# Towards an Open Web Index
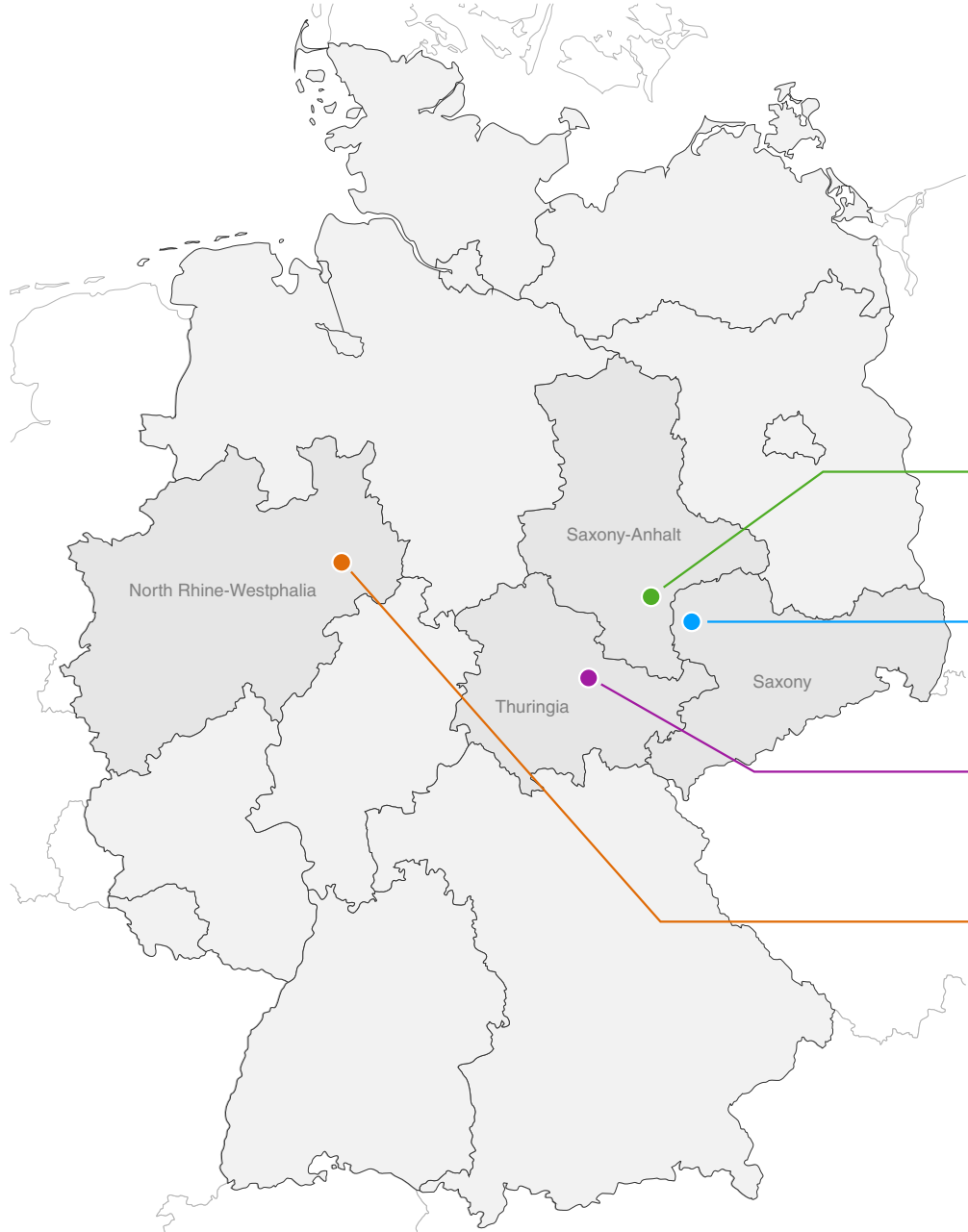## Lessons from the Past

Michael Völske
Bauhaus-Universität Weimar
webis.de

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

# Open Web Search

# Open Web Search

## Challenges

Independence

Scale

User data

Market penetration

Funding

Transparency

# Open Web Search
## Challenges

Independence

Scale

User data

Market penetration

Funding

Transparency

# Independence

## The Search Market in 2020
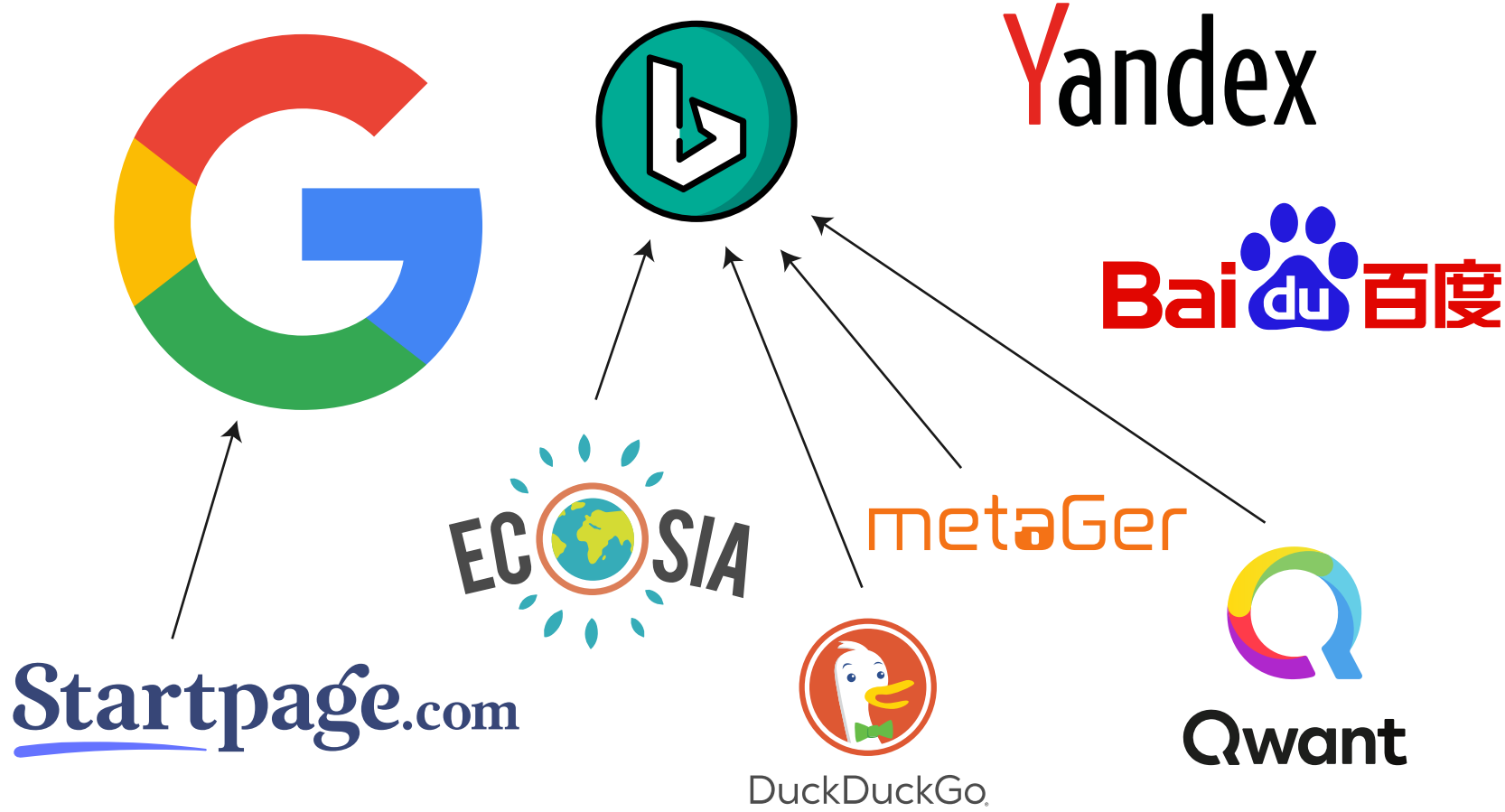
# Independence
## The Search Market in 2020

# Independence
## The Search Market in 2020

# Independence
## The Search Market in 2020

Völske@Webis

Scale

# The Global Datasphere

# The Global Datasphere

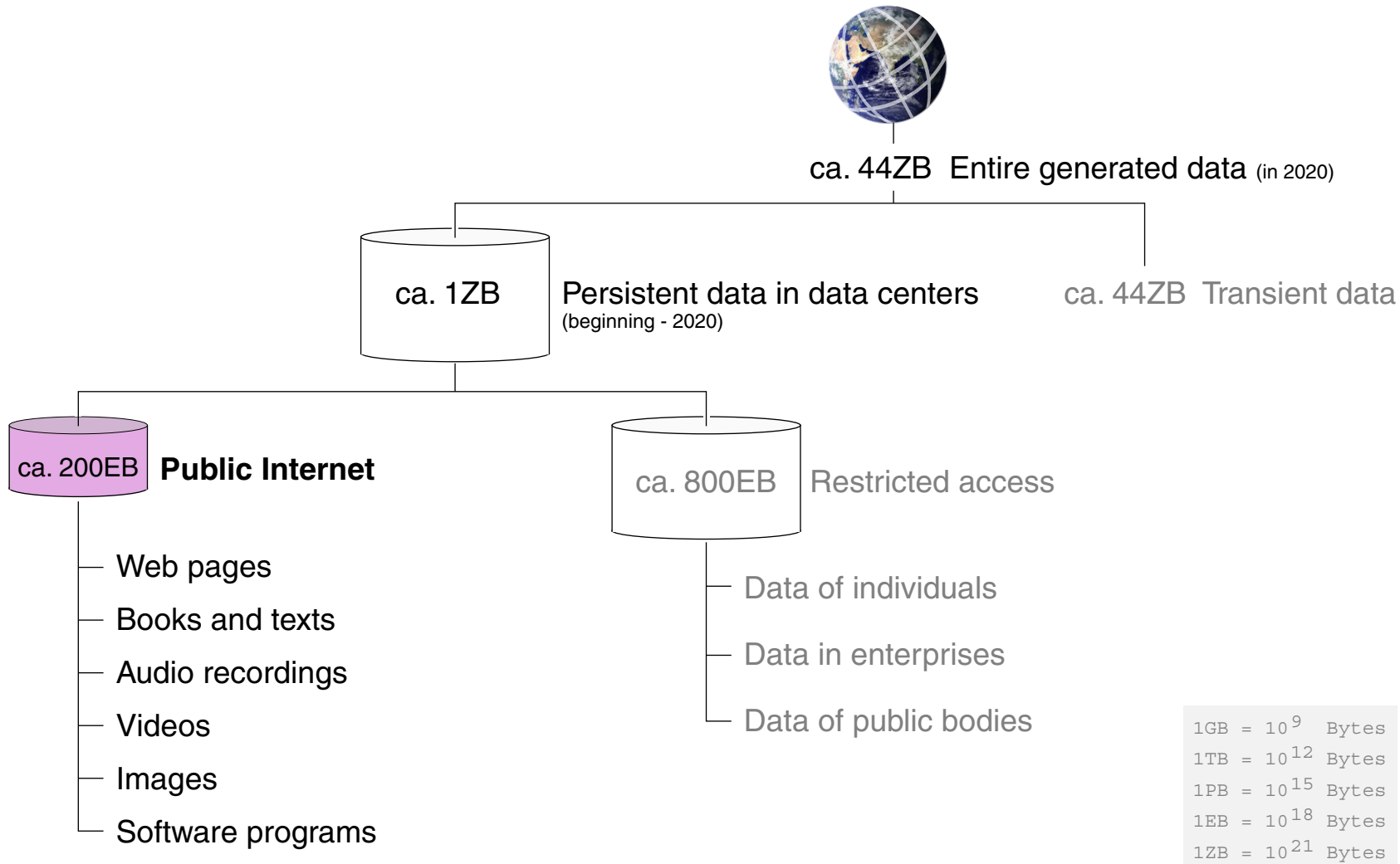*"A measure of all new data captured, created, and replicated in a single year."*

[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, …"*

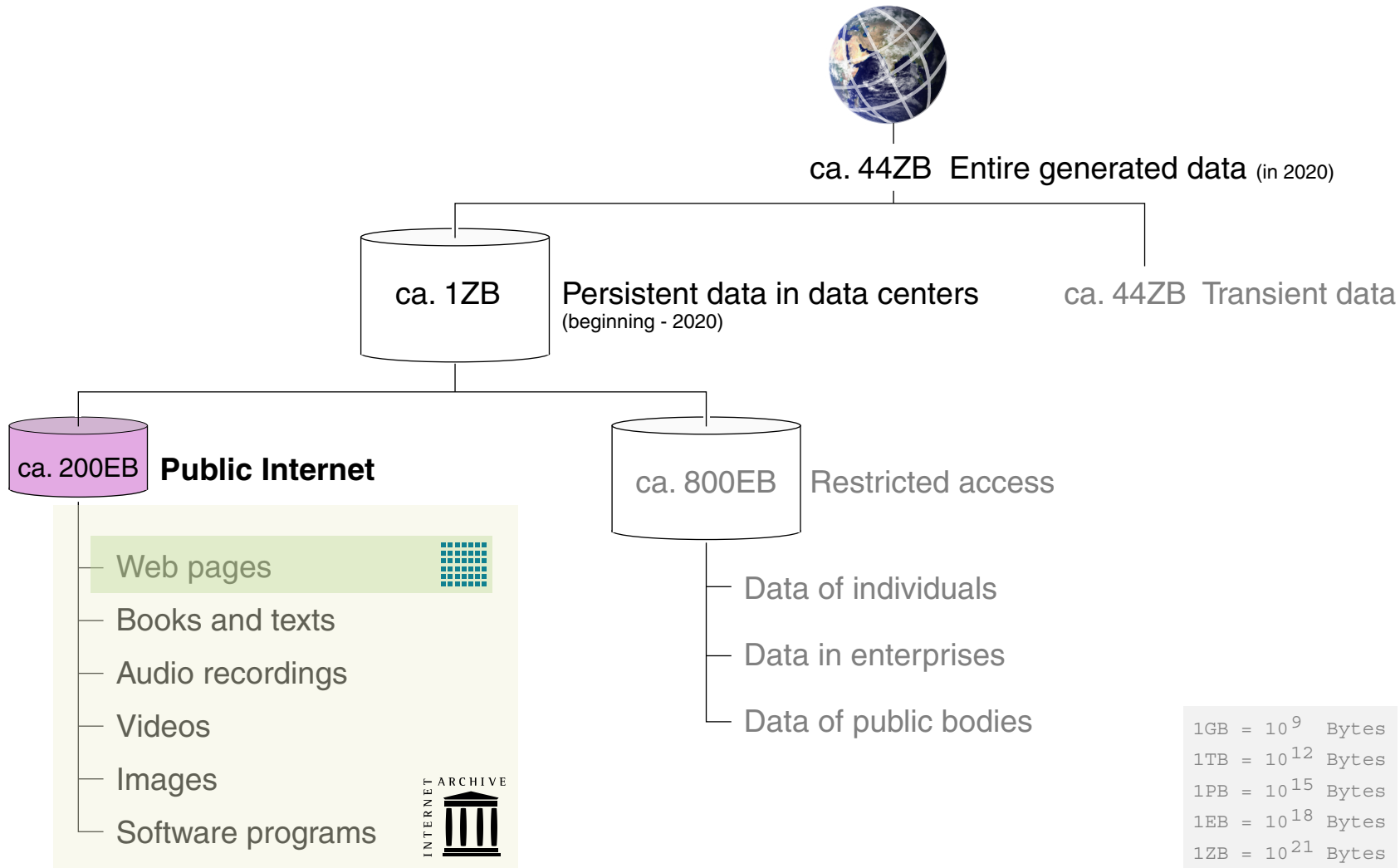[IDC, 2012]

Völske@Webis

# The Global Datasphere

ca. 44ZB  Entire generated data (in 2020)

ca. 1ZB  Persistent data in data centers
(beginning - 2020)

ca. 44ZB  Transient data

ca. 200EB  **Public Internet**

ca. 800EB  Restricted access

— Web pages

— 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
```

Völske@Webis

# The Global Datasphere

ca. 44ZB  Entire generated data (in 2020)

ca. 1ZB  Persistent data in data centers (beginning - 2020)

ca. 44ZB  Transient data

ca. 200EB  **Public Internet**

- Web pages
- 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
```

Völske@Webis

# The Global Datasphere
## Relating Data Source Sizes

200EB: Public Internet

15EB: Google

30PB: Web pages @ Internet Archive

8PB: Web pages @ Webis

200TB: Wikipedia including Wikimedia

200GB: English Wikipedia including media

2GB: All English Wikipedia article texts

| $10^9$ | $10^{12}$ | $10^{15}$ | $10^{18}$ | $10^{21}$ | Bytes |
|--------|-----------|-----------|-----------|-----------|-------|
| Giga | Tera | Peta | Exa | Zetta | |

# The Global Datasphere

## Relating Data Source Sizes

200EB: Public Internet

15EB: Google

30PB: Web pages @ Internet Archive

8PB: Web pages @ Webis

200TB: Wikipedia including Wikimedia

200GB: English Wikipedia including media

2GB: All English Wikipedia article texts

$10^9$        $10^{12}$        $10^{15}$        $10^{18}$        $10^{21}$    Bytes

Giga         Tera         Peta         Exa         Zetta

# The Global Datasphere

## *Where* is the Data Stored?



Legend:
- Consumer
- Enterprise
- Public Cloud

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

# The Global Datasphere

## *Where* is the Data Stored?



Consumer
Enterprise
Public Cloud

Among others:

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

User Data

# User Data

## Benefits of Implicit Relevance Feedback



**Basis:** Agichtein et al., Improving Web Search Ranking by Incorporating User Behavior Information, SIGIR (2006)

Völske@Webis

# User Data

## Risks: Record Linkage

| Timestamp | UID | URL |
|---|---|---|
| 2020-10-13 14:06 | A1234 | sportsnews.com/sports/ |
| 2020-10-13 14:07 | A1234 | news.com/myaccount/solso |
| 2020-10-13 14:07 | A1234 | ecommerce.com/receipt_id=1234054 |
| 2020-10-13 14:08 | A1234 | sportsnews.com/sports |
| 2020-10-13 14:06 | B5678 | sportsnews.com/sports |
| 2020-10-13 14:07 | B5678 | sportsnews.com/sports |

Völske@Webis

# User Data

| Domain | Unique Visitors |
|---|---|
| sportsnews.com | 2 |
| news.com | 1 |
| ecommerce.com | 1 |

| Timestamp | UID | URL |
|---|---|---|
| 2020-10-13 14:06 | A1234 | sportsnews.com/sp |
| 2020-10-13 14:07 | A1234 | news.com/myaccount/soiso |
| 2020-10-13 14:07 | A1234 | ecommerce.com/receipt_id=1234054 |
| 2020-10-13 14:08 | A1234 | sportsnews.com/sports |
| 2020-10-13 14:06 | B5678 | sportsnews.com/sports |
| 2020-10-13 14:07 | B5678 | sportsnews.com/sports |

# User Data
## Risks: Record Linkage

| Domain | Unique Visitors |
|---|---|
| sportsnews.com | 2 |
| news.com | 1 |
| ecommerce.com | 1 |

| Timestamp | UID | URL |
|---|---|---|
| 2020-10-13 14:06 | A1234 | sportsnews.com/sp... |
| 2020-10-13 14:07 | A1234 | news.com/myaccount/soiso |
| 2020-10-13 14:07 | A1234 | ecommerce.com/receipt_id=1234054 |
| 2020-10-13 14:08 | A1234 | sportsnews.com/sports |
| 2020-10-13 14:06 | B5678 | sportsnews.com/sports |
| 2020-10-13 14:07 | B5678 | sportsn... |

| URL | UID |
|---|---|
| sportsnews.com/sports | A1234 |
| news.com/myaccount/alice | A1234 |
| ecommerce.com/receipt_id=1234054 | A1234 |

# User Data

Risks: Record Linkage

| Timestamp | UID | URL |
|---|---|---|
| 2020-10-13 14:06 | A1234 | sportsnews.com/sp... |
| 2020-10-13 14:07 | A1234 | news.com/myaccount/soiso |
| 2020-10-13 14:07 | A1234 | ecommerce.com/receipt_id=1234054 |
| 2020-10-13 14:08 | A1234 | sportsnews.com/sports |
| 2020-10-13 14:06 | B5678 | sportsnews.com/sports |
| 2020-10-13 14:07 | B5678 | sportsn... |

| Domain | Unique Visitors |
|---|---|
| sportsnews.com | 2 |
| news.com | 1 |
| ecommerce.com | 1 |

| URL | UID |
|---|---|
| sportsnews.com/sports | A1234 |
| news.com/myaccount/alice | A1234 |
| ecommerce.com/receipt_id=1234054 | A1234 |

# User Data
## Risks: Record Linkage

| Domain | Unique Visitors |
|---|---|
| sportsnews.com | 2 |
| news.com | 1 |
| ecommerce.com | 1 |

| Timestamp | UID | URL |
|---|---|---|
| 2020-10-1... | ... | sportsnews.com/sp... |
| 2020-10-1... | ... | ...count/soiso |
| 2020-10-... | ... | receipt_id=1234054 |
| 2020-10-13 14:05 | A1234 | sportsnews.com/sports |
| 2020-10-13 14:06 | B5678 | sportsnews.com/sports |
| 2020-10-13 14:07 | B5678 | sportsn... |

| URL | UID |
|---|---|
| sportsnews.com/sports | A1234 |
| news.com/myaccount/alice | A1234 |
| ecommerce.com/receipt_id=1234054 | A1234 |

# User Data
## Traditional Data Collection

# User Data
## Privacy-preserving Data Collection

CLIQZ [2015 – 2020]

Human Web & Human Proxy Network

# User Data

## Privacy-preserving Data Collection

**CLIQZ** [2015 – 2020]

Human Web & Human Proxy Network

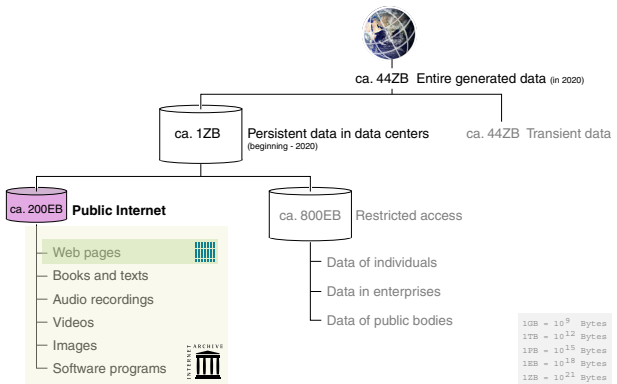# Summary

Independence

Scale

User data

Market penetration

Funding

Transparency

Challenges

# Summary



Challenges

Independence

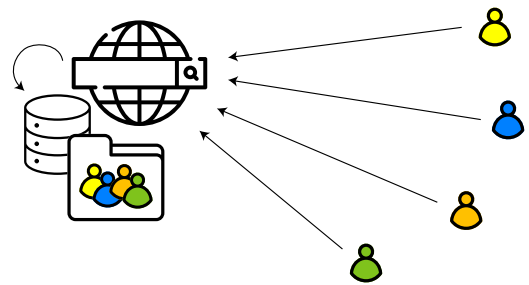# Summary


Challenges


Independence


Scale

# Summary


Challenges


Independence


Scale


User Data

Thank You!