

Web Archive Analytics

Benno Stein
Bauhaus-Universität Weimar
webis.de

INFORMATIK2020 · Karlsruhe · September 30, 2020

Outline

- ① The Global Datasphere
- ② The Internet Archive
- ③ Web Archive Analytics @ Webis
- ④ 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]



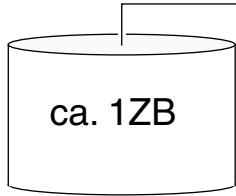
“... images and videos on mobile phones uploaded to YouTube, digital movies populating the pixels of our high-definition TVs, security footage at airports and major events such as the Olympic Games, subatomic collisions recorded by the Large Hadron Collider at CERN, banking data swiped in an ATM, transponders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

[IDC, 2012]

The Global Datasphere in 2020

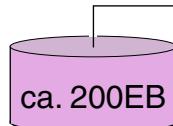


ca. 59ZB Entire data generated in 2020

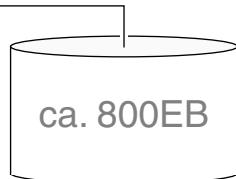


Persistent data in data centers
(beginning - 2020)

ca. 59ZB Transient data



Public access



Restricted access

Web pages (< 1EB)

- Books and texts
- Audio recordings
- Videos
- Images
- Software programs

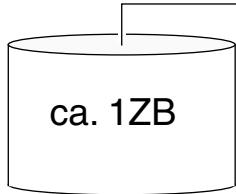
- Data of individuals
- Data in enterprises
- Data of public bodies

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

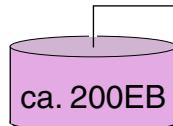


ca. 59ZB Entire data generated in 2020

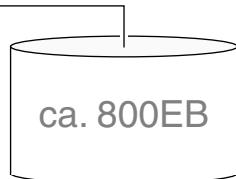


Persistent data in data centers
(beginning - 2020)

ca. 59ZB Transient data



Public access



Restricted access

Web pages (< 1EB)



- Books and texts
- Audio recordings
- Videos
- Images
- Software programs

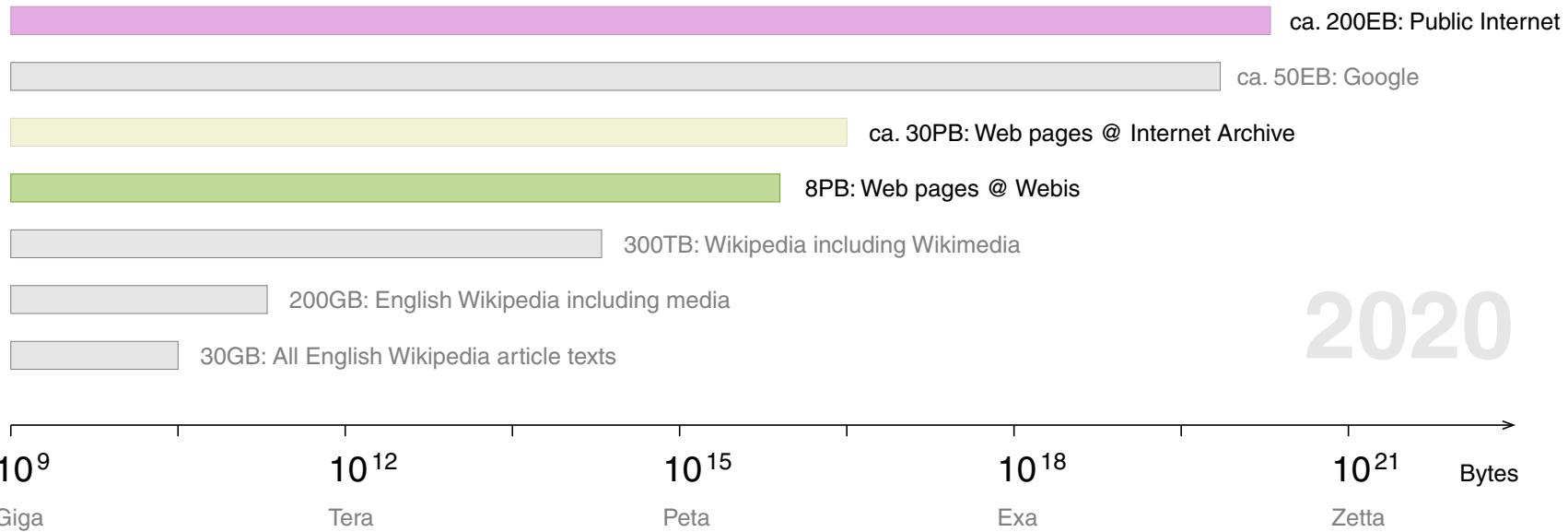


- Data of individuals
- Data in enterprises
- Data of public bodies

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

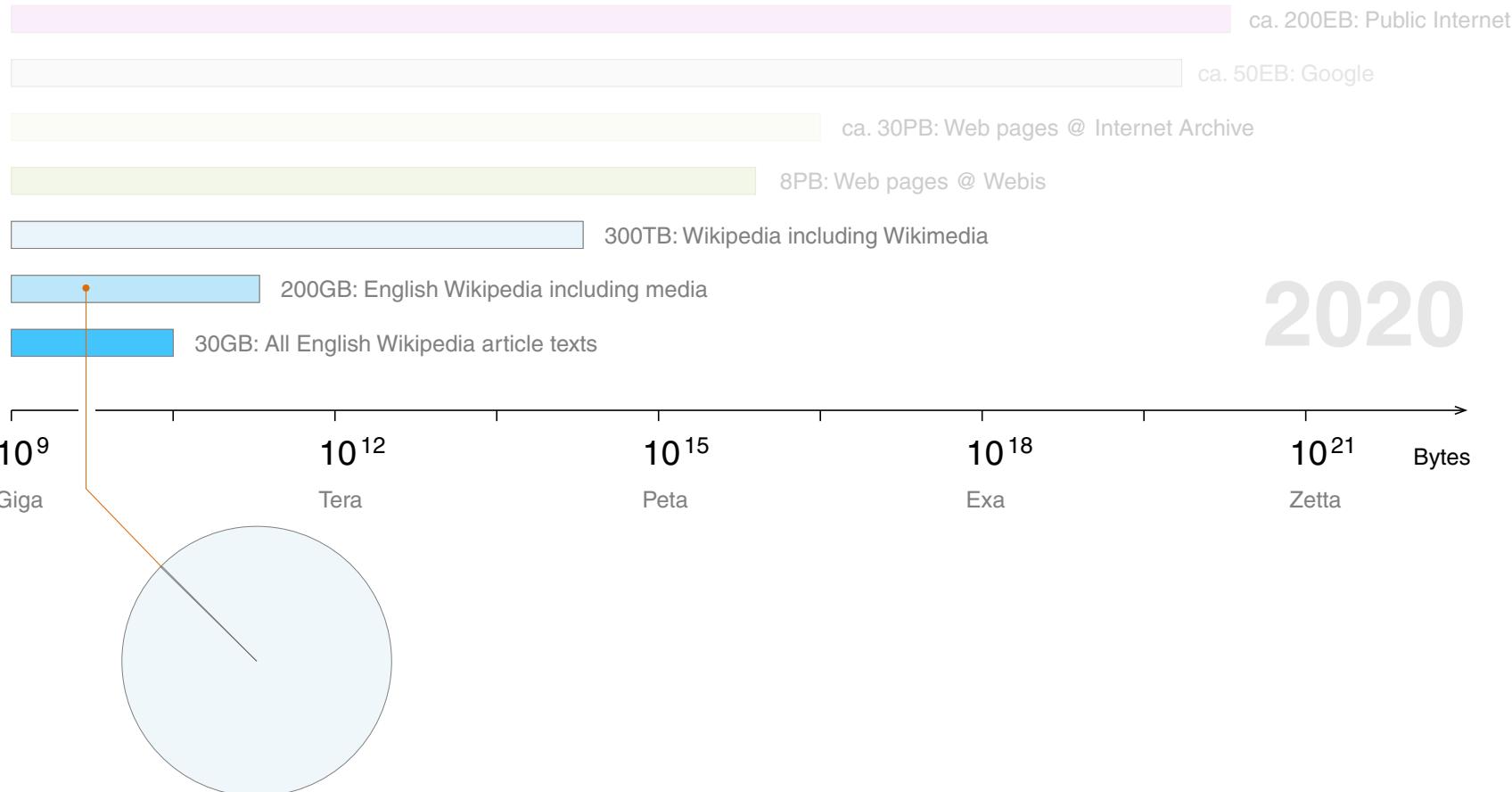
The Global Datasphere in 2020

Relating Data Source Sizes



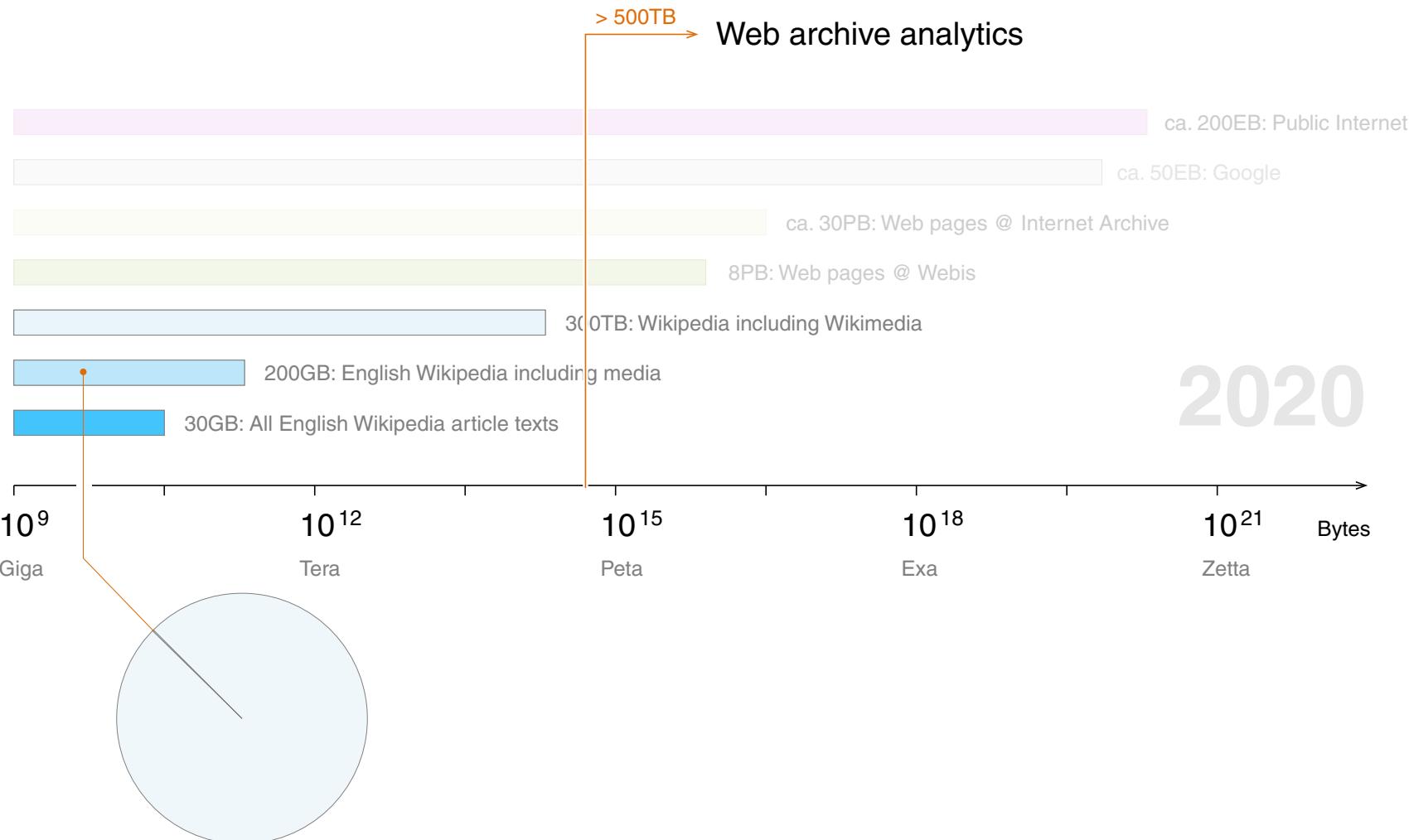
The Global Datasphere in 2020

Relating Data Source Sizes



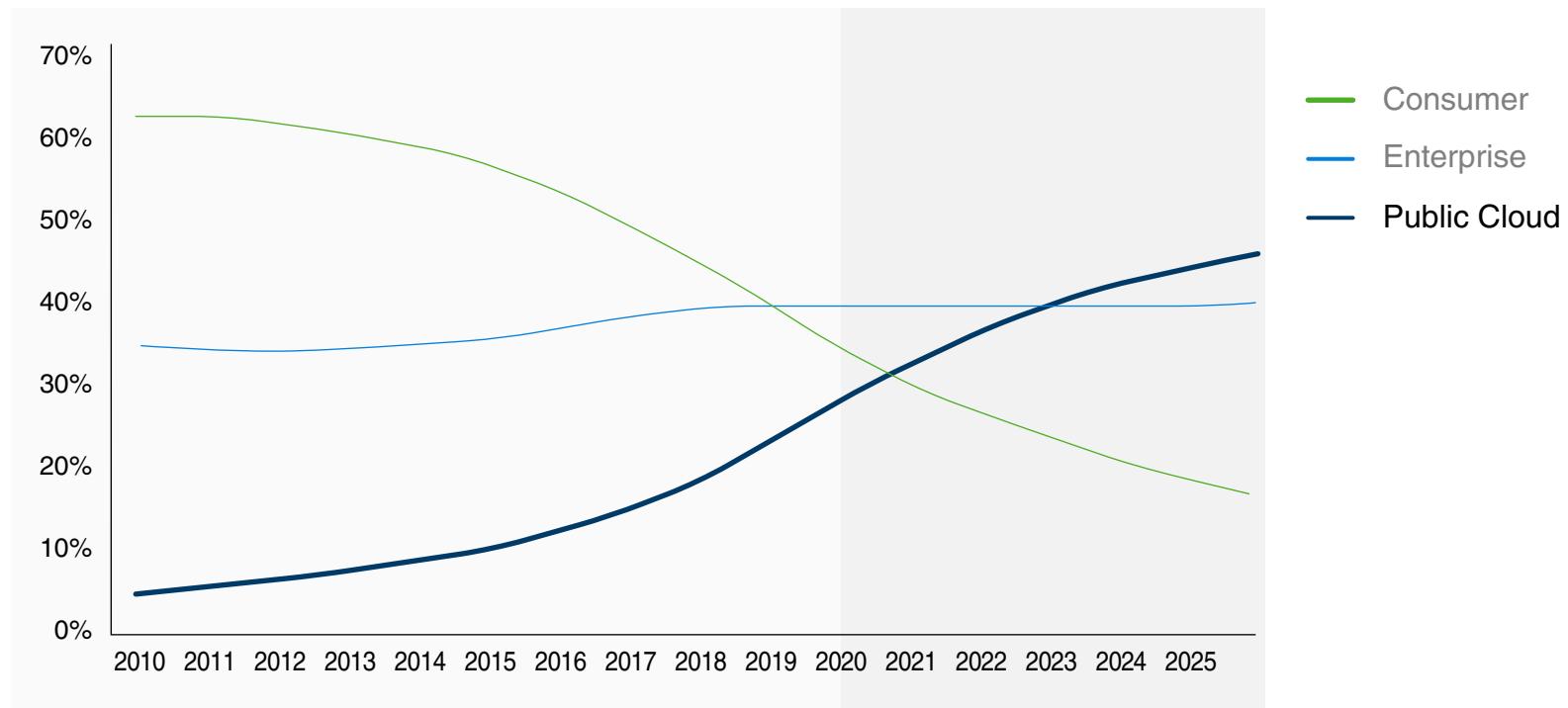
The Global Datasphere in 2020

Relating Data Source Sizes



The Global Datasphere in 2020

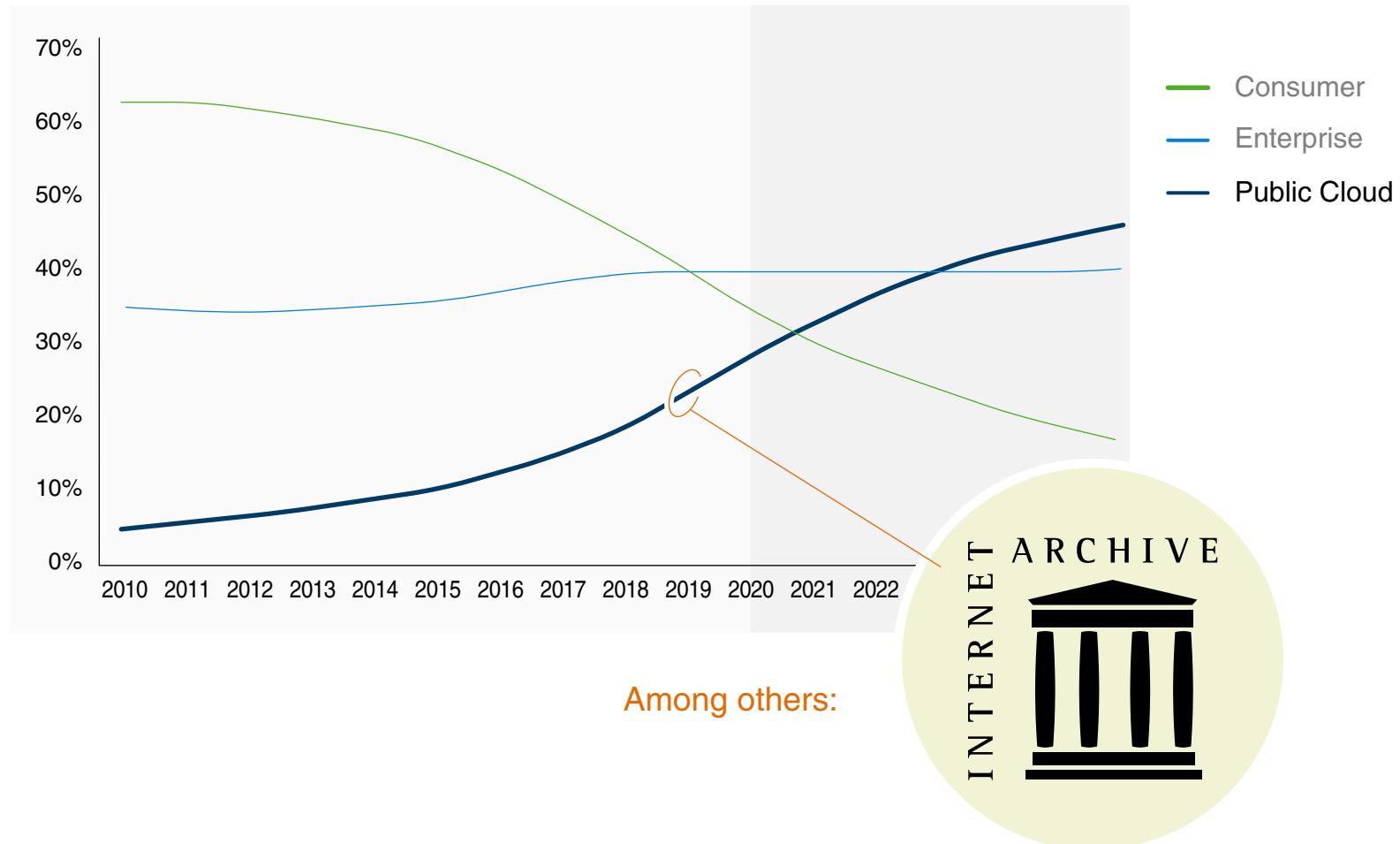
Where is the Data Stored?



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?

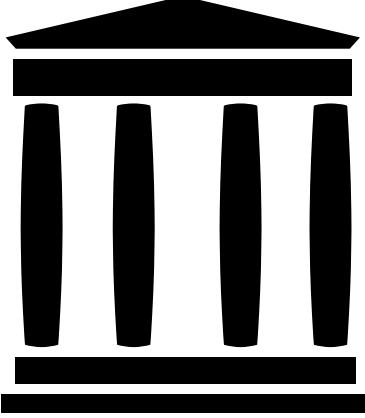


Basis: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, May 2020.



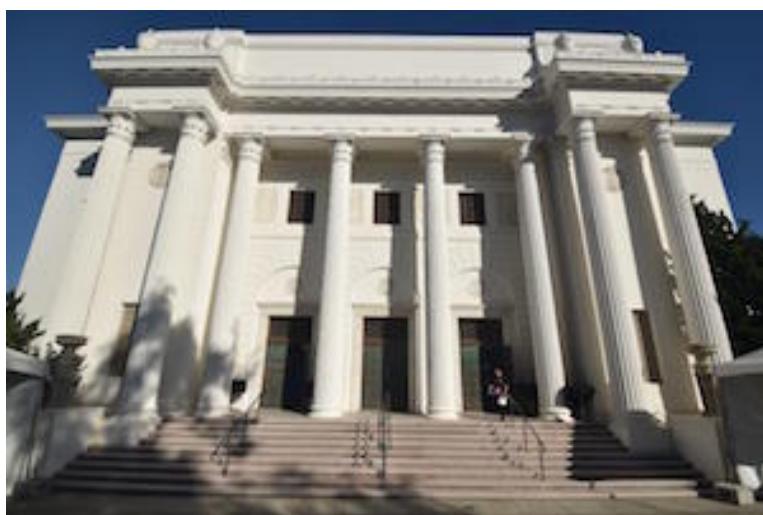
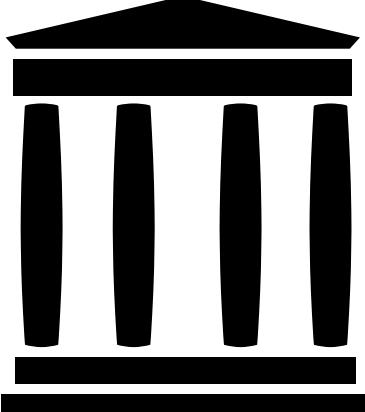
The Internet Archive

INTERNET ARCHIVE



- 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

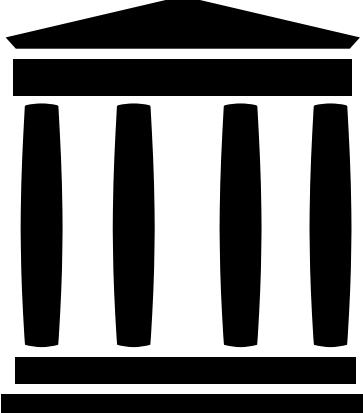
INTERNET ARCHIVE



- Founded 1996 by Brewster Kahle
- For all things digital:
 - 477 billion web pages (ca. 30PB) – accessible via the **Wayback Machine**
 - 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

INTERNET ARCHIVE
wayback machine

INTERNET ARCHIVE



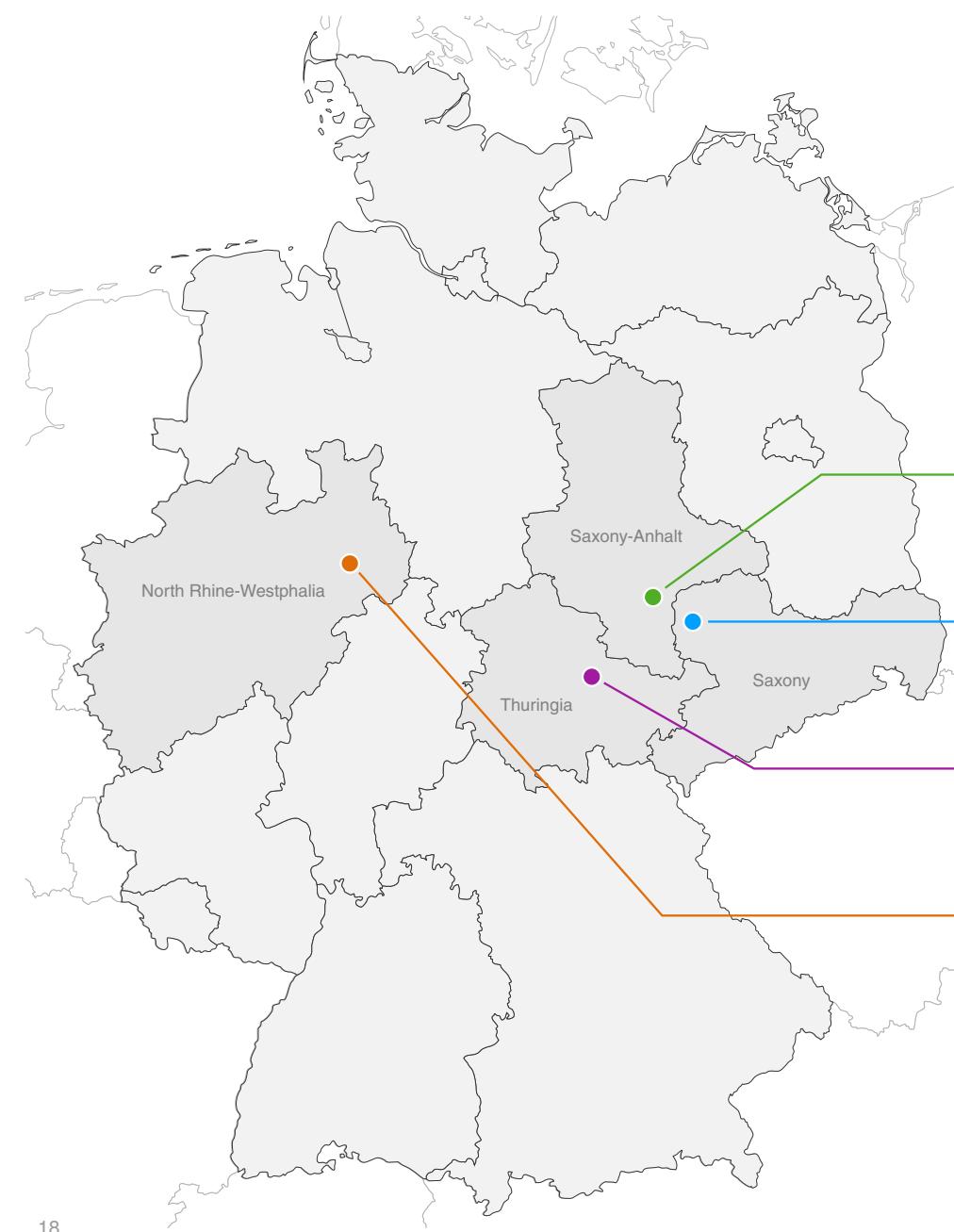
Mission: “Universal access to all knowledge.”

- 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]



Web Archive Analytics @ Webis



MLU Halle-Wittenberg

Prof. Dr. Matthias Hagen

Leipzig University

Prof. Dr. Martin Potthast

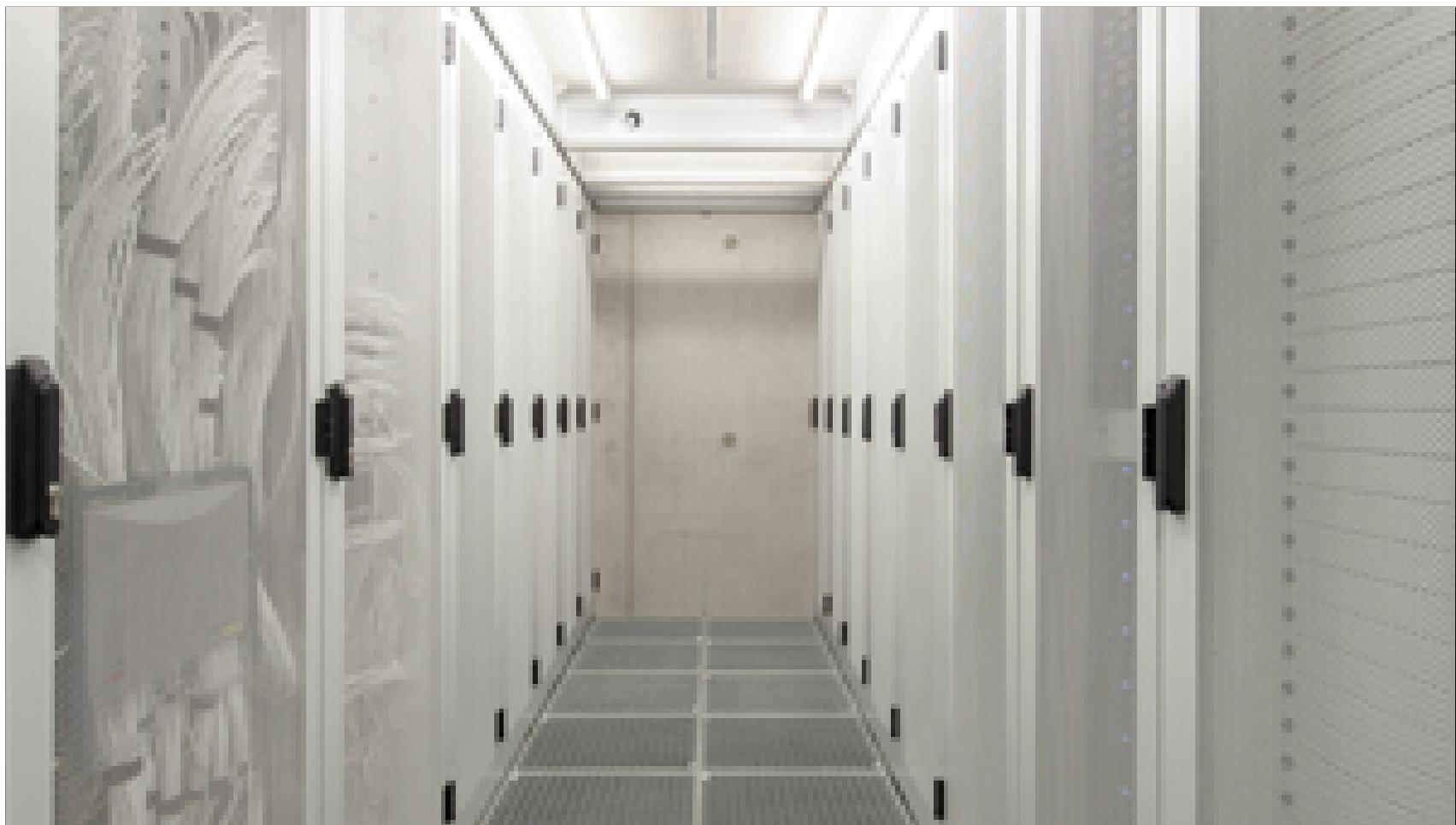
Bauhaus-Universität Weimar

Prof. Dr. Benno Stein

Paderborn University

Prof. Dr. Henning Wachsmuth

Webis Data Center (Digital Bauhaus Lab)



Webis Data Center (Digital Bauhaus Lab)

	α -web [2009]	β -web [2015]	γ -web [2016 + 2020]	δ -web [2018]	ϵ -web [2020]
Nodes	44	135	9	78	55
Disk [PB]	0.2	4.1	0.01	12	0.1
Cores	176	1,740	384 + 227,328	1,248	1,100
	$\cong 3.2 \text{ TFLOPs}$	$\cong 67.4 \text{ TFLOPs}$	$\cong 690.0 \text{ TFLOPs}$	$\cong 119.8 \text{ TFLOPs}$	$\cong 44 \text{ TFLOPs}$
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 + 2020]	δ -web [2018]	ϵ -web [2020]
Nodes	44	135	9	78	55
Disk [PB]	0.2	4.1	0.01	12	0.1
Cores	176	1,740	384 + 227,328	1,248	1,100
	$\approx 3.2 \text{ TFLOPs}$	$\approx 67.4 \text{ TFLOPs}$	$\approx 690.0 \text{ TFLOPs}$	$\approx 119.8 \text{ TFLOPs}$	$\approx 44 \text{ TFLOPs}$
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

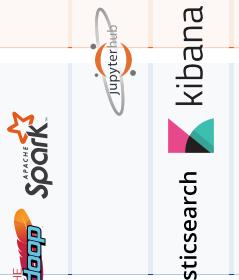
Webis Analytics Stack

Vendor stack

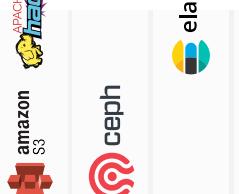
Data Consumption Layer



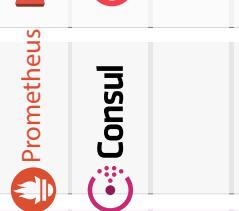
Data Analytics Layer



Data Management Layer



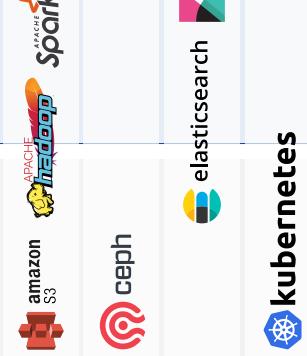
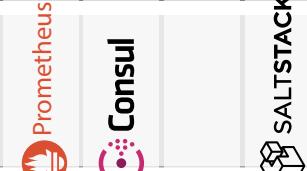
Hardware Layer



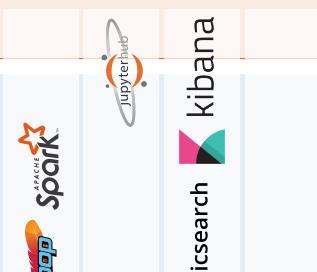
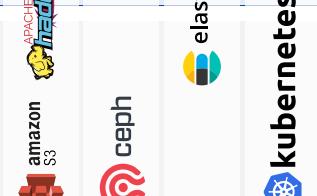
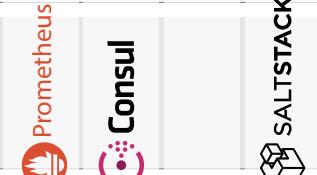
Data Acquisition Layer



Webis Analytics Stack

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

Webis Analytics Stack

	Task Stack	Technology stack	Vendor stack
Data Consumption Layer	- Query and explore - Visualize and interact - Explain and justify	- Visual analytics - Immersive technologies - Intelligent agents	
Data Analytics Layer	- Diagnose and reason - Structure identification - Structure verification	- Distributed learning - State-space search - Symbolic inference	
Data Management Layer	- Provenance tracking - Normalization - Cleansing	- Key-value store - RDF triple store - Graph store - Object store	
Hardware Layer	- Monitoring - Replication	- Orchestration - Parallelization - Virtualization	
Data Acquisition Layer	- Replay - Collect - Log	- Distant supervision - Crowdsourcing - Crawling and archiving	

Webis Analytics Stack

	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 		Experts: <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



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 Wikidata.

- 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 “**Tqbfjot1d**” 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.

□ Who Wrote the Web?

Applying author identification technology at web-scale.

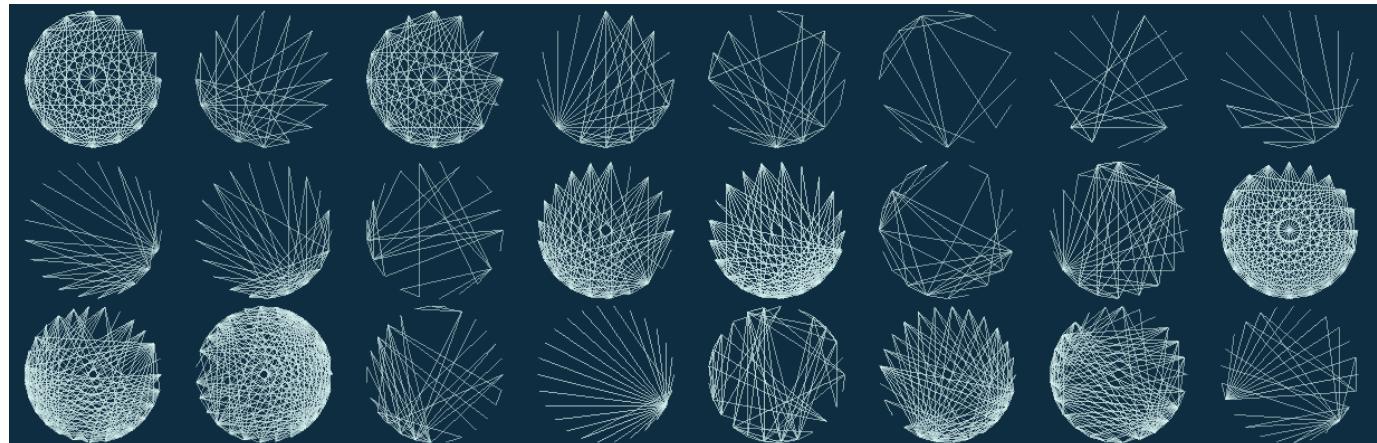
□ Text Reuse Analytics

Goals: (1) Finding the Wikipedia within the Wikipedia.

(2) Quantifying the impact of Wikipedia on the “rest of the web.”

□ Text Reuse Illustration

Example: Visualizing article similarities in Wikipedia.



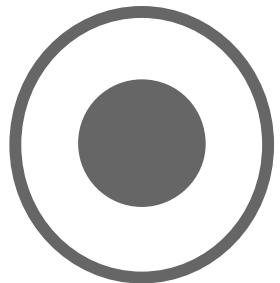
Riemann et al.:
*Visualizing Article
Similarities in
Wikipedia.*
EuroVis 2016

❑ Abstractive Snippet Generation

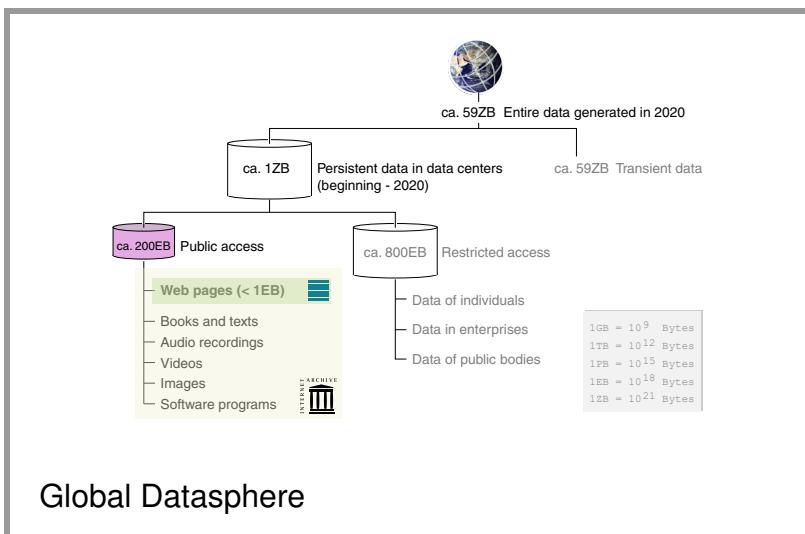
Approach: Use of anchor contexts to generate abstractive snippets with a pointer-generator network, exploiting ClueWeb09, ClueWeb12, and the DMOZ Open Directory Project.

❑ Learn Automatic Summarization

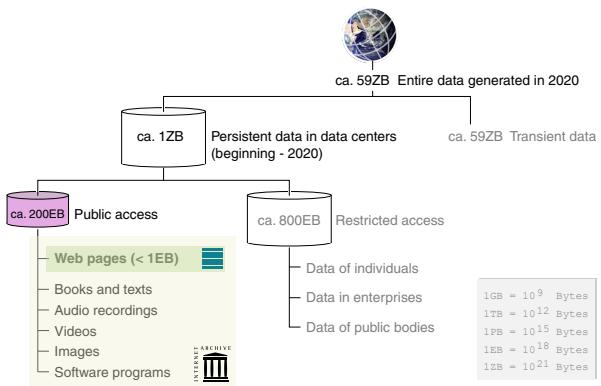
Approach: Exploit author-provided summaries, taking advantage of the common practice of appending a “TL;DR” to long posts.



Summary



Summary

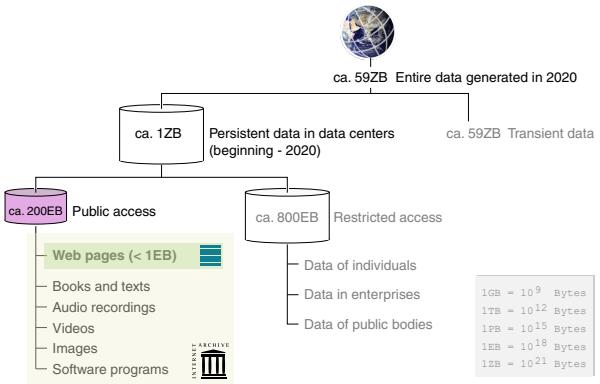


Global Datasphere



Internet Archive

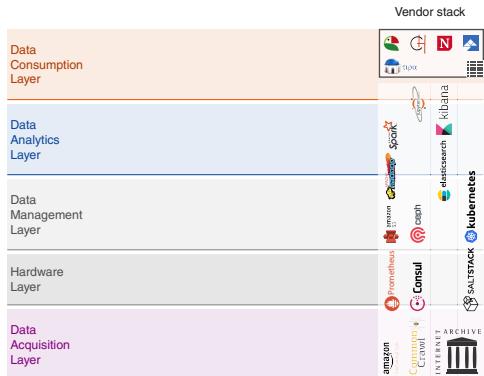
Summary



Global Datasphere

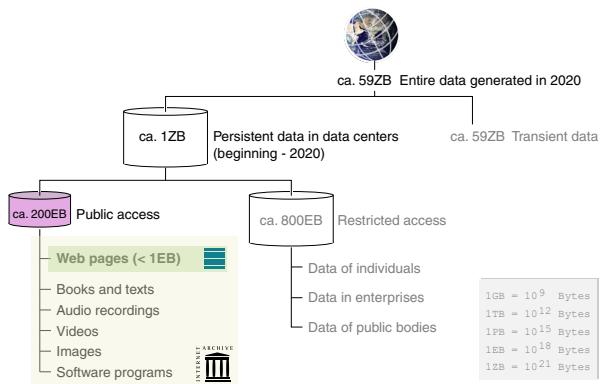


Internet Archive



Webis Analytics Stack

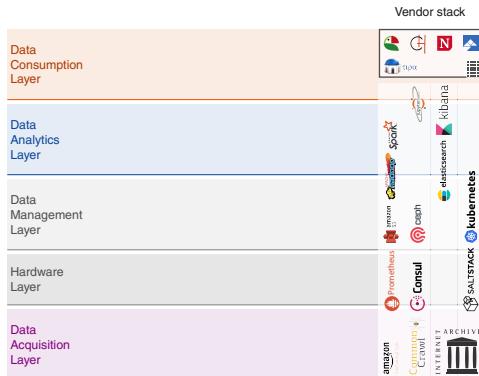
Summary



Global Datasphere



Internet Archive



Webis Analytics Stack



Webis Archive Research

Thank You!