak**

Phrase search engine for text correction and idiomatic writing.



**Picapica**

Search engine for text reuse detection.

- ❑ **Detect and Visualize Vandalism in Social Software**

Approach: Spatio-temporal analysis of reverted Wikipedia edits.

- ❑ **“Celebrity” Profiling**

Goal: Following personal traits on the Internet.

- ❑ **Hyperpartisan News Detection**

Goal: Analyzing political bias and illustrating provenance on the Internet.

Archival support

Argumentation

Language models

Search engines

Social sciences

Text reuse

Text synthesis

## ❑ Who Wrote the Web?

Applying author identification technology at web-scale.

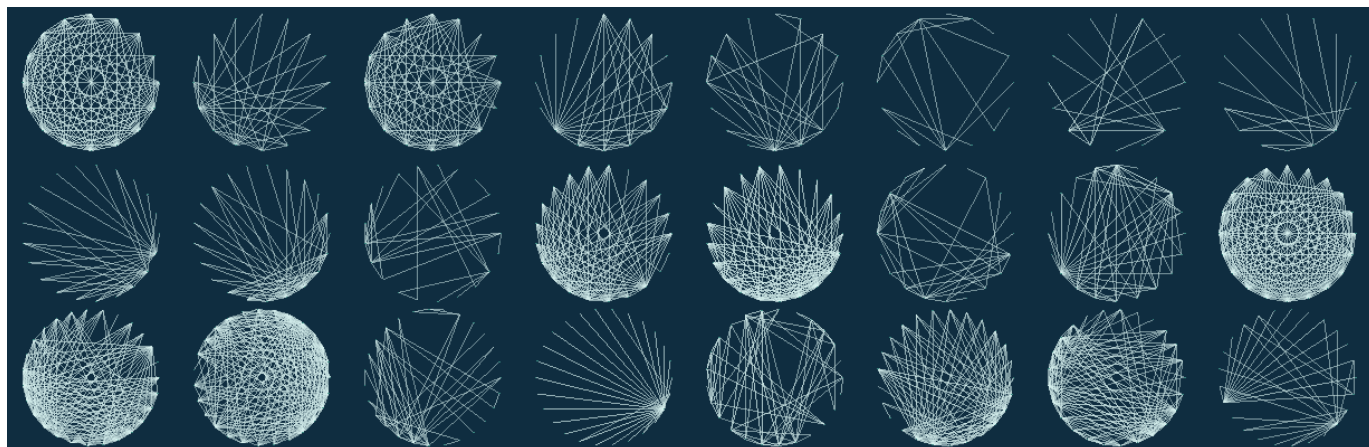
## ❑ Text Reuse Analytics

Goals: (1) Finding Wikipedia text reuse (on the web).

(2) Quantifying the prevalence of scientific text reuse.

## ❑ Text Reuse Illustration

Example: Visualizing article similarities in Wikipedia.



Riemann et al.:  
*Visualizing Article  
Similarities in  
Wikipedia.*  
EuroVis 2016

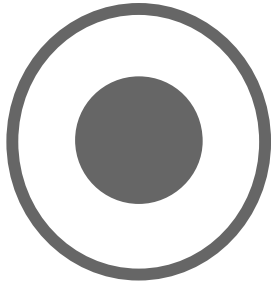
## ❑ Abstractive Snippet Generation

Approach: Use of anchor contexts to generate abstractive snippets with a pointer-generator network, exploiting ClueWeb09, ClueWeb12, and the DMOZ Open Directory Project.

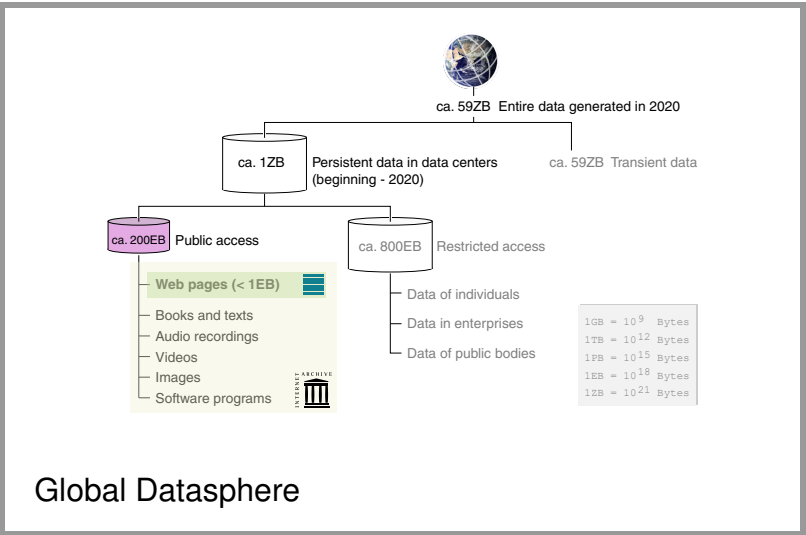
## ❑ Learn Automatic Summarization

Approach: Exploit author-provided summaries, taking advantage of the common practice of appending a “TL;DR” to long posts.



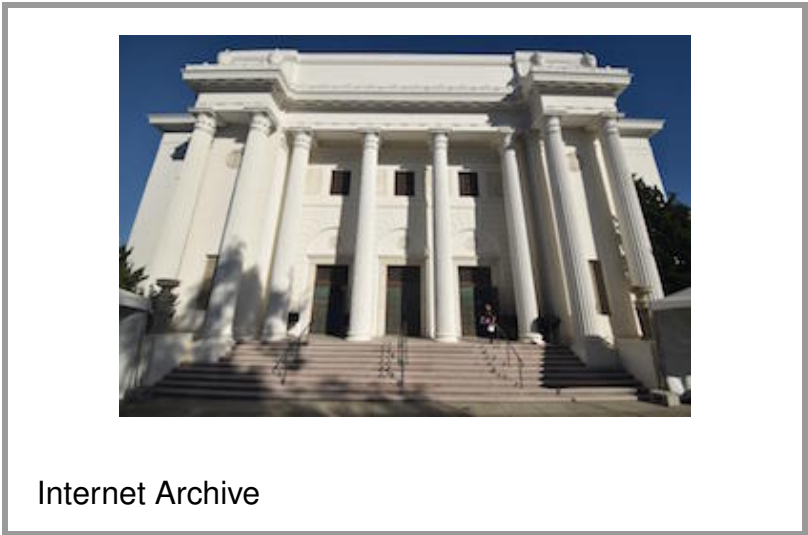
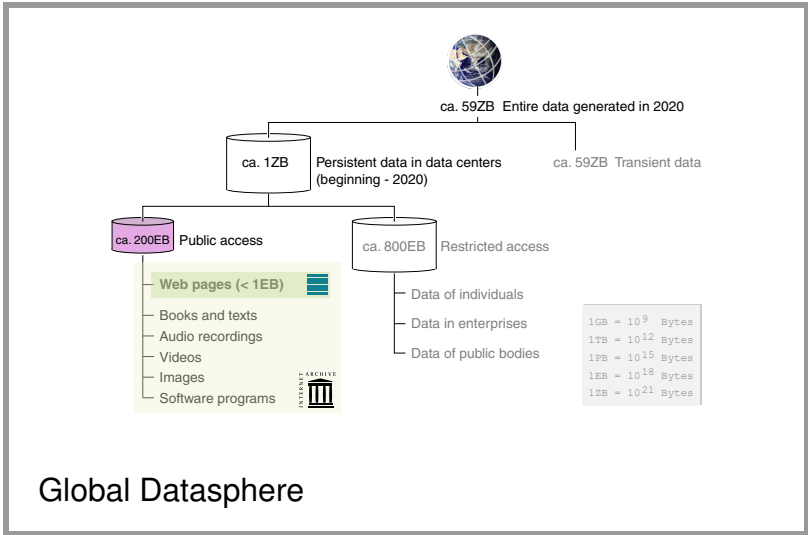


# Summary

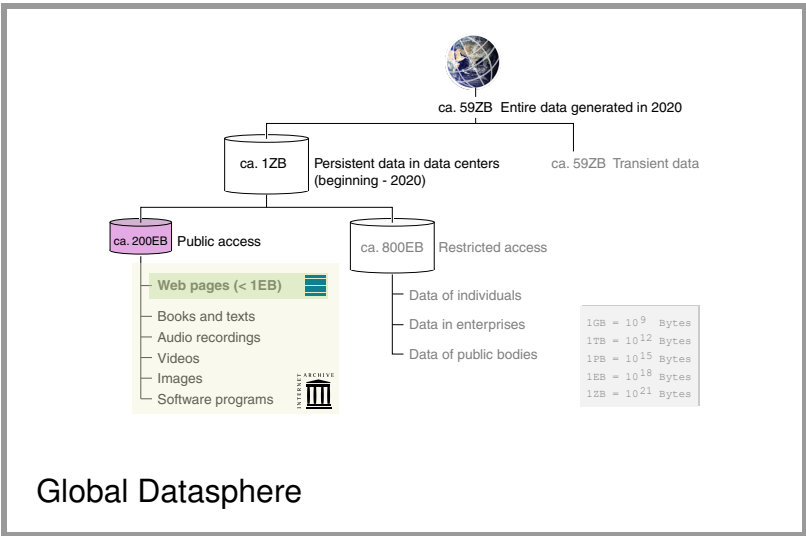


Global Datasphere

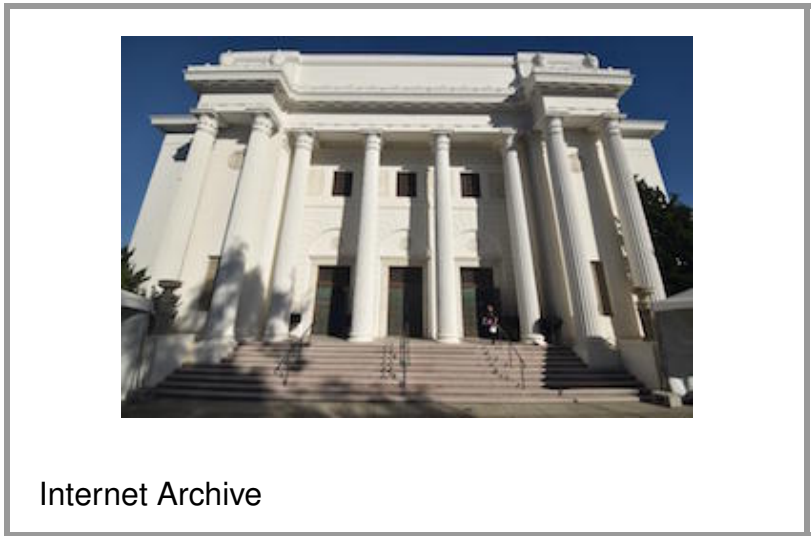
# Summary



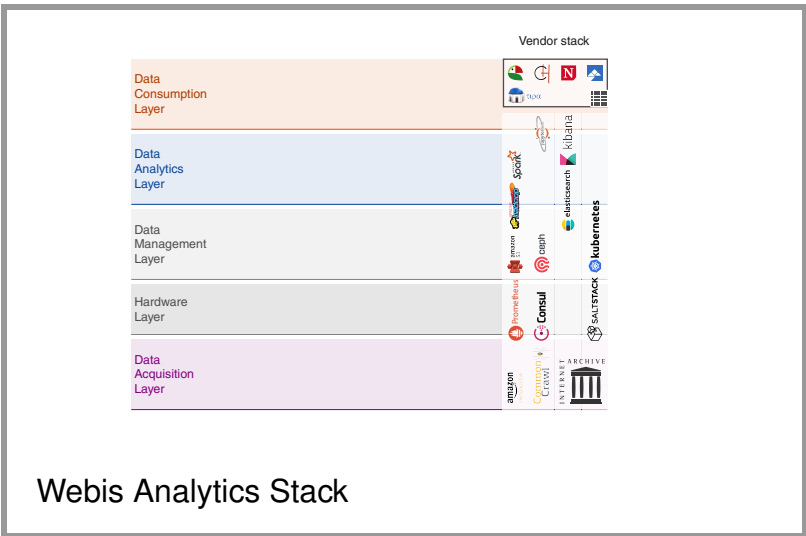
# Summary



Global Datasphere

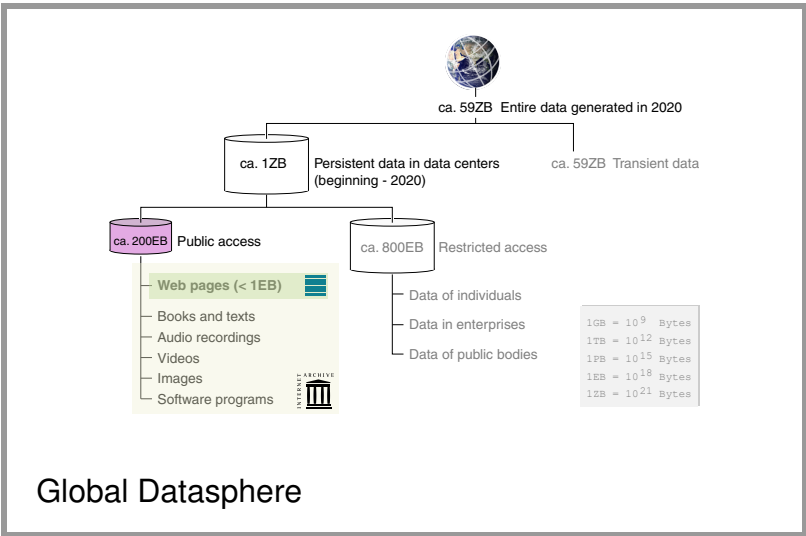


Internet Archive

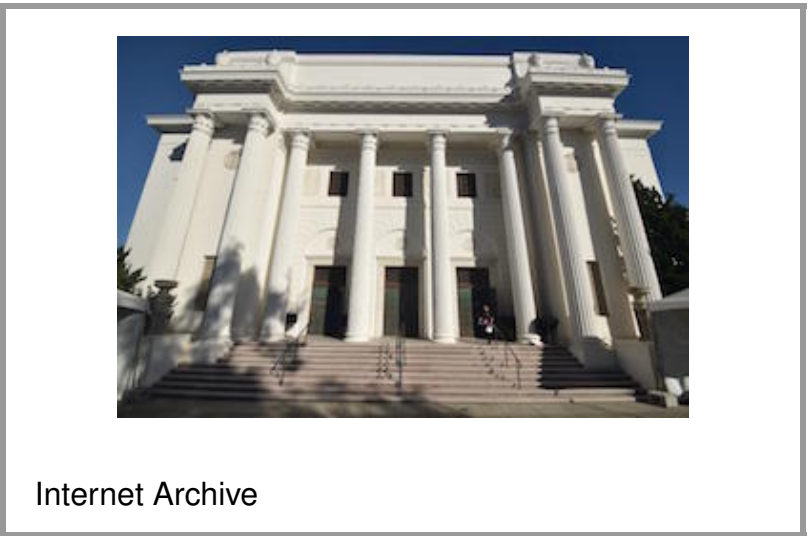


Webis Analytics Stack

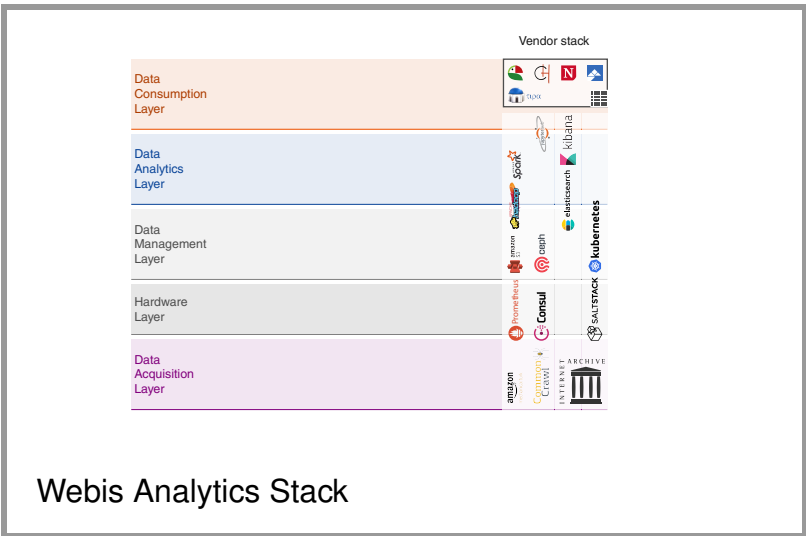
# Summary



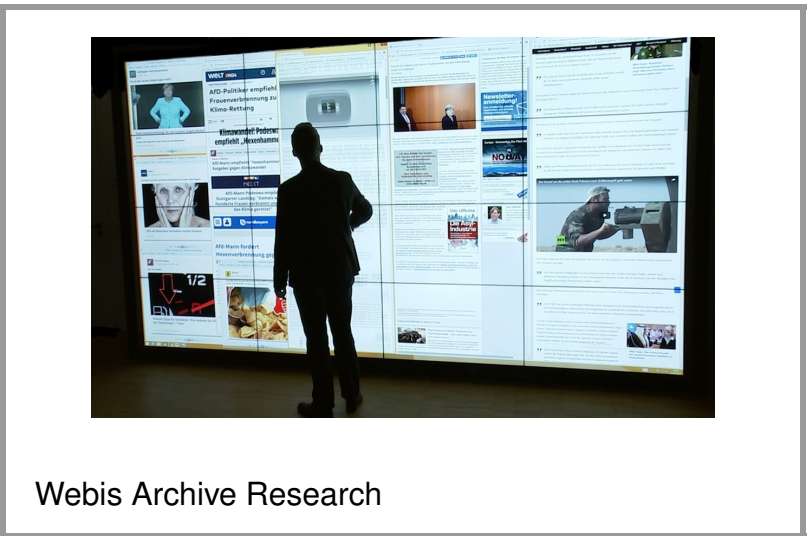
Global Datasphere



Internet Archive



Webis Analytics Stack



Webis Archive Research

Thank You!