

# **Day 6: Data management and ethics**

## Guiding Principles for Data Management II

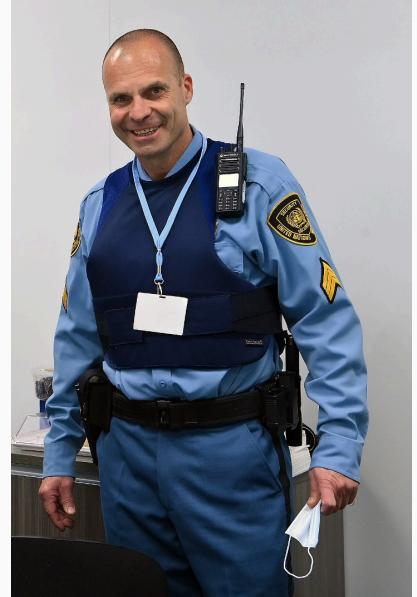
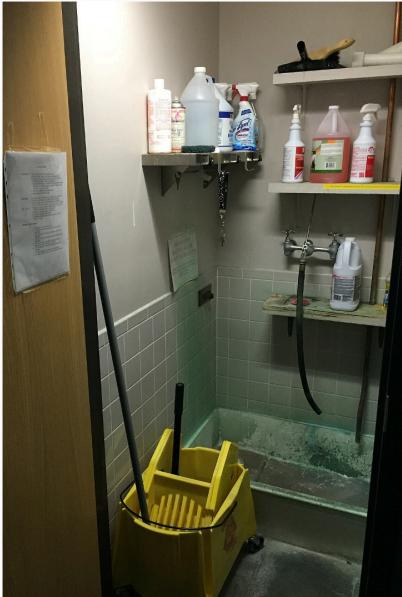
---

Sebastian Ramirez-Ruiz  
Hertie School

# What people think working with data looks like...

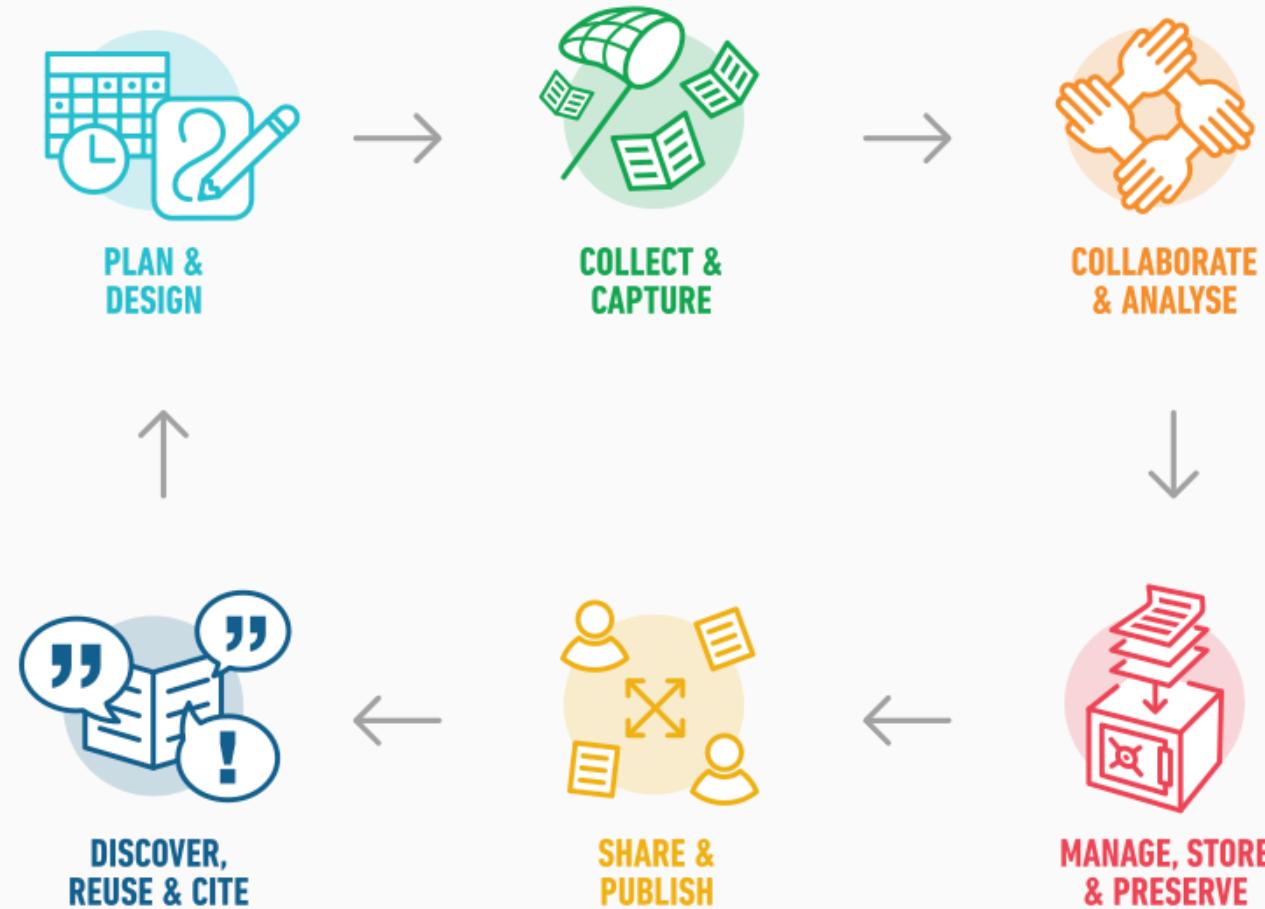


# How it really is...



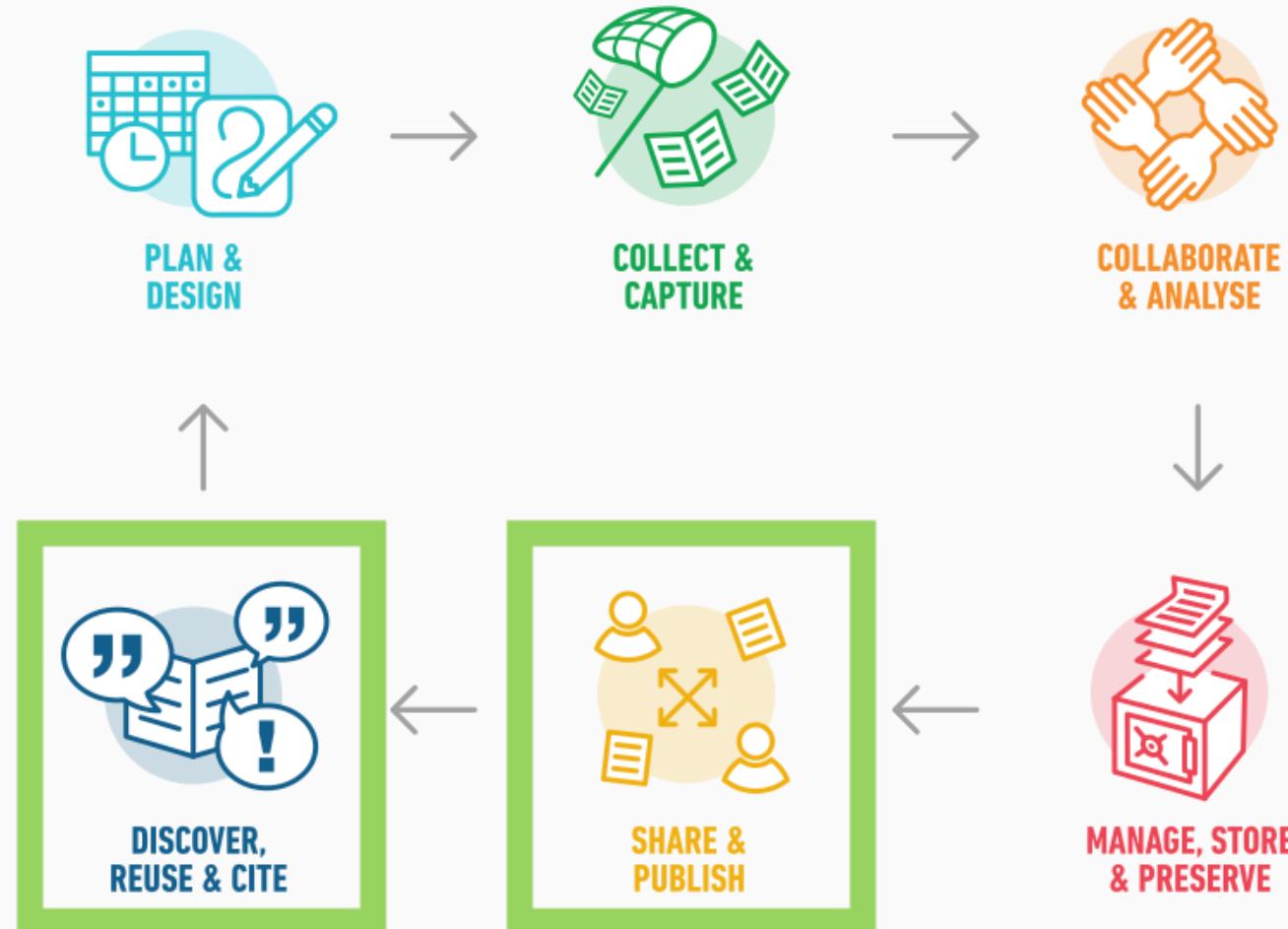
You have to wear many hats...

# Research data management (RDM) lifecycle



1. Archiving and publishing
2. Open government data
3. Discovering

# Research data management (RDM) lifecycle



## Raise your hand if you:

- Can state the difference between data archiving and data publishing
- Can list some of the benefits of data publishing
- Can differentiate between different data publication services (data journals, self-archiving, a data repository)
- Know which data repositories fit the needs of different data sources
- Know of ways to promote your government data publication

# **Archiving and publishing**

---

High-quality data have the potential to be reused in many ways. *Archiving* and *publishing* your data properly is at the core of making your data FAIR and will enable both your future self as well as others to get the most out of your data.

## Archiving data for future reference

**Data archiving** is about storing and preserving data for the long term. When you archive your data, you make sure you can read and access the data later on. You can then also allow access to others for verification purposes when such a request arrives. In all cases, you should store your data safely, in a suitable file format, with adequate documentation.

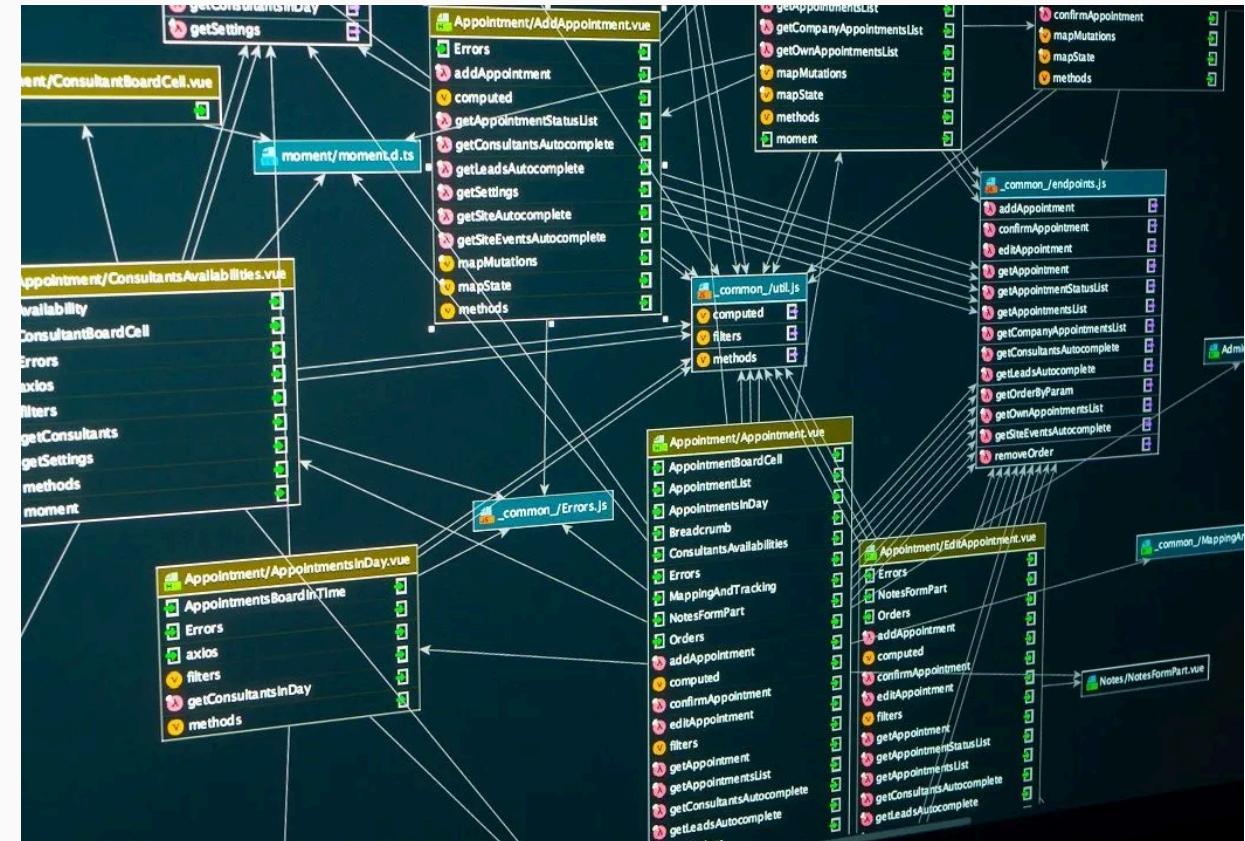


# Data publishing

## Publishing data for reuse

To make your data reusable for purposes beyond the one for which you collected them, you can

**publish your data**. Publishing your data is the act of publicly disclosing the research data you have collected, making them findable, accessible, and reusable.



## Scientific progress

Data archiving and publication has direct benefits for the research itself (more robust), for the discipline and for science in general by enabling new collaborations, new data uses and establishing links to the next generation of researchers.

## Norms of the field

Depending on field may be more prone. General practice in *computational social science*.

## Drivers

- Funders
- Publishers



**Should you publish your data or shouldn't you?**

**Should you publish your data or shouldn't you? And if so, which part of it?**

## Should you publish your data or shouldn't you? And if so, which part of it?

- Sometimes this question is **straightforward to answer:**
  - the funder of the research demands it
  - strong researcher inclination
- Still, *not all data are created equal* and data publishing does involve an investment of researcher resources. **Some datasets have a more obvious reuse potential than others.**

## Does your dataset have reuse potential?

- Does your data have potential value in terms of reuse,
  - national/international standing and quality,
  - historical importance,
  - uniqueness,
  - originality,
  - size,
  - scale,
  - costs of data production or innovative nature of the research?
- Could you foresee that secondary analyses on your data would benefit science? Or policy analysis?

If your answer to any of these is yes, your dataset has serious reuse potential.

## Is your dataset reusable?

To be suitable for reuse, your dataset must be

### **functionally usable.**

- Can the data be accessed and utilized? (e.g., *readable format*)
- Is there sufficient metadata available to allow future users to comprehend your data? (e.g., *documentation*)
- Are there any legal restrictions that prevent the data from being published?

**If you have addressed these practical considerations  
and your data hold potential value for reuse, you are  
'good to go!'**

# Example from our own experience

## Publishing Combined Web Tracking and Survey Data

Simon Munzert<sup>1</sup>

Sebastian Ramirez-Ruiz<sup>1</sup>

Oliver Watteler<sup>2</sup>

Johannes Breuer<sup>2,3</sup>

Veronika Batzdorfer<sup>2</sup>

Christina Eder<sup>2</sup>

Deborah Wiltshire<sup>2</sup>

Pablo Barberá<sup>4</sup>

Andrew Guess<sup>5</sup>

JungHwan Yang<sup>6</sup>

<sup>1</sup> Hertie School

<sup>2</sup> GESIS - Leibniz Institute for the Social Sciences

<sup>3</sup> Center for Advanced Internet Studies (CAIS)

<sup>4</sup> University of Southern California

<sup>5</sup> Princeton University

<sup>6</sup> University of Illinois at Urbana-Champaign

Source : Munzert, Simon, et al. "Publishing Combined Web Tracking and Survey Data." (2023).

<https://doi.org/10.31219/osf.io/y4v8z>

# How do I protect the privacy of my research subjects?



POLITICS  
& Media



## Table of Contents

<b>1. Demographics</b>	<b>8</b>
Wave number (wave)	8
Respondent ID (person_id)	8
Survey start time (starttime)	8
Survey end time (endtime)	9
Weight (weight)	9
Gender (gender)	9
Year of birth (birthyr)	10
Race (race)	10
Education (educ)	10
Marital status (marstat)	11
Children under age 18 in household (child18)	11
Employment (employ)	12
Household income (faminc_new)	12
State (inputstate)	13
Religion (religpew)	15
Born Again or Evangelical Christian (pew_bornagain)	15
Religious service attendance (pew_churatd)	16
Religion importance (pew_religimp)	16
Prayer frequency (pew_prayer)	17

Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

# How do I protect the privacy of my research subjects?

```
R> dat$url[1:20]

## [1] "mail.google.com/mail/u/0#inbox"
## [2] "bing.com/search?q=fu%C3%9Fnagel+w%C3%A4chst+nach+oben+und+verhornt&form=EDGHPT&qs=PF&cvid=a97dc782996"
## [3] "preisvergleich-pilot.de/sport-und-freizeit.html"
## [4] "http://tile-service.weather.microsoft.com:443"
## [5] "http://3c.1und1.de:443"
## [6] "https://www.trendfrage.de/"
## [7] "questler.de/intern/paidmail"
## [8] "http://www.ebesucher.de/surfbar/swhtmw.106_tobias_wutz"
## [9] "http://www.google.de:443"
## [10] "http://www.ebesucher.de/surfbar/dermaerker"
## [11] "tracking.dpd.de/parcelstatus?query=01505160077643&locale=en_DE"
## [12] "ma-shops.de/shops/search.php?searchstr=kreuzz%C3%BCge&catid=904&submitBtn=Finden&lang=de&PHPSESSID="
## [13] "studien.usuma.de/cc3/survey/intweb.dll/project/wsusuma/CAWI170607_0G1/username=GMI&passwd=WW2017&PIN=2"
## [14] "amazon.de/Samsung-J5550-Fernseher-Triple-Tuner/dp/B00U57X80I/ref=sr_1_4?s=home-theater&ie=UTF8&qid=15"
## [15] "keypanel.de/?ID=13"
## [16] "http://outlook.office.com:443"
## [17] "ci-marketing.de/_p.php?userid=3673&mailid=229933&mc=6495367958"
## [18] "http://www.farmerama.com:443"
## [19] "https://www.amazon.de/b/ref=gbph_ftr_m-8_ee12_page_263?node=7361369031&pf_rd_p=7bde80db-c7fc-47f3-a66"
## [20] "http://www.webmie.com:443"
```

Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

# How do I protect the privacy of my research subjects?

```
R> # parse Google requests
R> google_maps_df <- filter(dat, str_detect(url, "/maps"))
R>
R> # how much of the URLs are Google Maps URLs?
R> nrow(google_maps_df)

## [1] 12206

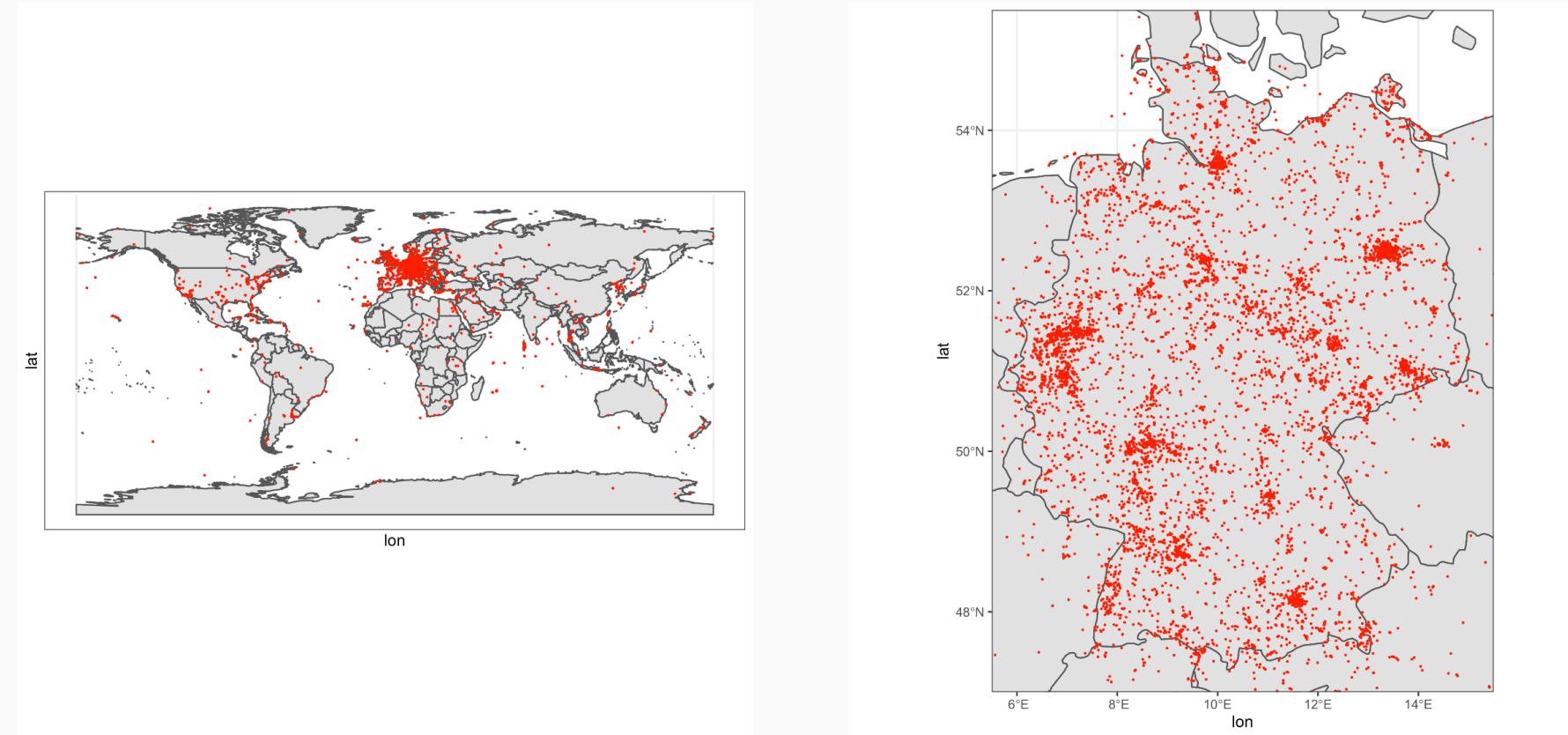
R> nrow(google_maps_df)/nrow(dat) # about 0.6%

## [1] 0.006104038

R> # extract coordinates
R> google_maps_df$coord_raw <- str_extract(google_maps_df$url, "@[-[:digit:]]{1,4}\\.\\.[[:digit:]]+,-[-[:digit:]]{1,4}
R> google_maps_df$lat <- str_extract(google_maps_df$coord_raw, "[-[:digit:]]{1,4}\\.\\.[[:digit:]]+") %>% as.numeric
R> google_maps_df$lon <- str_extract(google_maps_df$coord_raw, "[-[:digit:]]{1,4}\\.\\.[[:digit:]]+\$") %>% as.numeric
```

Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

# How do I protect the privacy of my research subjects?



Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

# How do I protect the privacy of my research subjects?

```
R> # parse ebay sites
R> ebay_df <- filter(dat, str_detect(url, "ebay"))
R>
R> # plz
R> ebay_plz <- ebay_df %>% filter(domain == "ebay-kleinanzeigen.de", str_detect(path, "^s\\-[[:digit:]]{5}")) %>%
R> length(ebay_plz)

## [1] 319

R> ebay_plz[1:20]

##  [1] "86720"  "92660"  "26655"  "97789"  "52457"  "21149"  "24983"  "78234"  "66346"
## [10] "08523"  "14715"  "01936"  "21149"  "12529"  "14715"  "72415"  "78244"  "21149"
## [19] "24306"  "86720"
```

Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

# How do I protect the privacy of my research subjects?

```
R> # parse ebay sites
R> ebay_df <- filter(dat, str_detect(url, "ebay"))
R>
R> # plz
R> ebay_plz <- ebay_df %>% filter(domain == "ebay-kleinanzeigen.de", str_detect(path, "^s\\-[[:digit:]]{5}")) %>%
R> length(ebay_plz)

## [1] 319

R> ebay_plz[1:20]

##  [1] "86720"  "92660"  "26655"  "97789"  "52457"  "21149"  "24983"  "78234"  "66346"
## [10] "08523"  "14715"  "01936"  "21149"  "12529"  "14715"  "72415"  "78244"  "21149"
## [19] "24306"  "86720"
```

Source Munzert, Simon, Barberá, Pablo, Guess, Andrew M., & Yang, JungHwan (2022). Media Exposure and Opinion Formation in an Age of Information Overload (MEOF). GESIS, Cologne. <https://doi.org/10.4232/1.13981>.

Let's take one minute to discuss in which ways **publishing these data can be challenging.**

---

If you publish your data in a data repository, you may  
need to apply a licence agreement to your data

**If you publish your data in a data repository, you may  
need to apply a licence agreement to your data**

- A licence agreement is a legal arrangement between the depositor of the data set and the repository
  - It stipulates what users are allowed to do with the data

If you publish your data in a data repository, you may need to apply a licence agreement to your data

- A licence agreement is a legal arrangement between the depositor of the data set and the repository
  - It stipulates what users are allowed to do with the data

To make re-use as likely as possible you may want to choose a licence which:

- Makes data available to the **widest audience** possible
- Makes the **widest range of uses** possible

# Overview of Creative Commons licences

Licence	Can I copy & redistribute the work?	Is it required to attribute the author?	Can I use the work commercially?	Am I allowed to adapt the work?	Can I change the licence when redistributing?
CC0	Y	N	Y	Y	Y
CC BY	Y	Y	Y	Y	Y
CC BY-SA	Y	Y	Y	Y	N
CC BY-ND	Y	Y	Y	N	Y
CC BY-NC	Y	Y	N	Y	Y
CC BY-NC-SA	Y	Y	N	Y	N
CC BY-NC-ND	Y	Y	N	N	Y

This table is inspired by the CESSDA Guide

**Publishing data in a data repository does not automatically make them openly accessible.**

Publishing data in a data repository does not automatically make them openly accessible. **(Sensitive) personal data can still be protected by limiting access to the data.**

Publishing data in a data repository does not automatically make them openly accessible. (Sensitive) personal data can still be protected by limiting access to the data. **Access controls can permit control down to an individual file level, meaning that mixed levels of access control can be applied to a data collection.**

Publishing data in a data repository does not automatically make them openly accessible. (Sensitive) personal data can still be protected by limiting access to the data. Access controls can permit control down to an individual file level, meaning that mixed levels of access control can be applied to a data collection.

## Many data repositories operate a three-tiered approach to data access:

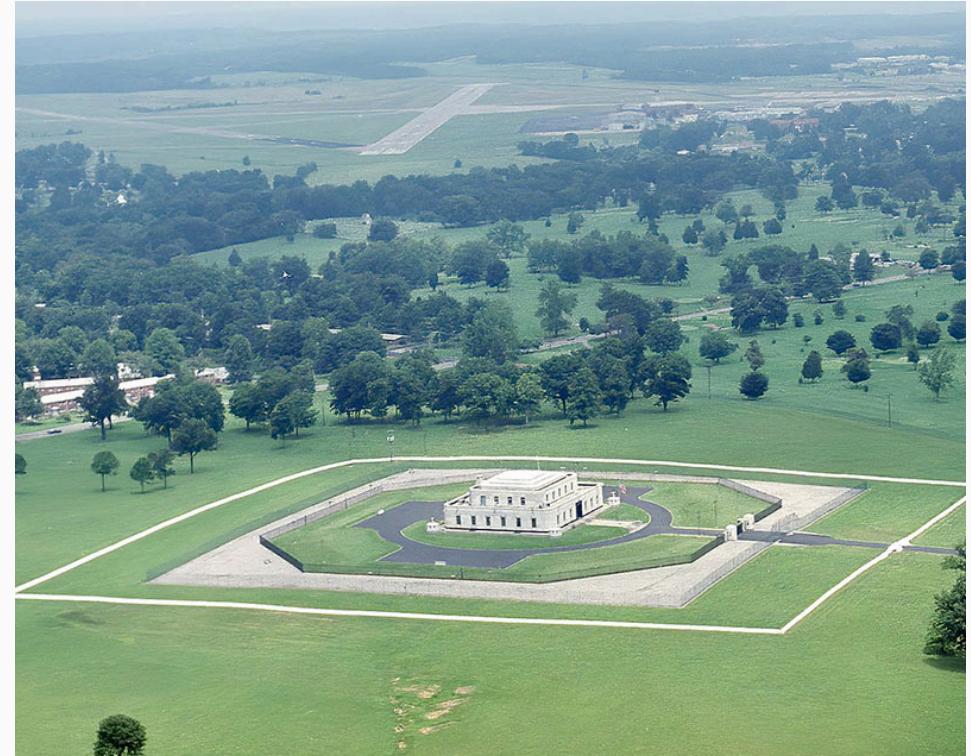
- **Open access**
  - Data that can be accessed by any user whether they are registered or not. Data in this category should not contain personal information unless consent is given (see 'Informed consent').
- **Access for registered users** (safeguarded)
  - Data that is accessible only to users who have registered with the archive. This data contains no direct identifiers but there may be a risk of disclosure through the linking of indirect identifiers.
- **Restricted access**
  - Access is limited and can only be granted upon request. This access category is for the most sensitive data that may contain disclosive information. Restricted access requires the long-term commitment of the researcher or person responsible for the data to handle the upcoming permission requests.
- **Embargo**
  - Besides offering the opportunity for restricted access 'for eternity' most data repositories allow you to place a temporary embargo on your data. During the embargo period, only the description of the dataset is published. The data themselves will become available in open access after a certain period of time.

## Special access

- Access through a *safe connection*
- Access in a *safe environment*

## Open metadata for (sensitive) personal data

This was one of the solutions for our survey and web-tracking data archiving and publishing effort.



## Special access

- Access through a *safe connection*
- Access in a *safe environment*

## Open metadata for (sensitive) personal data

This was one of the solutions for our survey and web-tracking data archiving and publishing effort.

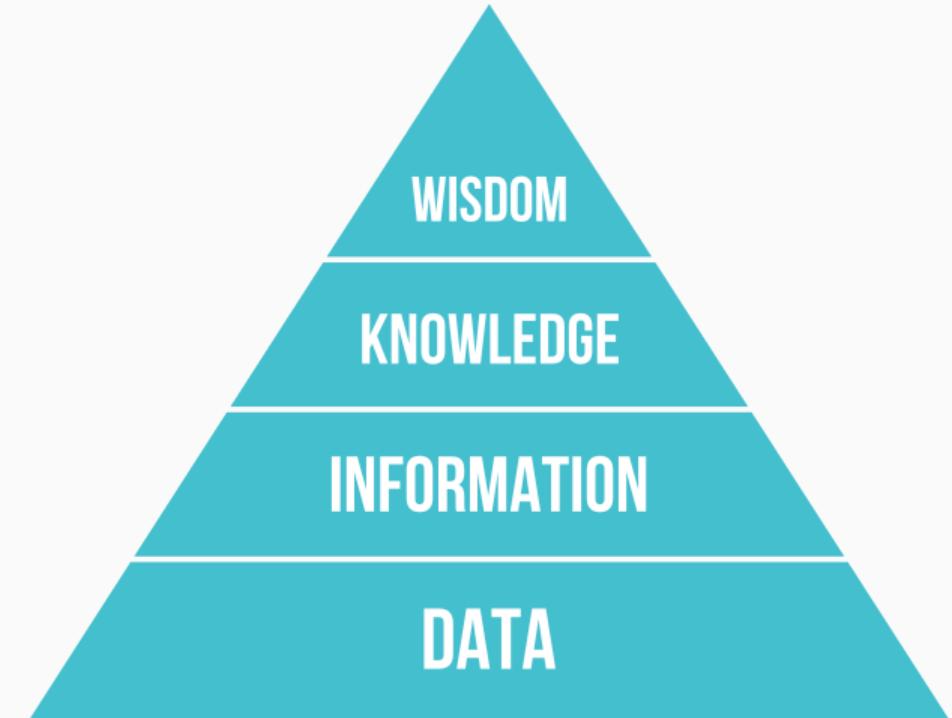


# **Open Government Data**

---

## What are data?

- Data are disembodied **facts, signs, and symbols.**
- We often define them by their **source** (e.g. *administrative, historical, medical, etc.*) and their **formats** (e.g., numerical, textual, still image, geospatial, audio, video, and software.)
- It is thought about as the basis of the **knowledge hierarchy** under the *DIKW* pyramid.



# What is open data?

## Open data

- Data that can be freely *used*, *modified*, and *shared* by anyone

## What characteristics does it have?

- **Non-proprietary**: Available in formats that are open and standard without any restrictions
- **Machine-Readable**: Provided in formats that can be easily processed by computers
- **No Restrictions**: Free from any legal, financial, or technical barriers.
- **Comprehensive**: Covers all necessary aspects to provide a complete understanding

The screenshot shows the homepage of data.gov.uk. At the top, there's a dark header with the text "data.gov.uk | Find open data" and a "Menu ▾" button. Below the header, a blue banner says "BETA This is a new service – your [feedback](#) will help us to improve it". The main title "Find open data" is centered above a search bar with the placeholder "Search data.gov.uk" and a magnifying glass icon. Below the search bar, there are three categories: "Business and economy" (with a link to small businesses, industry, imports, exports and trade), "Government" (with a link to staff numbers and pay, local councillors and department business plans), and "Towns and cities" (with a link to housing, urban planning, leisure, waste and energy consumption).

# Why Open Government Data (OGD)?

- Public organizations **produce** and **collect** a broad range of different types of data to perform their tasks.
  - Demographic information
  - Socioeconomic data
  - Epidemiological data
  - Healthcare use
  - Industry data
  - Crime data
  - Performance
  - ...
- **These data tend to be extensive and central**

**Government data can be a significant resource for increased transparency**

# The actors behind Open Government Data

There are two key civil society actors advocating for increased openness of information, documents and datasets held by public bodies

- **Right to information movement:** promotes a public right of access to information from a **human rights perspective**
- **Open government data movement:** uses predominantly **social and economic arguments** to encourage the opening up of government data, such as:
  - Information hitting the public domain can create conditions for inclusive service delivery and more participation
  - Can stimulate the economy by allowing third parties to innovate using public data

Many argue that these data should be public domain, since **governments hold predominantly data emanating from citizens.**

# The Open Government Data Principles

These principles were presented in 2007 during an **Open Government Working Group Meeting**

Government data shall be considered open if it is made public in a way that **complies with the principles** below:

- **Complete** : All public data are made available. Public data are data that is not subject to valid privacy, security or privilege limitations.
- **Primary** : Data are as collected at the source, with the highest possible level of granularity, not in aggregate or modified forms.
- **Timely** : Data are made available as quickly as necessary to preserve the value of the data.
- **Accessible** : Data are available to the widest range of users for the widest range of purposes.
- **Machine processable** : Data are reasonably structured to allow automated processing.
- **Non-discriminatory** : Data are available to anyone, with no requirement of registration.
- **Non-proprietary** : Data are available in a format over which no entity has exclusive control.
- **License-free** : Data are not subject to any copyright, patent, trademark or trade secret regulation. Reasonable privacy, security and privilege restrictions may be allowed.

10 principles proposed by Former US Federal Chief Information Officer **Vivek Kundra** before the US House of Representatives in 2011

- **Build end-to-end digital processes** : Automate transfer of data between systems to increase productivity, protect data integrity, and speed data dissemination. Capitalize on game-changing technologies to increase transparency.
- **Build once, use often** : Architect systems for re-use and share platforms to reduce costs, streamline systems and processes, reduce errors, and foster collaboration.
- **Tap into golden sources of data** : Pull data directly from authoritative sources to improve data quality, shorten processes and protect data integrity.
- **Release machine-readable data and encourage third party applications** : Make data machine-readable to allow the public to easily analyse, visualise and use government information.
- **Use common data standards** : Develop and use uniform, unique identifiers and data standards to ease the flow of data and reduce system complexity.

# Principles for Improving Federal Transparency in the US

10 principles proposed by Former US Federal Chief Information Officer **Vivek Kundra** before the US House of Representatives

- **Validate data up front** : **Correct errors during collection** and at the point of entry to block bad data from ever entering the system.
- **Release data in real time and preserve for future use** : **Release data as quickly as feasible** to enhance its relevance and utility while maintaining future accessibility.
- **Reduce burden** : **Collect data once and use it repeatedly**. Pull from existing data sets to reduce costs and burden and to increase productivity and uniformity.
- **Protect privacy and security** : **Safeguard the release of information to increase public trust**, participation, preserve privacy, and protect national security. Open Government doesn't mean vulnerable government.
- **Provide equal access and incorporate user feedback** : Provide a common view of data to all **stakeholders to foster collaboration**. Incorporate user feedback to help identify high-value, meaningful data sets, set priorities, to continuously drive and improve future planning and processes.

## Business Models Archetypes for value creation

Some actors make a business argument for Open Government Data. They identify some archetypes of business models that can create revenue from this practice:

- **Suppliers** that publish data, including the public sector;
- **Aggregators** that pool publicly available data and combine it to produce useful insights to be used by the various users;
- **Apps developers** that enable users to make more informed decisions (e.g. apps building on crime data, transport data);
- **Enrichers** that are large and established businesses producing significant quantities of open data and combine it with their own proprietary sources to provide services (e.g. insurers, retailers); and
- **Enablers** that are organisations that don't make money out of open data but provide platforms and technologies that others can use (e.g. websites that enable data sources of all types to make subsets of their data available to seek solutions from the public).

## Government

### Benefits

- **Efficiency:** Improves government operations and resource allocation.
- **Fraud Reduction:** Helps reduce fraud and errors.
- **Innovative Services:** Enables smarter, personalized public services.
- **Transparency:** Strengthens accountability and legitimacy.
- **Quality Interactions:** Enhances interactions between government and users.

### Examples

- **Tax Gap Reduction:** More accurate tax collection.
- **Public Service Delivery:** Innovative apps and services

## Citizens

### Benefits

- **Participation:** Encourages public participation in designing responses to public needs.
- **Informed Choices:** Provides information for making informed personal decisions.
- **Quality of Life:** Potentially improves overall quality of life through better services.

### Examples

- **Fix My Street (UK):** Citizens report local issues.
- **Chicago 311:** Government portal for public service requests.
- **Personal Data Access:** Crime rates, emissions, education stats.

## Civil society

### Benefits

- **Informed advocacy:** Facilitates the identification of public needs.
- **Watchdogs:** Open government data allows civil society organizations (CSOs) to track government activities, expenditures, and decisions.
- **Facilitating partnerships:** Can encourage collaboration between CSOs, government agencies, and other stakeholders

### Examples

- **Transparency International:** Corruption tracking.
- **Code for America:** Government data to build digital tools that improve public services.
- **Global Witness:** Investigate and expose illegal activities related to natural resources (land use, mining permits, deforestation)

## 1. Data to fact

Individuals seek specific facts in newly open datasets.

### Applications:

- Civic or bureaucratic engagement
- Business planning
- Personal choices

### Methods:

- Online interfaces
- Downloading spreadsheets

## 2. Data to information

Creation of static representations and interpretations of data.

### Manifestations:

- Visualizations
- Blog posts
- Infographics
- Written reports

## 3. Data to Interface

Interactive access and exploration of datasets.

### Examples:

- Searchable mapping mash-ups
- Tools for browsing large datasets and crowdsourcing feedback

### Components:

- Static data interpretations
- Summary statistics
- Algorithmically derived assessments

## 4. Data to Data

Sharing derived data in new formats or augmented forms.

### Methods:

- Reformatting original datasets
- Combining with other data
- Creating APIs
- Easy downloads of dataset subsets

### Advantages:

- Real-time updates
- Avoids re-downloading entire files

## 5. Data to Service

Services where OGD supports operations behind the scenes.

### Examples:

- Routing messages reporting issues like potholes to the appropriate authority
- Use of boundary data for service efficiency

# The five star open data scheme

The five-star data openness scheme was developed and presented in 2010 by **Tim Berners-Lee**  
**a co-creator of the world wide web.**

- ★ data resources that have been made available on the network under the terms of an open license (in any format);
- ★★ resources made available in the form of structured data (for instance, an spreadsheet instead of a scan)
- ★★★ data in the document saved in an open format such as CSV;
- ★★★★ URI-tagged resources that are totally searchable;
- ★★★★★ data combined with different data that provides context to it.

- **Policy challenges**

- Disclosure policies, limits in data transparency, copyright issues
- Lack of procedures and standards

- **Technical Challenges**

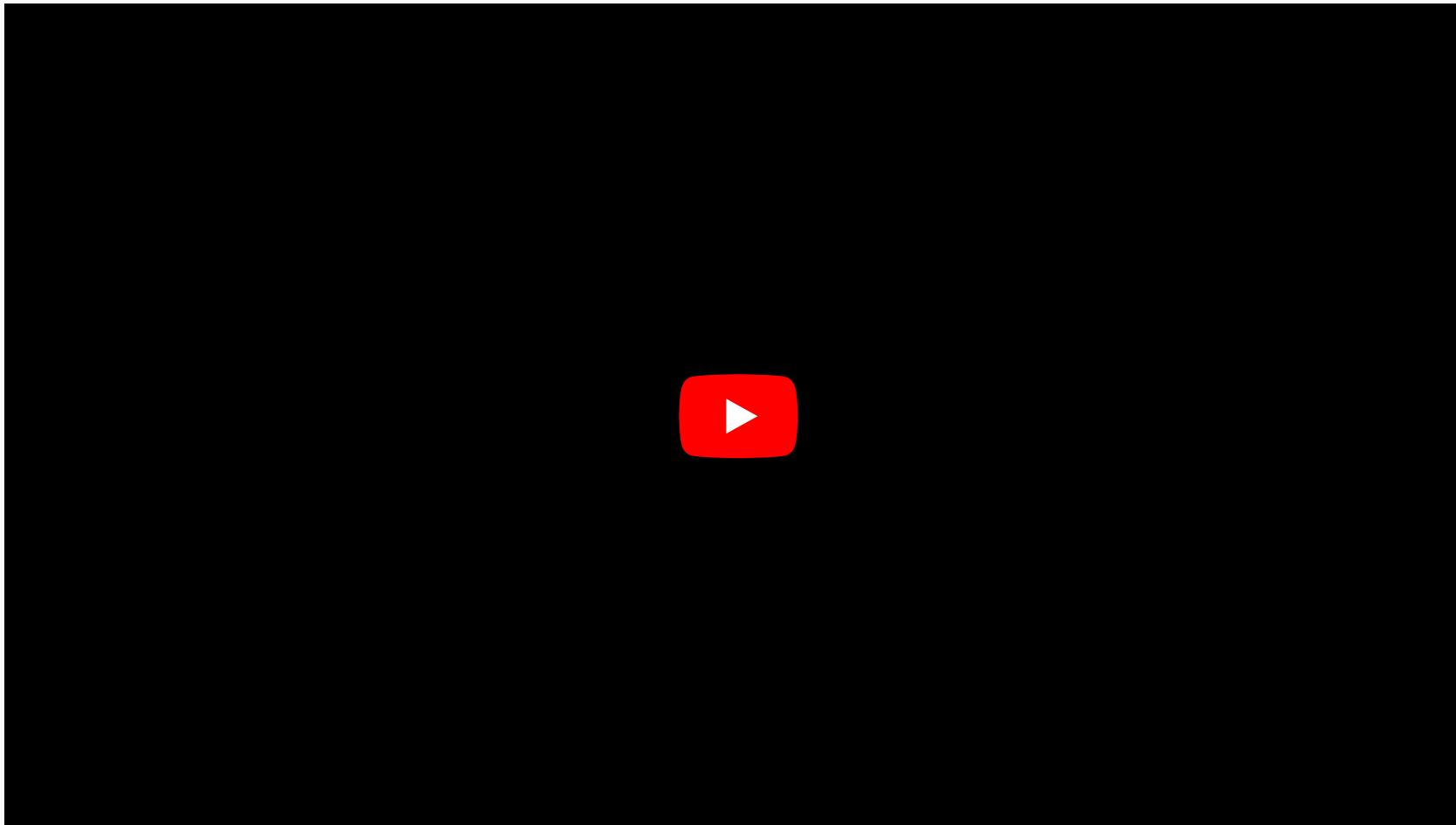
- Variability in data sets, formats, and standards
- Complexity of data access and source identification
- Need for IT infrastructure and privacy enhancement

- **Economic and Financial Challenges**

- Financial burden on governments
- Training, technology, and infrastructure costs
- Proprietary software and sensitive data issues

- **Organizational Challenges**
  - Institutional structures and leadership
  - Independent oversight and data transparency
  - Supporting data release workflows and cultural orientation
- **Cultural Challenges**
  - Public interest and awareness
  - Stakeholder buy-in
- **Legal Challenges**
  - Variability in legal frameworks across countries
  - National security, privacy, and commercialization issues
  - Ambiguities in intellectual property rights

# Responses to PSI Directive



The screenshot shows the CKAN for Government homepage. At the top, there is a navigation bar with links for Home, Features, Showcase, Solutions (with a dropdown arrow), Support (with a dropdown arrow), Blog, Events, FAQ, Docs (with a dropdown arrow), a 'View on Github' button, and a Login link. Below the navigation bar, the title 'CKAN for Government' is displayed in a large, bold, black font. To the left of the title, there is a paragraph of text: 'CKAN is used by national and regional government organisations throughout the European Union, the Americas, Asia and Oceania to power a variety of official and community data portals.' Below this text is a button labeled 'Schedule a chat'. To the right of the text is a large, semi-transparent image of a classical government building with a prominent blue dome and white columns. The image is partially obscured by a large, light gray triangle.

CKAN is used by national and regional government organisations throughout the European Union, the Americas, Asia and Oceania to power a variety of official and community data portals.

Schedule a chat

# Discovering

---

# Why data discover?

- Save costs and time

 HARVARD  
Dataverse

Add Data ▾ Search ▾ About User Guide Support Sign Up Log In

Deposit and share your data. Get academic credit.  
Harvard Dataverse is a repository for research data. Deposit data and code here.

Organize datasets and gather metrics in your own repository.  
A dataverse is a container for all your datasets, files, and metadata.

Publishing your data is easy on Harvard Dataverse!  
Learn about getting started creating your own dataverse repository here.

Add a dataset + Add a dataverse + Getting started ↗

---

Find data across research fields, preview metadata, and download files

Search over 178,300 datasets...  VIEW ALL DATA ↗

**Featured**  COVID-19 Data Collection  
A curated collection of COVID-19 data deposited in the Harvard Dataverse repository.

Browse by subject

Agricultural Sciences 4,870	Computer and Information Science 3,746	Medicine, Health and Life Sciences 10,415
Arts and Humanities 36,677	Earth and Environmental Sciences 9,508	Physics 1,700
Astronomy and Astrophysics 1,343	Engineering 2,253	Social Sciences 63,928
Business and Management 2,294	Law 5,818	



# Why data discover?

- Save costs and time
- Compare results and replicate studies

The screenshot shows the Harvard Dataverse homepage. At the top, there's a navigation bar with links for Add Data, Search, About, User Guide, Support, Sign Up, and Log In. Below the header, there are three main sections: "Deposit and share your data. Get academic credit.", "Organize datasets and gather metrics in your own repository.", and "Publishing your data is easy on Harvard Dataverse!". Each section has a "Add a dataset" or "Add a dataverse" button. A large search bar below these sections allows users to search over 178,300 datasets. A featured collection for "COVID-19 Data Collection" is highlighted with a thumbnail of the virus and a description. At the bottom, there's a "Browse by subject" section with links to various fields like Agricultural Sciences, Arts and Humanities, etc., followed by a "Feedback" button.

# Why data discover?

- Save costs and time
- Compare results and replicate studies
- Integrate into a larger environment

The screenshot shows the Harvard Dataverse homepage. At the top, there's a navigation bar with links for Add Data, Search, About, User Guide, Support, Sign Up, and Log In. Below the navigation, there are three main sections: "Deposit and share your data. Get academic credit.", "Organize datasets and gather metrics in your own repository.", and "Publishing your data is easy on Harvard Dataverse!". Each section has a "Add a dataset" or "Add a dataverse" button. A central search bar allows users to search over 178,300 datasets. Below the search bar, a "COVID-19 Data Collection" is highlighted as a featured curated collection. The page also includes a "Feedback" button at the bottom right.

## What do you need?

- Identify the domain
- Measures and constructs
- What data types

## Locate potential data hosts

- Ask (you've identified some people working on it)
- Data repositories
- Search engines and data aggregators (e.g., [Google Dataset search](#))
- Data catalogues ([GESIS](#)) and data journals ([Scientific Data](#))

## Set up a search query

## Ask for help

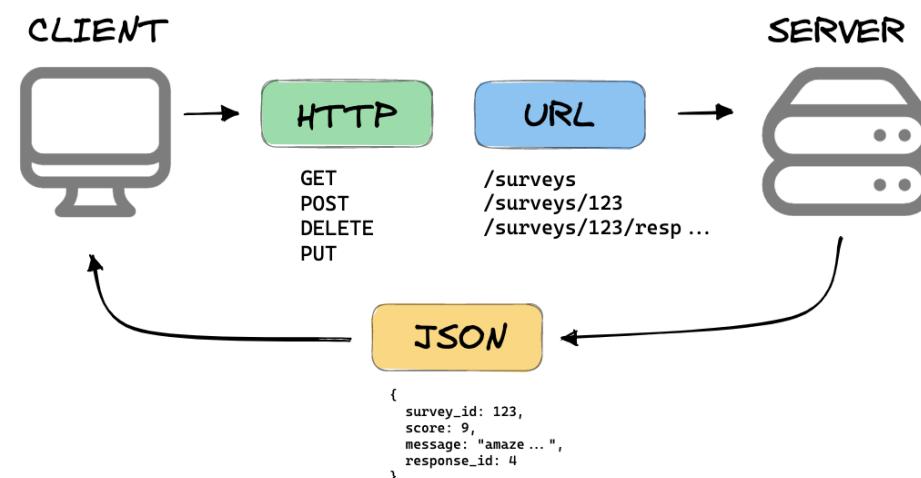
## Evaluate data quality

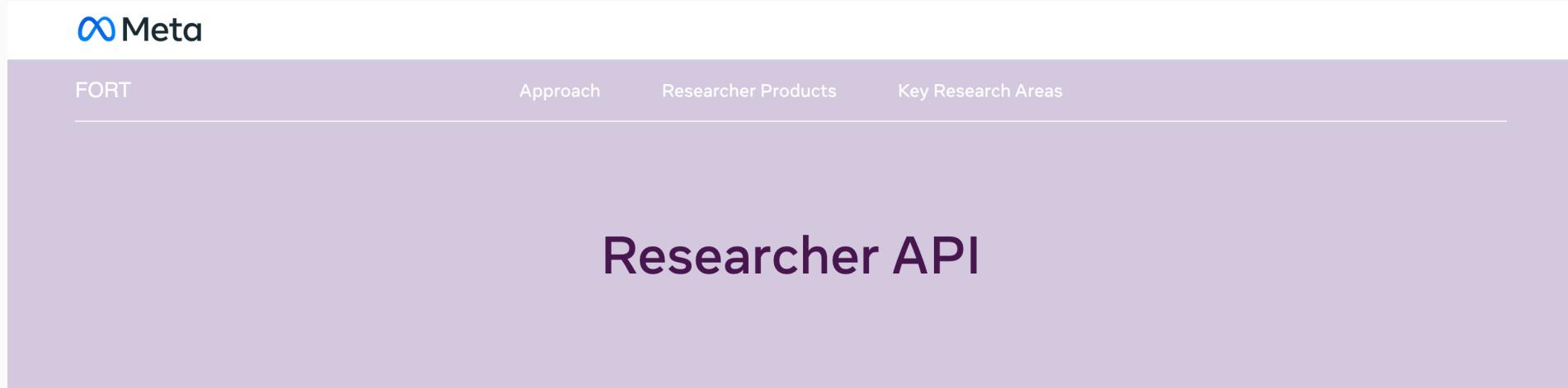
- **What** information was collected?
- **Who** collected the data? **When**? And, **where**?
- **Why** was the data created? E.g., different purposes for data collection are research, social policy, marketing etc.
- **How** was the data collected? You need detailed information about the **methodology**.
- How was the data **processed**?
  - Were there any changes in data? Who adjusted data in what way after it was collected?
  - To which manipulations was the data exposed?
- Were **consistency** and **logic checks** employed?
  - Is the data “clean”, i.e. were nonlogical and erroneous values deleted?
- What quality assurance procedures were used? Did researchers use verified measurement tools?
  - Documentation

# APIs (Application Programming Interface)

A set of **rules and protocols** for building and interacting with software applications. Act as intermediaries that allow different software systems to communicate with each other.

## WHAT IS A REST API?

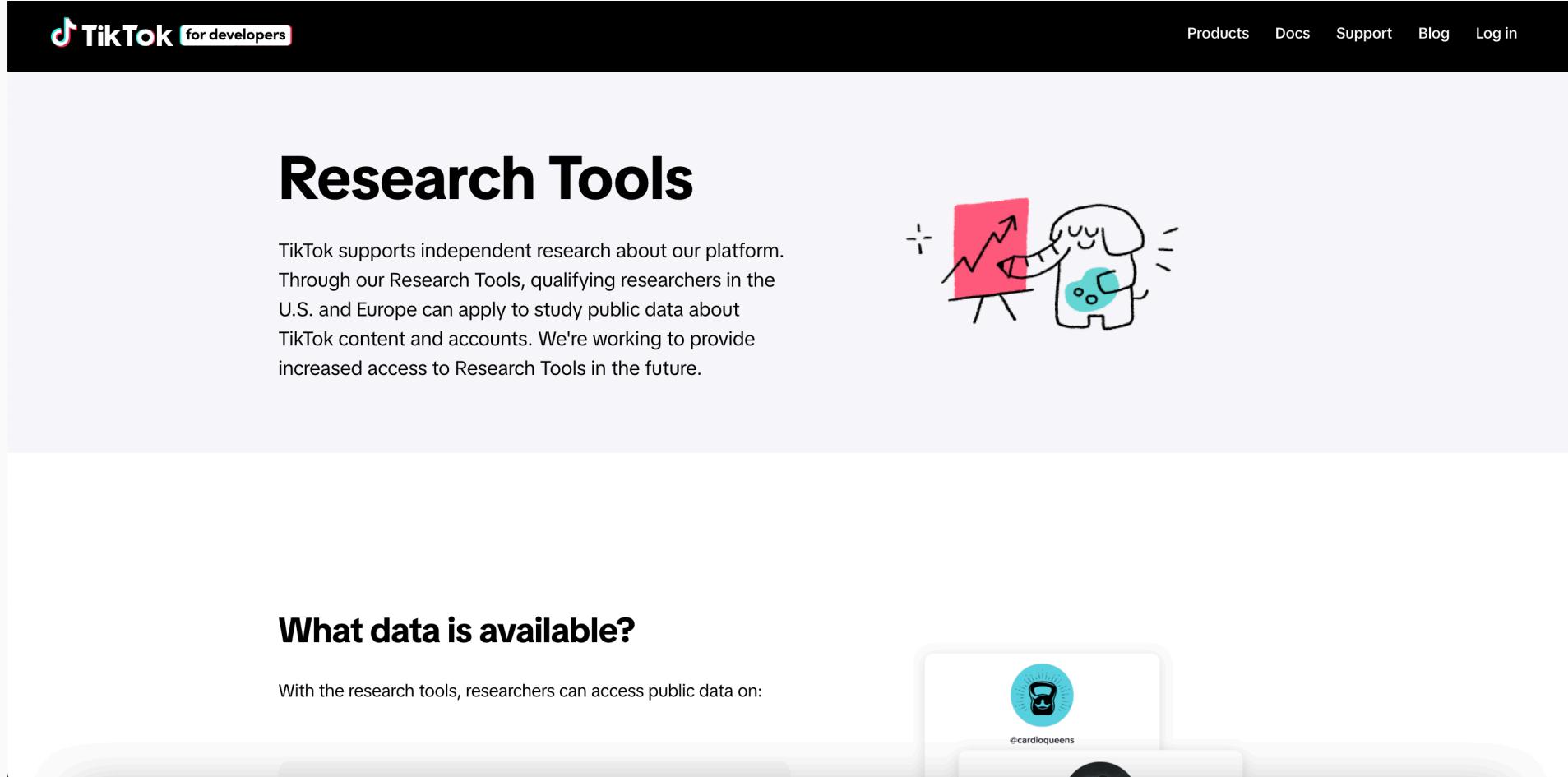




The screenshot shows the homepage of the Researcher API. At the top, there is a navigation bar with the Meta logo and links for "FORT", "Approach", "Researcher Products", and "Key Research Areas". Below the navigation bar, the text "Researcher API" is prominently displayed. A large image of a person with glasses is visible on the left side of the page. On the right side, there is a section titled "RESEARCHER API" with the subtext "One of the first Meta APIs designed for academics".

RESEARCHER API

One of the first Meta APIs designed for academics



The image shows a screenshot of the TikTok Research Tools landing page. At the top left is the TikTok logo with 'for developers' text. At the top right are links for 'Products', 'Docs', 'Support', 'Blog', and 'Log in'. The main title 'Research Tools' is in large bold black font. Below it is a paragraph of text: 'TikTok supports independent research about our platform. Through our Research Tools, qualifying researchers in the U.S. and Europe can apply to study public data about TikTok content and accounts. We're working to provide increased access to Research Tools in the future.' To the right of the text is a cartoon illustration of a white dog holding a red chart with an upward arrow. Below the main section is a heading 'What data is available?' and a subtext 'With the research tools, researchers can access public data on:'. At the bottom is a small image of a mobile device screen showing a profile picture of a person with a blue circle over their face, with the handle '@cardioqueens' below it.

TikTok for developers

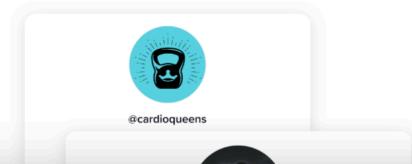
Products Docs Support Blog Log in

# Research Tools

TikTok supports independent research about our platform. Through our Research Tools, qualifying researchers in the U.S. and Europe can apply to study public data about TikTok content and accounts. We're working to provide increased access to Research Tools in the future.

## What data is available?

With the research tools, researchers can access public data on:



@cardioqueens

The screenshot shows the World Bank's website for 'Understanding Poverty / Research & Publications / Documents & Reports'. The main heading is 'Documents & Reports'. Below it, a section titled 'The World Bank Documents & Report API' is described. It mentions that the World Bank offers an API for searching and retrieving public documents. The 'Overview (A Simple Example)' section explains that the API uses URL query strings for HTTP GET requests. A list of items to be searched includes 'what is to be searched', 'what is to be returned for each matching record', and 'how what is returned is formatted (as JSON or XML)'. At the bottom, there is a note about a simple request for wind turbine documents and a 'Feedback Survey' button.

THE WORLD BANK  
IBRD • IDA

WHO WE ARE   WHAT WE DO   WHERE WE WORK   UNDERSTANDING POVERTY   WORK WITH US   WB Live  

Understanding Poverty / Research & Publications / Documents & Reports

## Documents & Reports

### The World Bank Documents & Report API

The World Bank offers an API that allows for the search and retrieval of the public, Bank documents available in the Documents & Reports site. Records can be retrieved in a format useful for research and for inclusion in web sites outside of Documents & Reports and the World Bank.

## Overview (A Simple Example)

The API is offered much like a REST API, accepting requests in the form a URL's query string as HTTP GET requests. The URL's query string determines:

- what is to be searched
- what is to be returned for each matching record
- how what is returned is formatted (as JSON or XML)

Feedback Survey

What follows is a simple request querying for records of documents related to wind turbines. It also requests that

# Questions?

---