

IR and NLP Research in the Webis Group

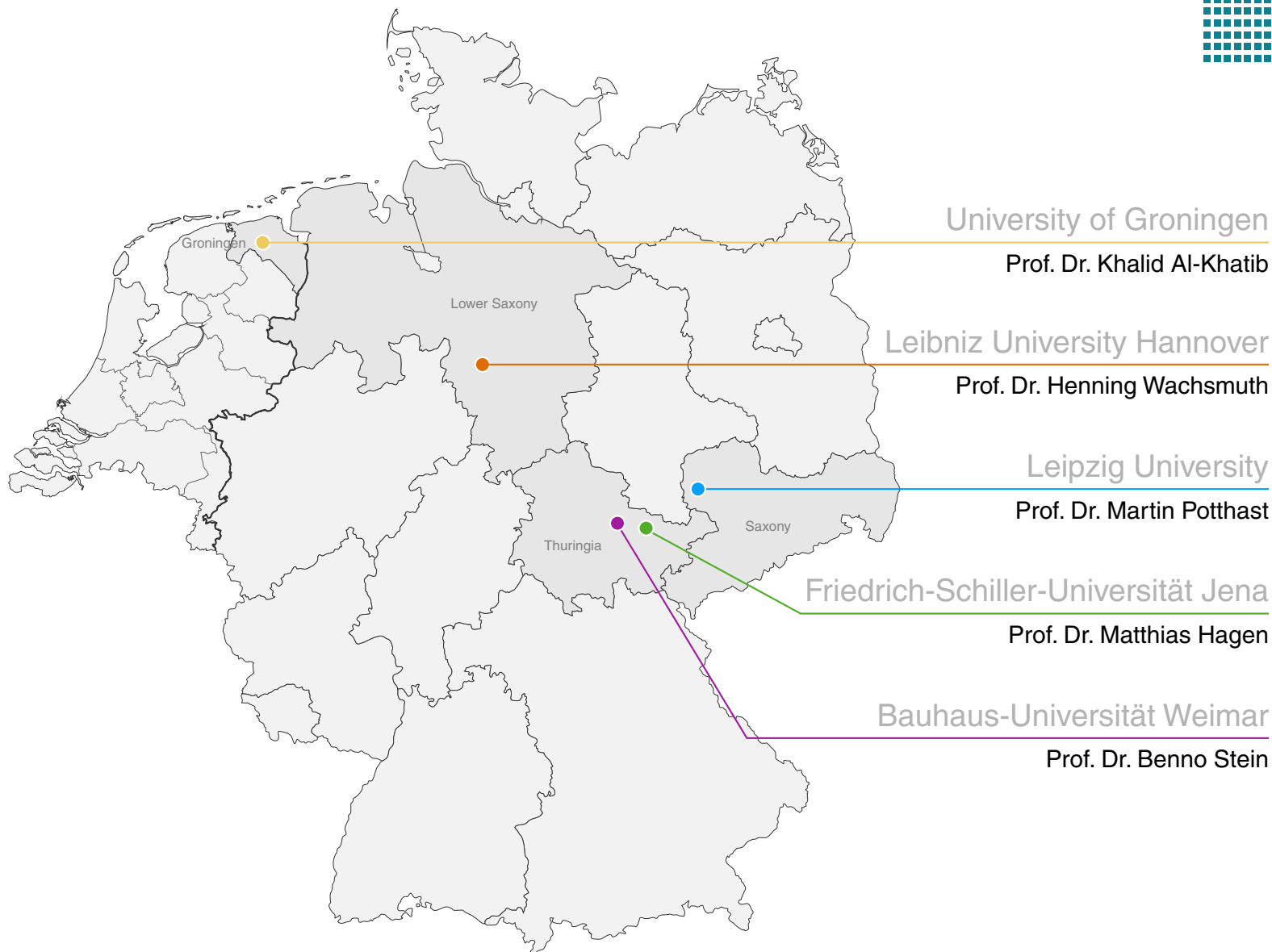
—Overview and Background—

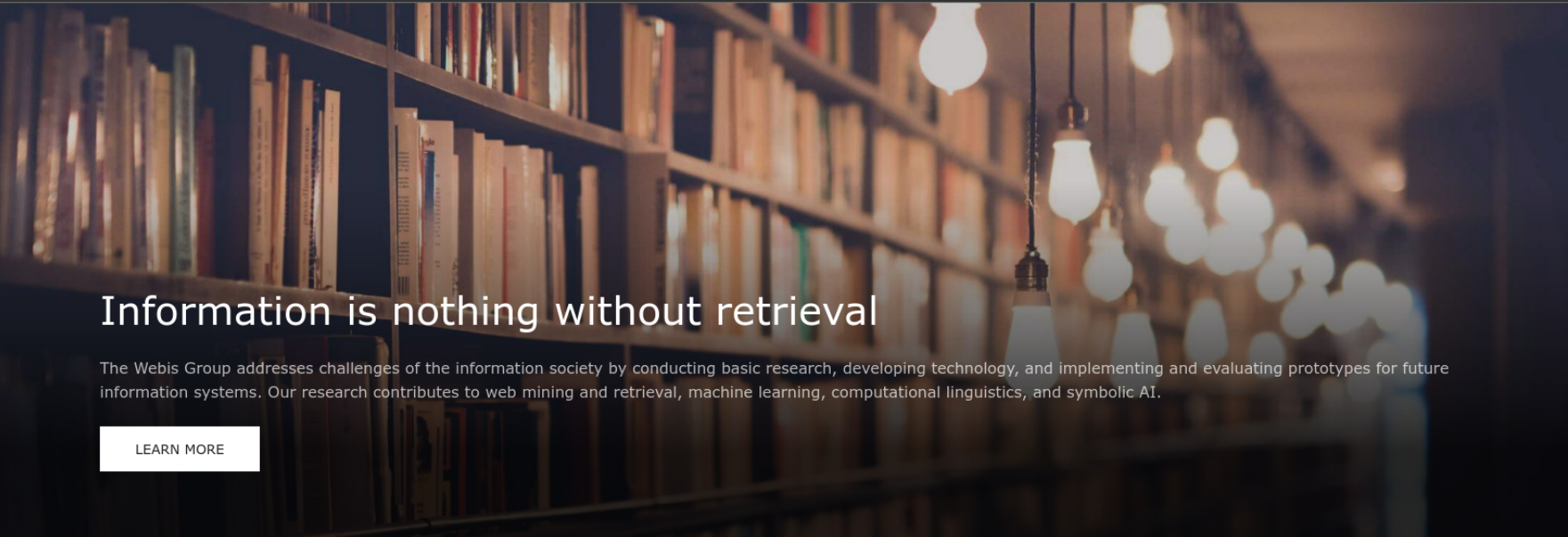
Benno Stein
Bauhaus-Universität Weimar
webis.de

Data Science Colloquium · Paderborn · December 16, 2022

Outline

- ① About us (Webis)
- ② Archive Data
- ③ Data Analytics @ Webis






Information is nothing without retrieval


The Webis Group addresses challenges of the information society by conducting basic research, developing technology, and implementing and evaluating prototypes for future information systems. Our research contributes to web mining and retrieval, machine learning, computational linguistics, and symbolic AI.

LEARN MORE

Search Engines

 **Argo**
Argument search

 **ChatNoir**
Web search

 **Netspeak**
Writing assistance

 **Picapica**
Plagiarism detection

	GRONINGEN	HANNOVER	JENA	LEIPZIG	WEIMAR
Home	Home	Home	Home	Home	Home
People	People	People	People	People	People
Teaching	Teaching	Teaching	Teaching	Teaching	Teaching
Research	Research	Research	Research	Research	Research



Archive Data: The Global Datasphere



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, ...”

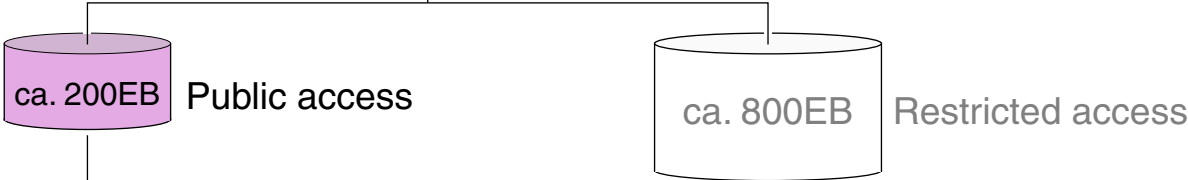
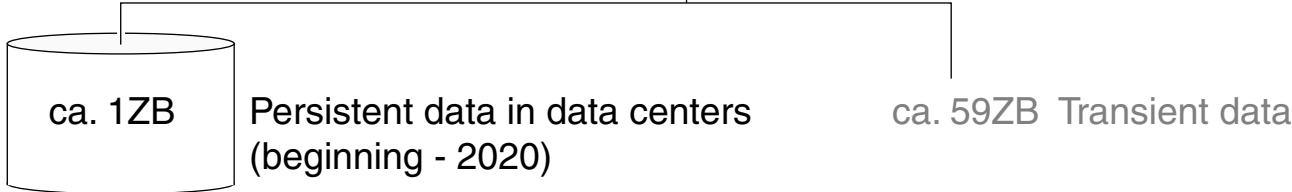
[IDC, 2012]

© 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

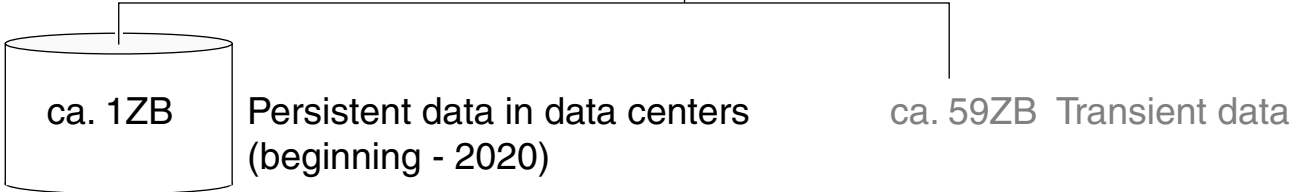
Data of public bodies

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



ca. 200EB Public access

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

ca. 800EB Restricted access

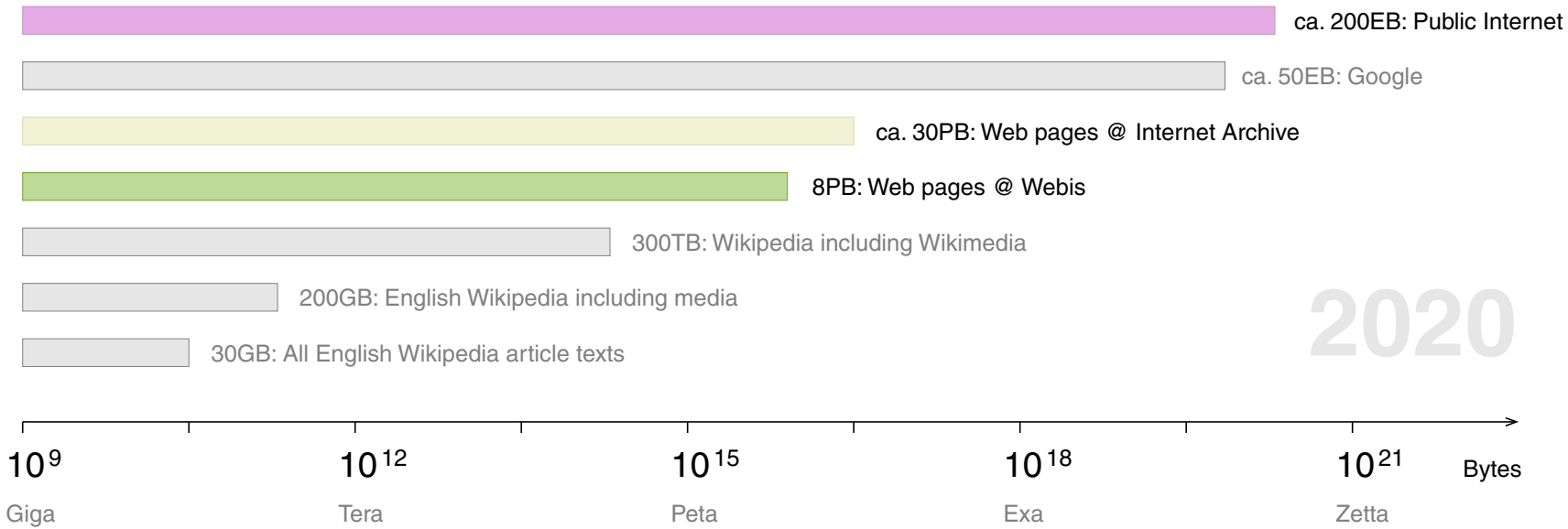
- Data of individuals
- Data in enterprises
- Data of public bodies

1GB	=	10 ⁹	Bytes
1TB	=	10 ¹²	Bytes
1PB	=	10 ¹⁵	Bytes
1EB	=	10 ¹⁸	Bytes
1ZB	=	10 ²¹	Bytes

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

The Global Datasphere in 2020

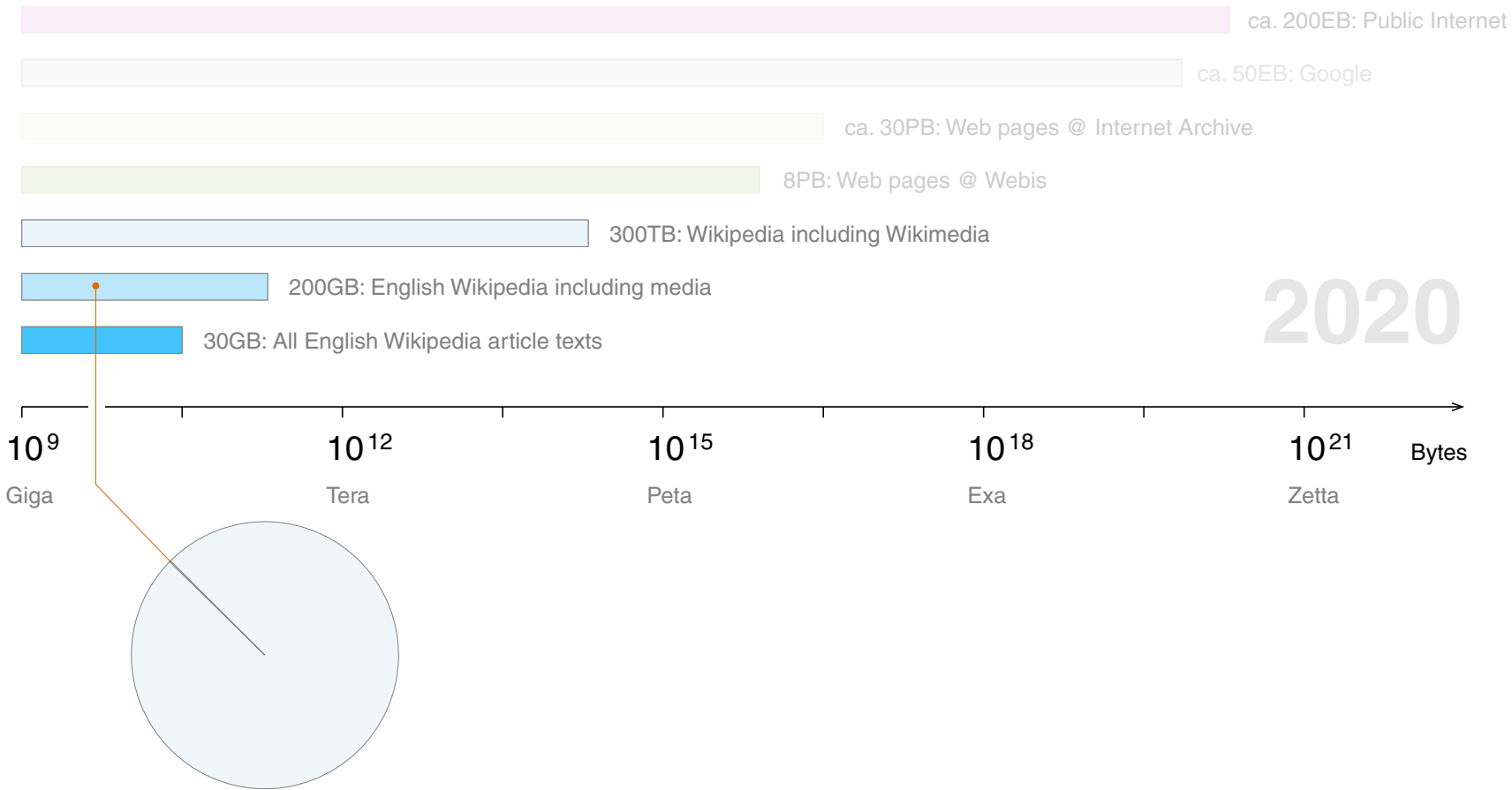
Relating Data Source Sizes [Vöslke et al. 2021]



2020

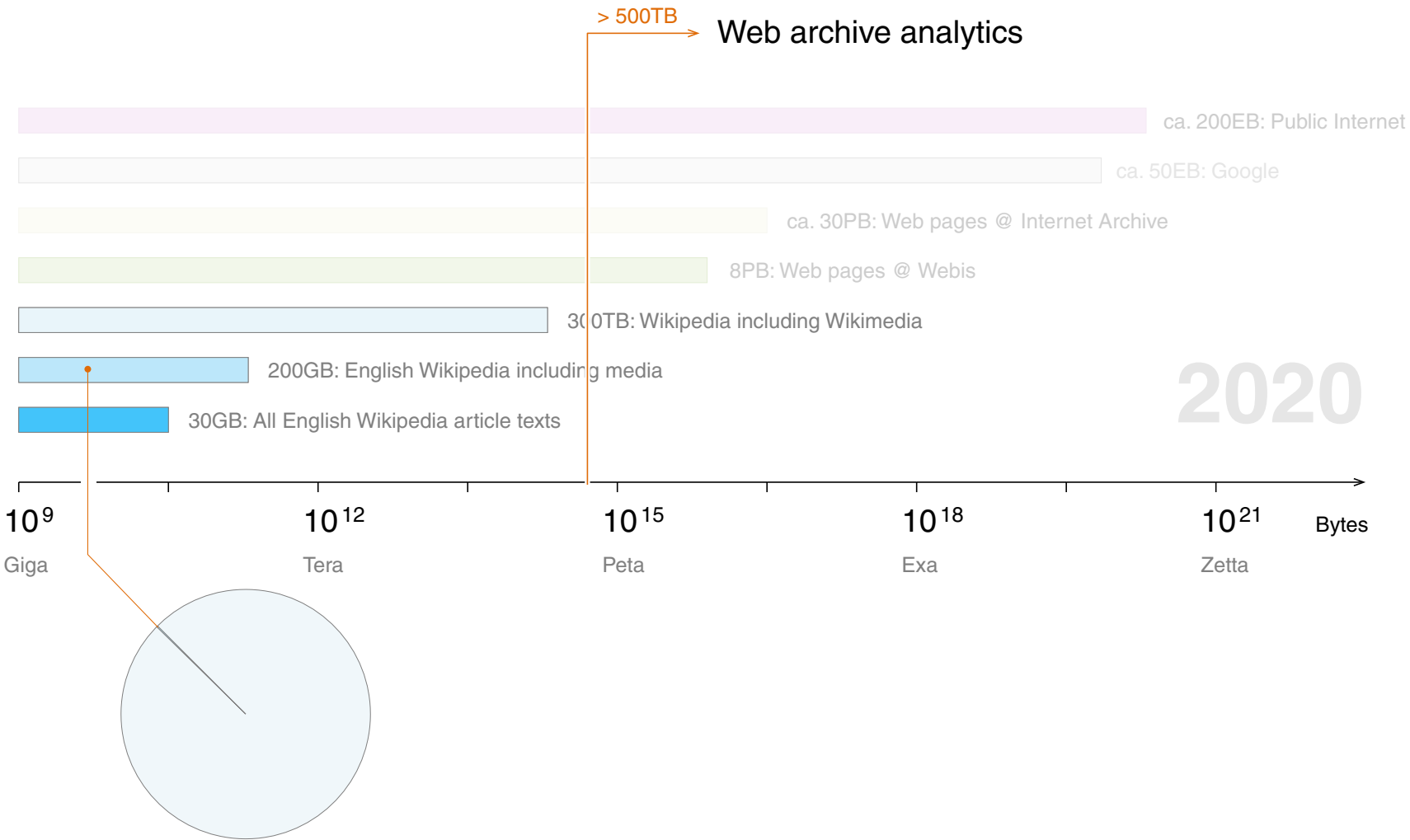
The Global Datasphere in 2020

Relating Data Source Sizes [Vöslke et al. 2021]



The Global Datasphere in 2020

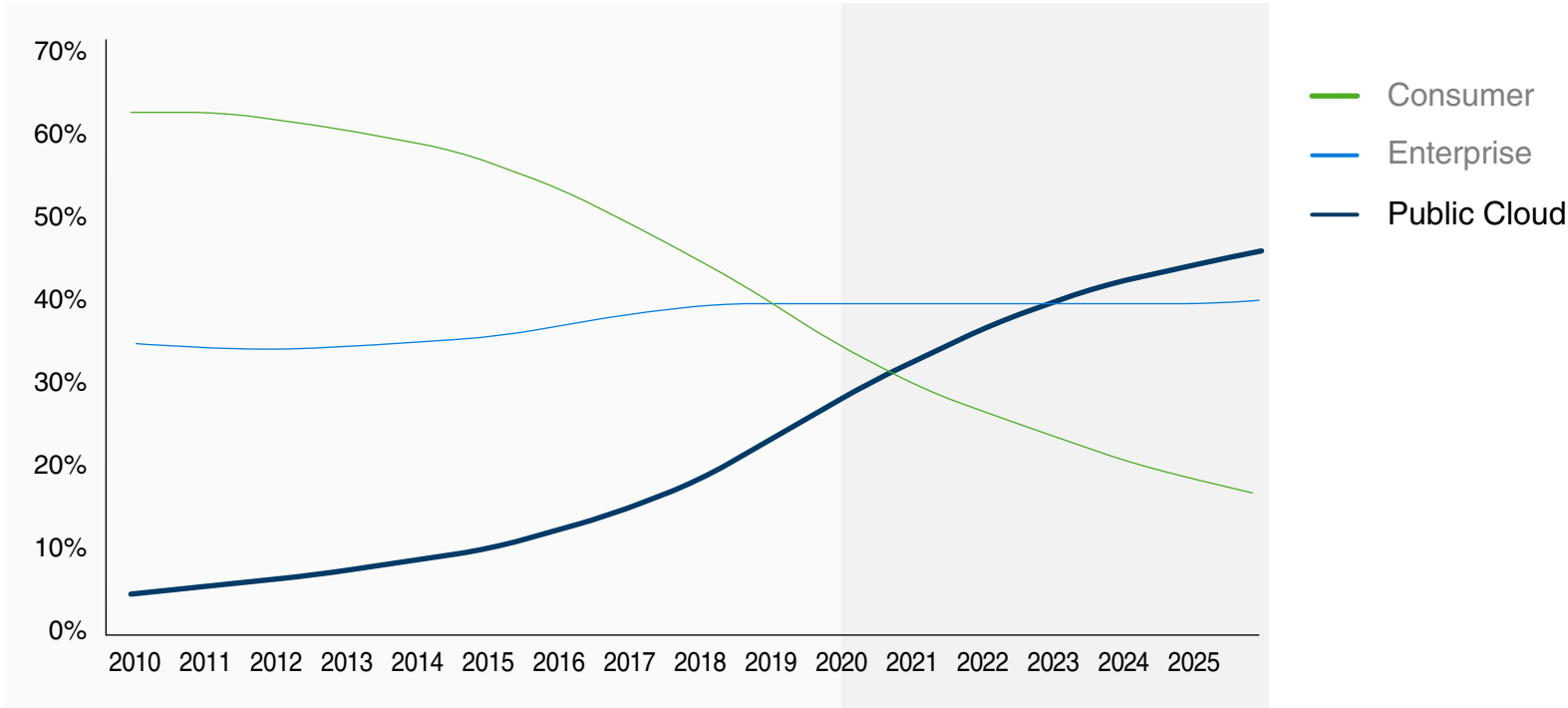
Relating Data Source Sizes [Vöslke et al. 2021]



2020

The Global Datasphere in 2020

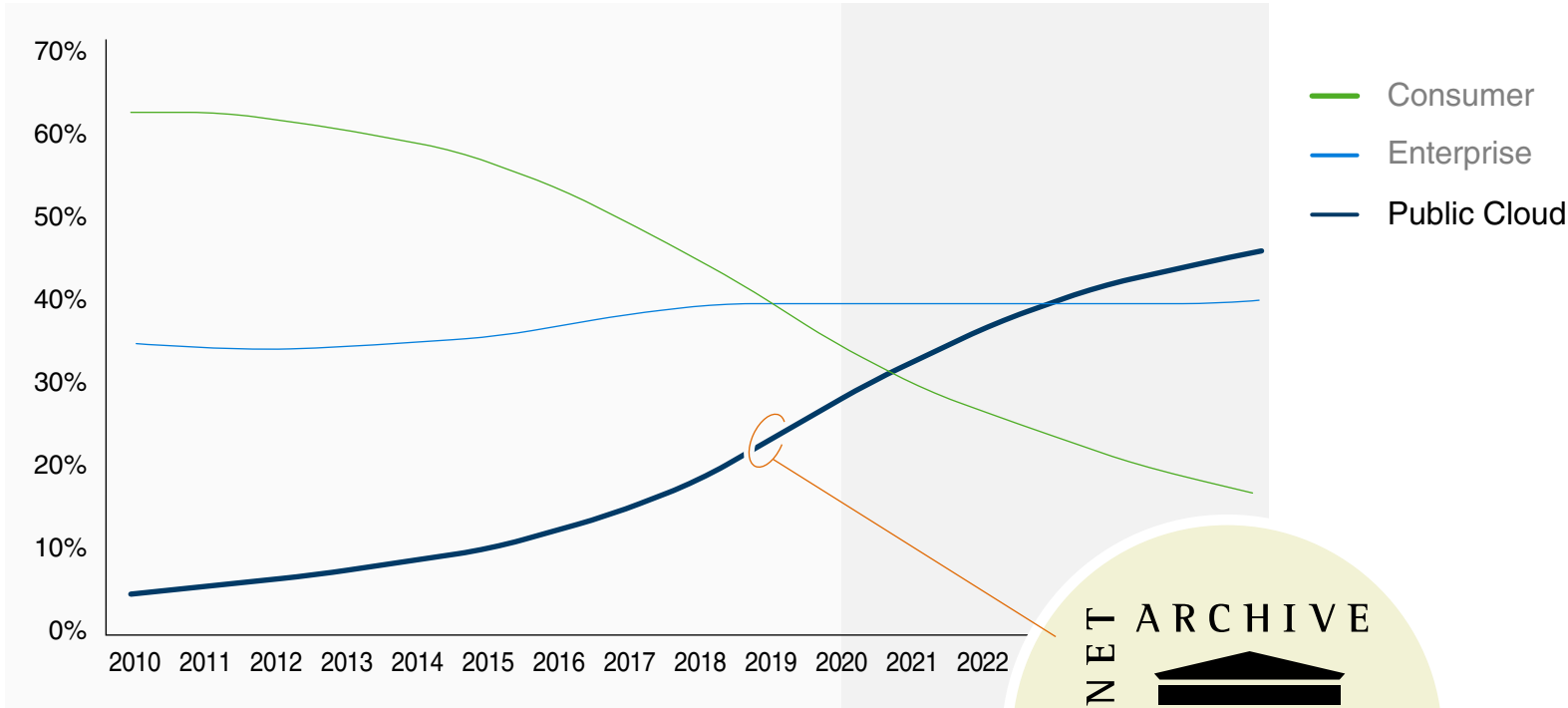
Where is the Data Stored?



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:



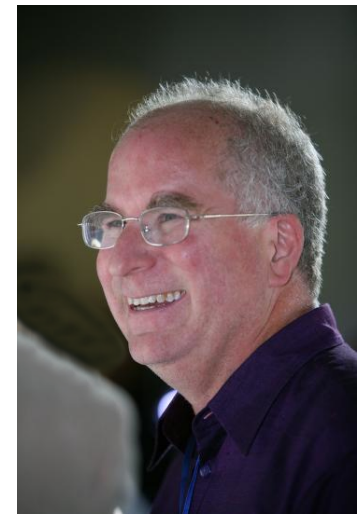
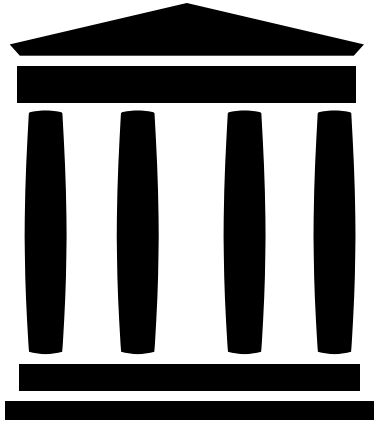
Basis: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, May 2020.




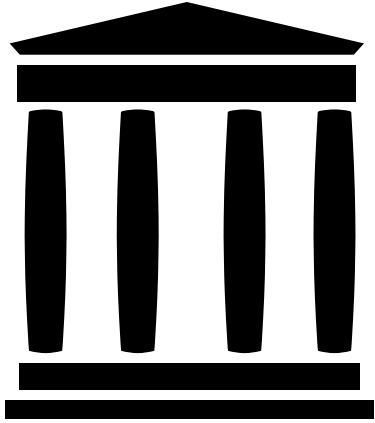
Archive Data: The Internet Archive



- ❑ 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]











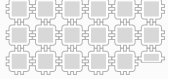
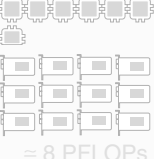
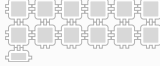


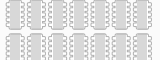





Data Analytics @ Webis: Platform and Stacks

Webis Data Center (Digital Bauhaus Lab)



Webis Data Center (Digital Bauhaus Lab)

	α -web [2009]	β -web [2015]	γ -web [2016 + 2021]	δ -web [2018]	ϵ -web [2020]
Nodes	44 	135 	9 	78 	55 
Disk [PB]	0.2 	4.1 	0.08 	12 	0.1 
Cores	176  $\cong 3.2$ TFLOPs	1,740  $\cong 67.4$ TFLOPs	672 + 227,328  $\cong 8$ PFLOPs	1,248  $\cong 119.8$ TFLOPs	1,100  $\cong 44$ TFLOPs
RAM [TB]	0.8 	28 	7.5 	10 	7 

Typical research:

α -Web. Teaching, Staging environment

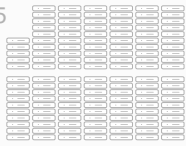









β -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

γ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

δ -Web. Web archiving, Virtualization (storage)

ϵ -Web. Search index construction, Argument search

Webis Data Center (Digital Bauhaus Lab)

	α -web [2009]	β -web [2015]	γ -web [2016 + 2021]	δ -web [2018]	ϵ -web [2020]
Nodes	44 	135 	9 	78 	55 
Disk [PB]	0.2 	4.1 	0.08 	12 	0.1 
Cores	176  $\cong 3.2$ TFLOPs	1,740  $\cong 67.4$ TFLOPs	672 + 227,328  $\cong 8$ PFLOPs	1,248  $\cong 119.8$ TFLOPs	1,100  $\cong 44$ TFLOPs
RAM [TB]	0.8 	28 	7.5 	10 	7 

Typical research:

α -Web. Teaching, Staging environment


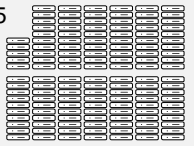

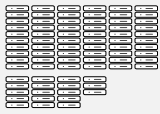







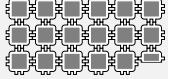
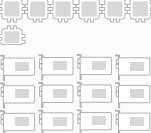



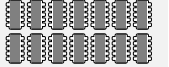



β -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

γ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

δ -Web. Web archiving, Virtualization (storage)

ϵ -Web. Search index construction, Argument search

Webis Data Center (Digital Bauhaus Lab)

	α -web [2009]	β -web [2015]	γ -web [2016 + 2021]	δ -web [2018]	ϵ -web [2020]
Nodes	44 	135 	9 	78 	55 
Disk [PB]	0.2 	4.1 	0.08 	12 	0.1 
Cores	176  ≈ 3.2 TFLOPs	1,740  ≈ 67.4 TFLOPs	672 + 227,328  ≈ 8 PFLOPs	1,248  ≈ 119.8 TFLOPs	1,100  ≈ 44 TFLOPs
RAM [TB]	0.8 	28 	7.5 	10 	7 

Typical research:

α -Web. Teaching, Staging environment


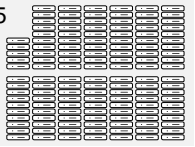

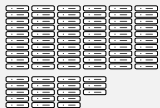







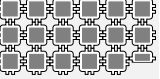
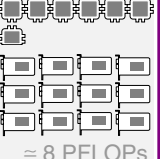



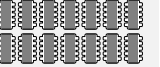



β -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

γ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

δ -Web. Web archiving, Virtualization (storage)

ϵ -Web. Search index construction, Argument search

Webis Data Center (Digital Bauhaus Lab)

	α -web [2009]	β -web [2015]	γ -web [2016 + 2021]	δ -web [2018]	ϵ -web [2020]
Nodes	44 	135 	9 	78 	55 
Disk [PB]	0.2 	4.1 	0.08 	12 	0.1 
Cores	176  $\cong 3.2$ TFLOPs	1,740  $\cong 67.4$ TFLOPs	672 + 227,328  $\cong 8$ PFLOPs	1,248  $\cong 119.8$ TFLOPs	1,100  $\cong 44$ TFLOPs
RAM [TB]	0.8 	28 	7.5 	10 	7 

Typical research:

α -Web. Teaching, Staging environment

β -Web. Web mining (map reduce), Authorship analytics, Virtualization (compute, web services)

γ -Web. Machine learning (embedding, deep learning), Text synthesis, Language modeling

δ -Web. Web archiving, Virtualization (storage)

ϵ -Web. Search index construction, Argument search

Webis Analytics Stack

Data
Consumption
Layer

Data
Analytics
Layer

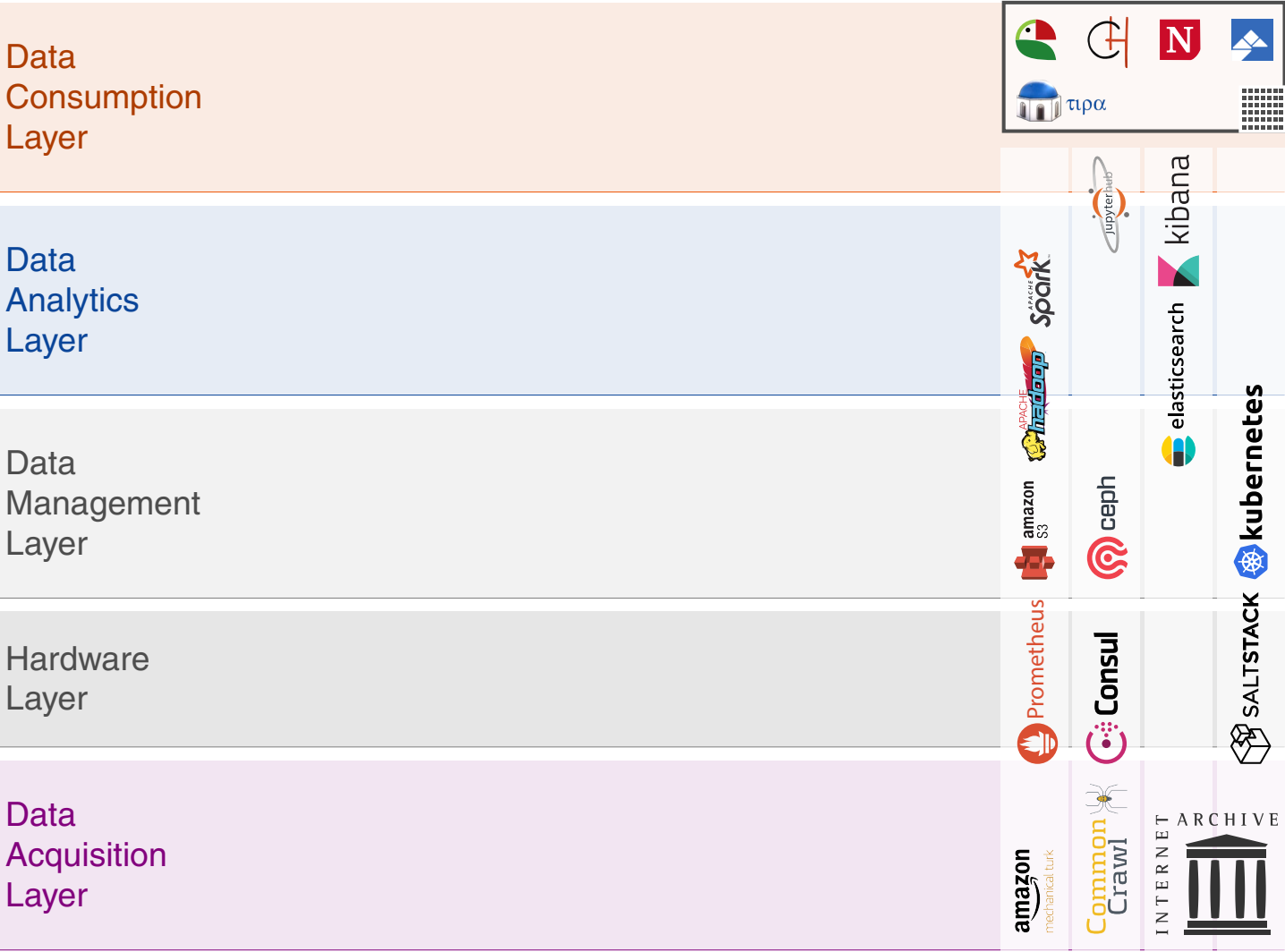
Data
Management
Layer

Hardware
Layer

Data
Acquisition
Layer

Webis Analytics Stack

Vendor stack



Webis Analytics Stack

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

Webis Analytics Stack

Task Stack

Technology stack

Vendor stack

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

Webis Analytics Stack

Task Stack

Technology stack

Vendor stack

Roles

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



Data Analytics @ Webis: Research

Search Engines
Language Models
Social Sciences
Argumentation
Text Reuse
Text Synthesis
Archival Support



args.me

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

Recent extension: Search for argumentative images.

Background: Argument ranking and search.



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.

□ 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.

Background: The BSI password creation advice.

Example:

The quick brown fox jumps over the lazy dog.”

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

❑ Detect and Visualize Vandalism in Social Software

Approach: Spatio-temporal analysis of reverted Wikipedia edits.

Service: Data analytics and visualization.

Background: Vandalism in Wikipedia.

❑ “Celebrity” Profiling

Goal: Following personal traits on the Internet.

❑ Hyperpartisan News Detection

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

- ❑ **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.

❑ Who Wrote the Web?

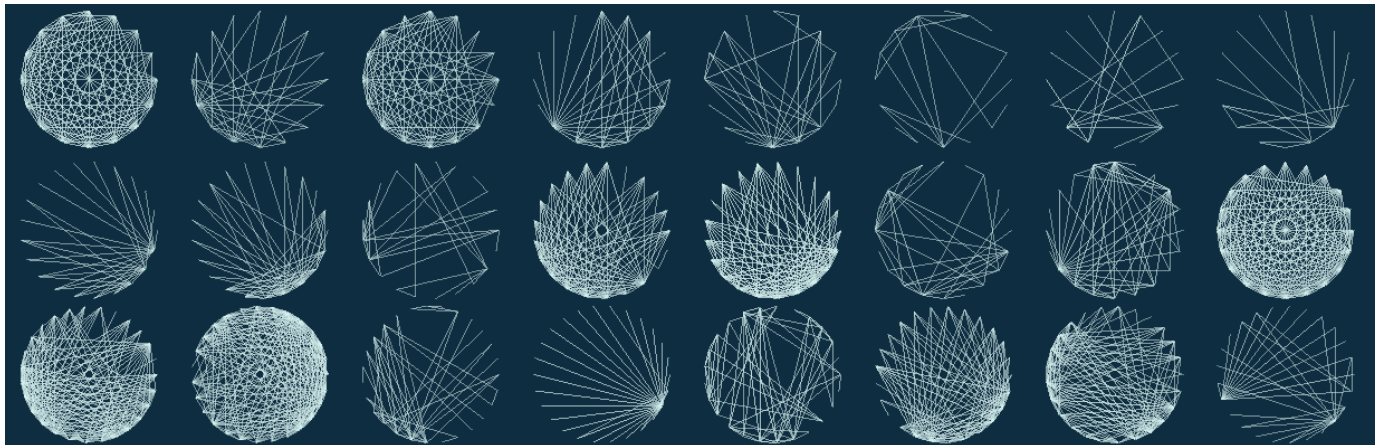
Applying author identification technology at web-scale.

❑ 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.

❑ Web Page Segmentation

Goal: Improve reliability of semantic web page segmentation.

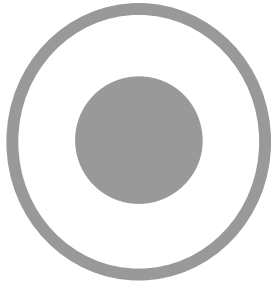
❑ Web Crawling Quality Analysis

Goals: (1) Detect incomplete crawls.

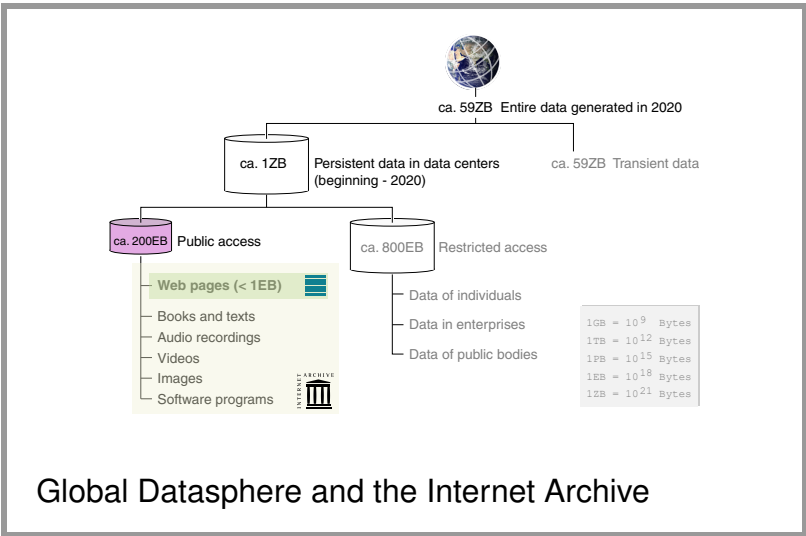
(2) Improve the web page reconstructability from crawls.

❑ Personal Web Archival

Goal: Technology for individual web archive creation and search.

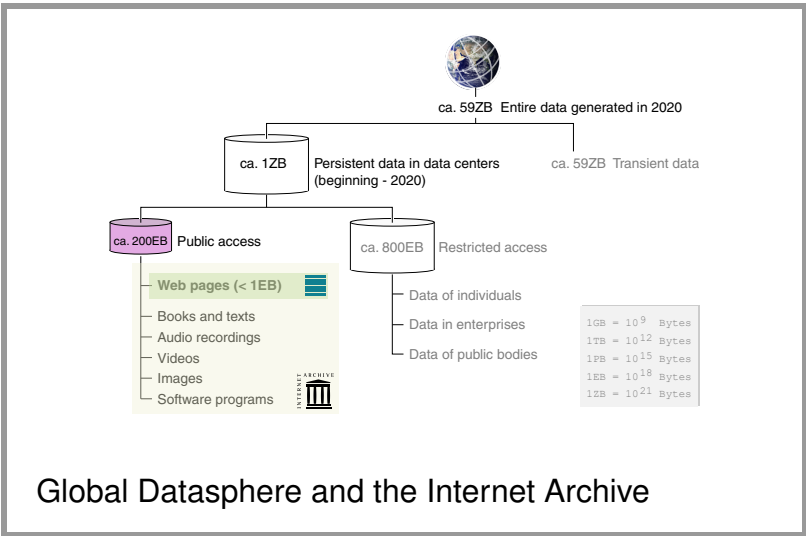


Summary

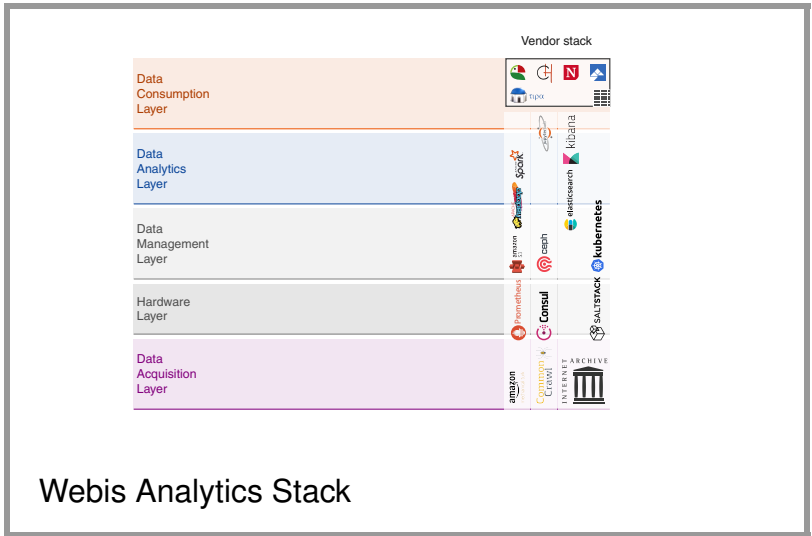


Global Datasphere and the Internet Archive

Summary

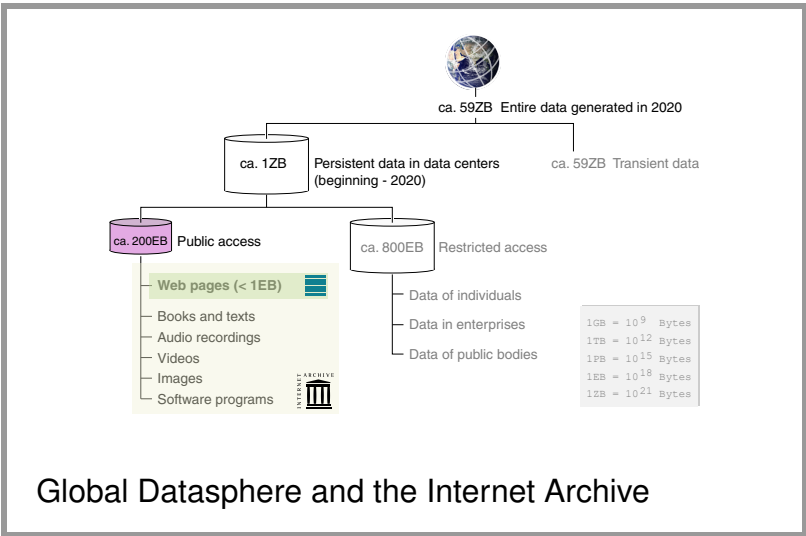


Global Datasphere and the Internet Archive

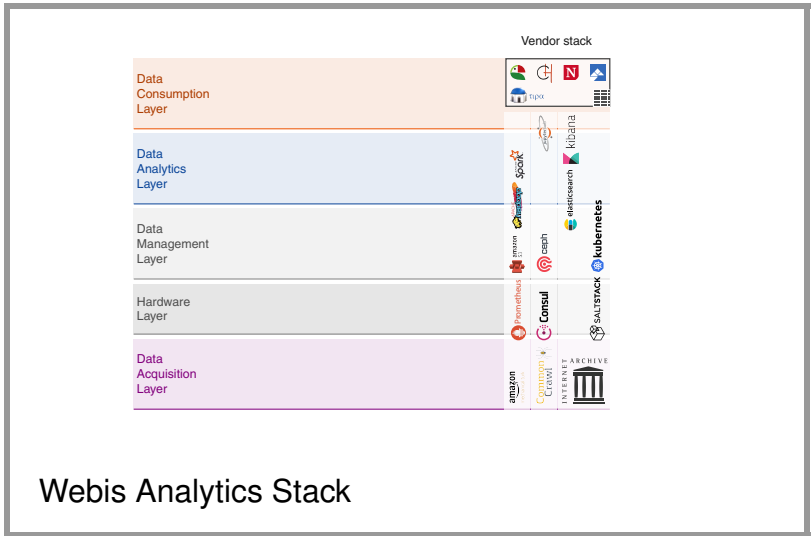


Webis Analytics Stack

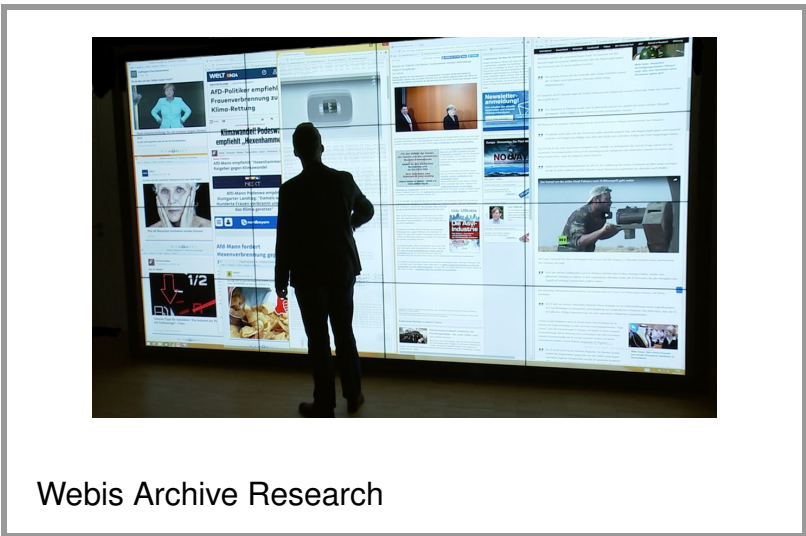
Summary



Global Datasphere and the Internet Archive

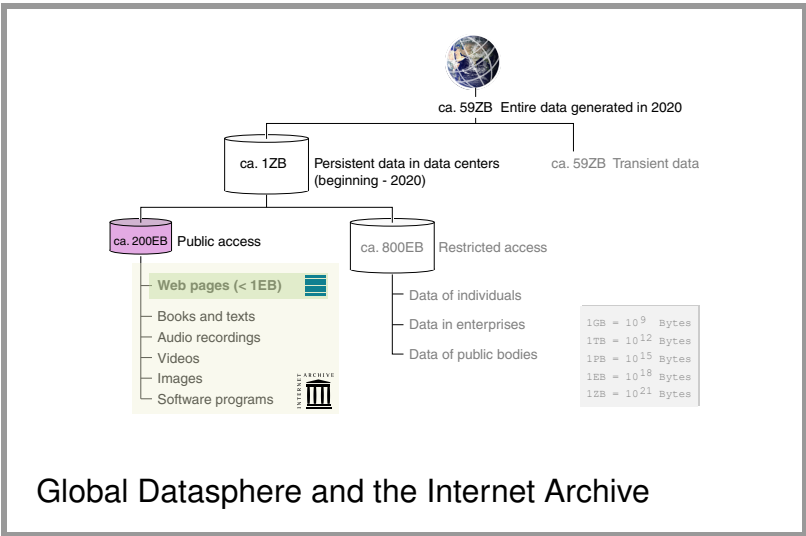


Webis Analytics Stack

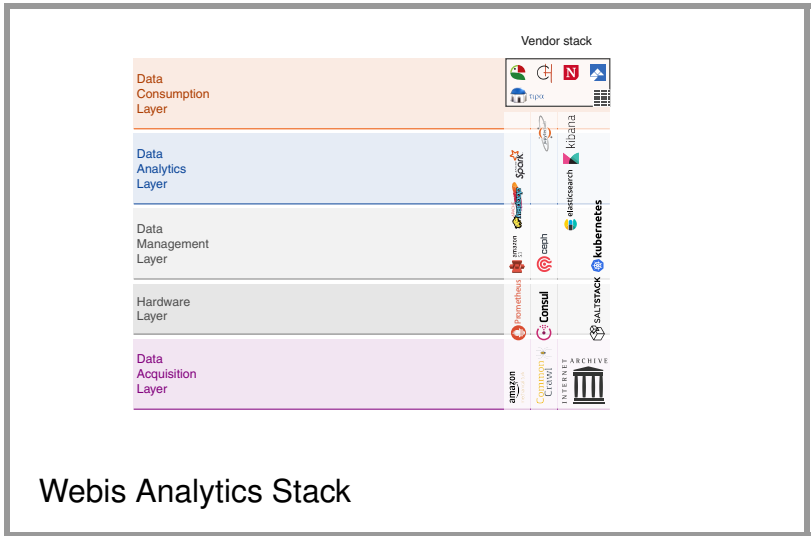


Webis Archive Research

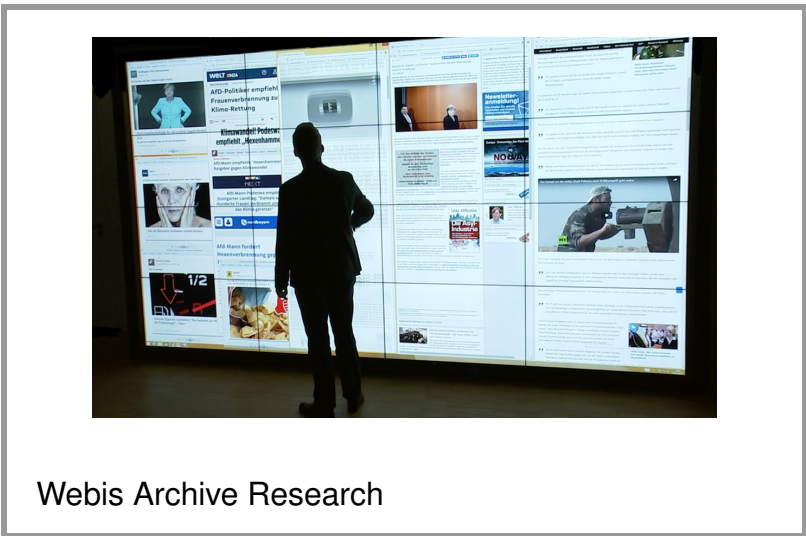
Summary



Global Datasphere and the Internet Archive



Webis Analytics Stack



Webis Archive Research



Webis Events

Thank You!