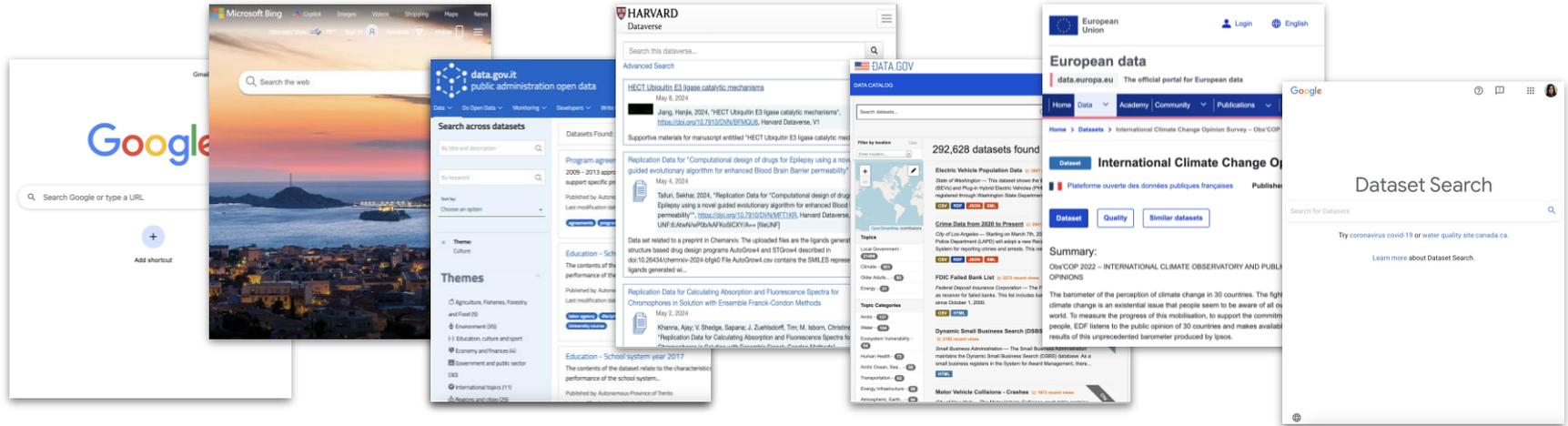# Discovering Datasets on the Web Scale: Challenges and Recommendations for Google Dataset Search

Katrina Sostek, Daniel M. Russell, Nitesh Goyal, Tarfah Alrashed, Stella Dugall, and Natasha Noy

# More Data on Web… More Dataset-Search Tools…
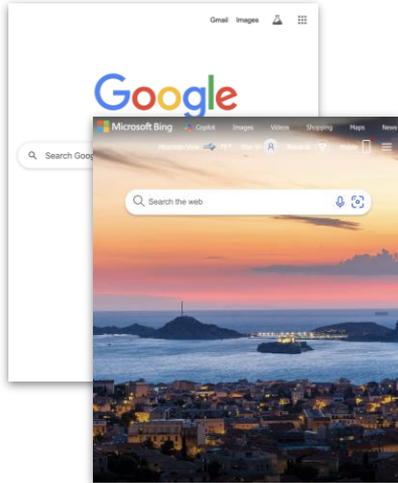
# The Scope of Dataset-Search Tools
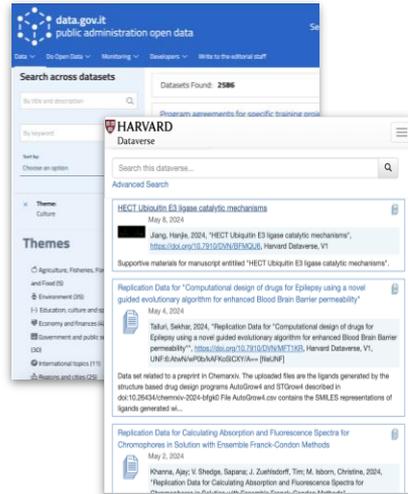
General Purpose Search Tools

Dataset-Specific Search Tools

General Purpose Web Search
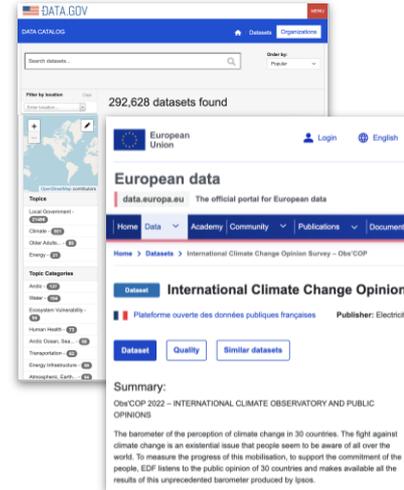
Dataset Repositories

Dataset Meta-Portals

Dataset-Specific Web Search



Google

# The Scope of Dataset-Search Tools

General Purpose Search Tools

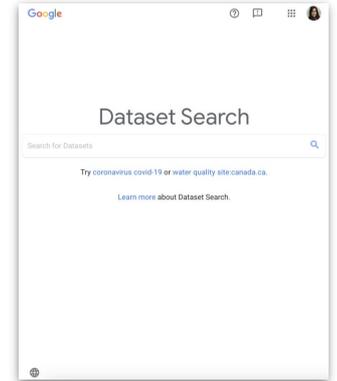Dataset-Specific Search Tools

General Purpose Web Search

Dataset Repositories

Dataset Meta-Portals

Dataset-Specific Web Search



Google

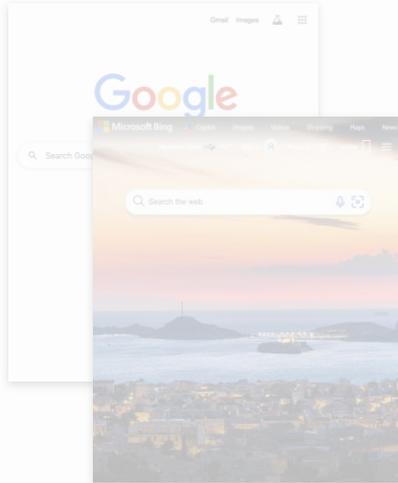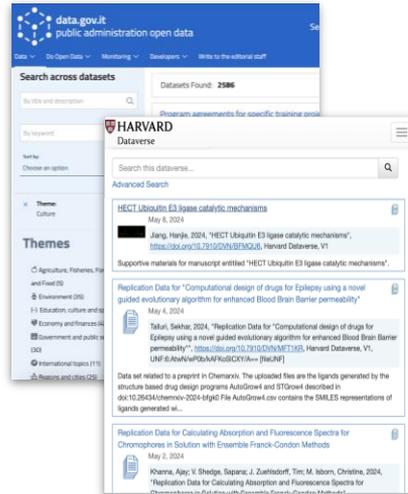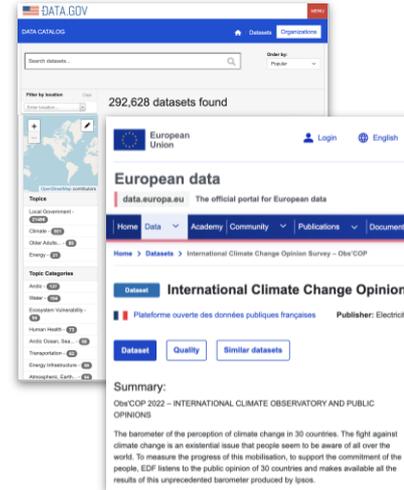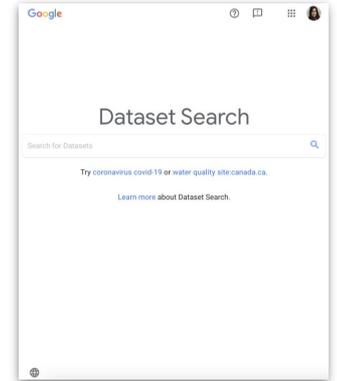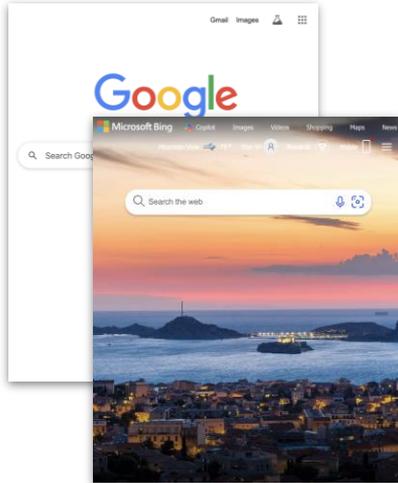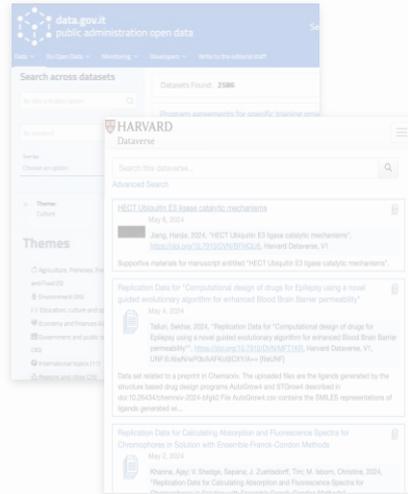# The Scope of Dataset-Search Tools

General Purpose Search Tools

Dataset-Specific Search Tools

General Purpose Web Search
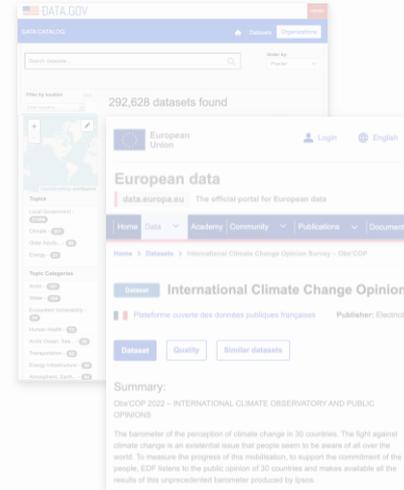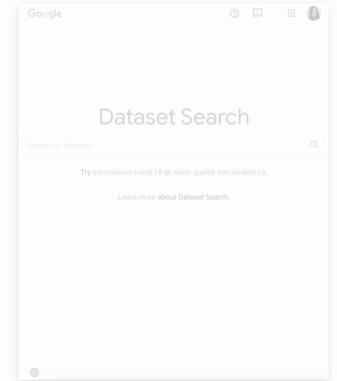
Dataset Repositories

Dataset Meta-Portals

Dataset-Specific Web Search

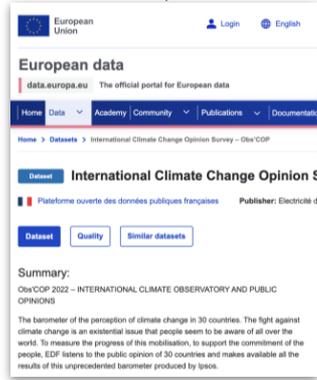# Previous Studies of Dataset-Discovery Tools

General Purpose Search Tools

Dataset-Specific Search Tools
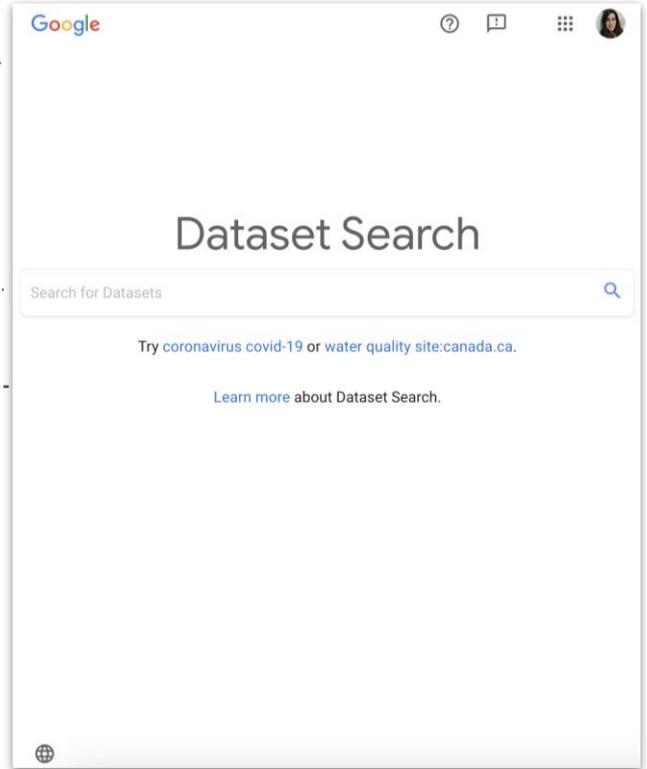
General Purpose Web Search

Dataset Repositories

Dataset Meta-Portals

Dataset-Specific Web Search

**User Interviews**

**User Surveys**

**Logs Analysis**

**Tool Reviews**

# The Scope of Google Dataset Search

Dataset Meta-Portals

Dataset Repositories

datasetsearch.research.google.com

Google

# Google Dataset Search

Research Questions

- How well does Dataset Search support user needs for dataset discovery as the only dataset-specific web search tool?

- What advantages and challenges do users face with Dataset Search due to its uniquely large scope and open approach?



Google

# User Interviews

## Participants

- We recruited 20 participants: 12 female, 10 male

- Seek datasets at least once a month

## Study Design

90-minute semi-structured virtual interviews.

- Background and motivations (~15 min)

- Recent dataset search challenges (~45 min)

- Dataset Search usage (~30 min)

10%

**Communication and Media**

Journalism

14%

**Business and Consulting**

Management Consulting, Market Research and Philanthropy

33%

**Science and Technology**

Technology, Biotechnology and Environmental Science

43%

**Health and Wellness**

Health/Medical Sciences, Psychology, Health Policy and Insurance
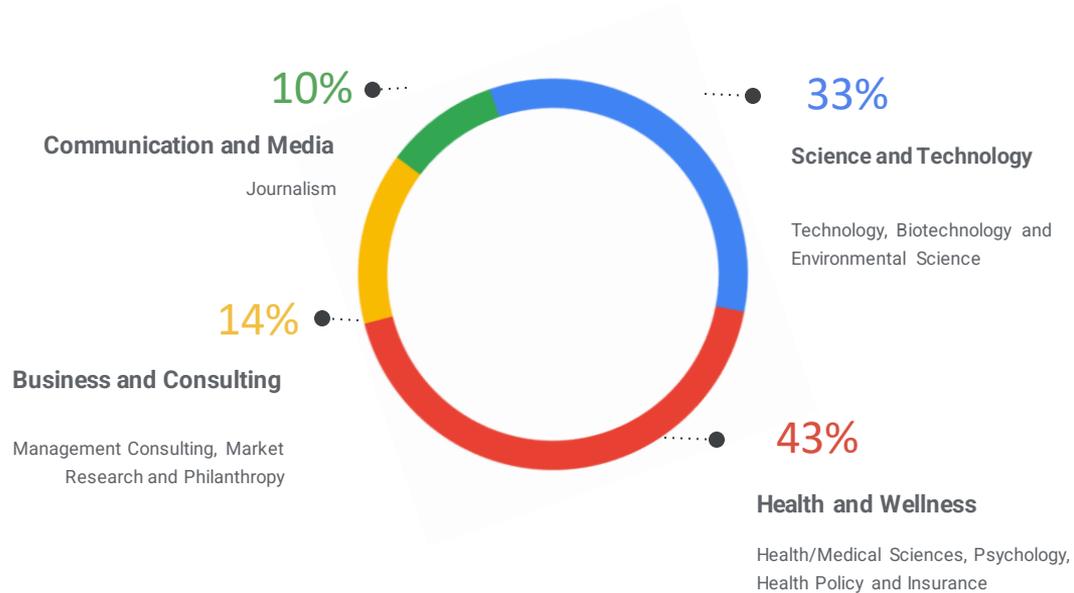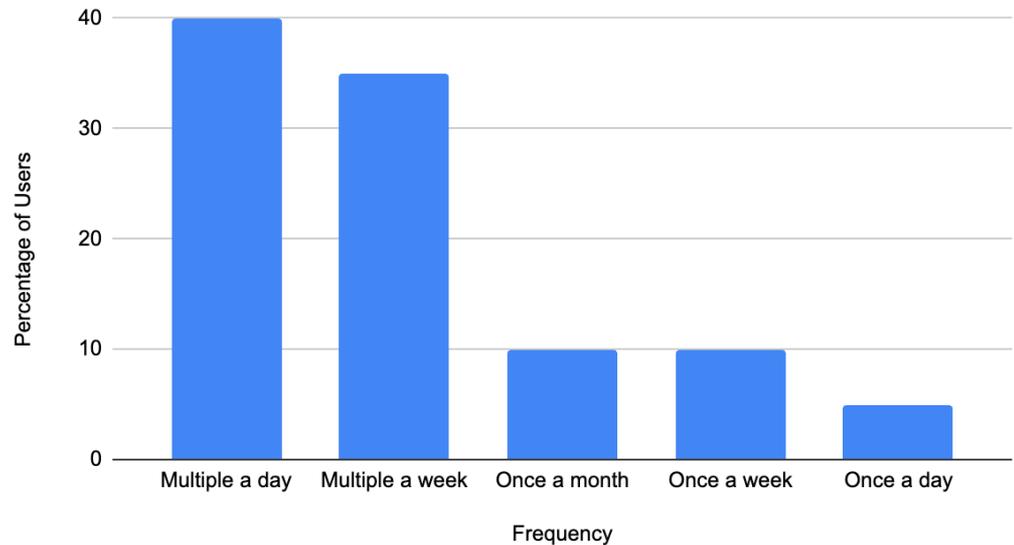
Google

# User Interviews

Participants

• We recruited 20 participants: 12 female, 10 male

• Seek datasets at least once a month

Study Design

90-minute semi-structured virtual interviews.

• Background and motivations (~15 min)

• Recent dataset search challenges (~45 min)

• Dataset Search usage (~30 min)

## How often do you search for datasets?



Google

# Findings

Our findings correspond to three themes:

1. Mental model development

Participants initially relied on familiar tool knowledge to understand the tool, updating their understanding as they used it. However, some initially confused its purpose with other tools.

2. Dealing with diverse datasets

Participants confirmed previous findings that incomplete and inconsistent metadata pose challenges, emphasizing the importance of accessing underlying data and ensuring trustworthiness. Dataset Search exacerbates these issues due to its wide range of sources and varying interfaces.

3. Data search proficiency

Some participants lacked familiarity with dataset-specific concepts, suggesting the need for tutorials and in-tool help to improve users' data search skills.

Google

# Building a Mental Model of the Tool

Challenge#1: Dataset Search and Google Web Search.

• Set user expectations that not all datasets on the web are in the tool.

• Fail gracefully when there are few results for a query.

• Compare mental model development in dataset web search between Google and standalone use.

Challenge#2: User expectations about the scope of the tool.

• Make it easier to find related artifacts in their tools.

• Make it easier to find datasets in the context of their search results.

• Infer the artifacts associated with datasets.



Results from Dataset Search in Google Web Search

Google

# Building a Mental Model of the Tool

Challenge#1: Dataset Search and Google Web Search.

- Set user expectations that not all datasets on the web are in the tool.

- Fail gracefully when there are few results for a query.

- Compare mental model development in dataset web search between Google and standalone use.

Challenge#2: User expectations about the scope of the tool.

- Make it easier to find related artifacts in their tools.

- Make it easier to find datasets in the context of their search results.

- Infer the artifacts associated with datasets.



Some datasets are linked to scientific articles in Google Scholar

# Building a Mental Model of the Tool

Challenge#3: Dataset replication across the web.

- Add indications in the UI explaining relationships between replicas

- Display dataset version types and sources in a clear, meaningful way.

- Infer the provenance of datasets

- Infer relationships between datasets



It also finds **versions** now!

Dataset Search finds **replicas** of datasets

Google

# Making Sense of Heterogeneous Datasets

Participants confirmed previous findings on dataset metadata limitations. While helpful for deciding on exploration, metadata lacked clarity for understanding datasets and determining usability. Dataset Search metadata lacks consistency due to diverse sources.

Challenge#1: Metadata shortcomings and inconsistencies

• Allow users to tailor their views and metadata selection based on their specific needs.

• Infer missing metadata fields.

Google

# Making Sense of Heterogeneous Datasets

Challenge#2: Downloading underlying data

- Add dataset download links, descriptions of content, previews, visualizations if not already supported.

- Infer formats and schema from underlying data, regardless of the accuracy of the dataset metadata.

Challenge#3: Building trust in an unfamiliar data sources.

- Provide signals about the trustworthiness or quality of data sources.

- Show more context about data sources to help users build trust.

- Enable users to customize views and results based on trust signals.

- Study how users use trustworthiness signals when datasets have incomplete metadata and artifacts.



Dataset Previews

# Learning How to Search for Data

Dataset searching requires data literacy and domain knowledge. Participants were unfamiliar with tool concepts like data formats and usage rights. Some discussed user learning improvements.

Dataset discovery is a skill and needs to be part of data literacy

- Create more educational materials and in-tool information about licenses, data formats, etc.

- Add educational information to help novice users, for example, specialized terminology, trustworthiness signals

- Develop educational content on web-scale dataset discovery for data literacy.
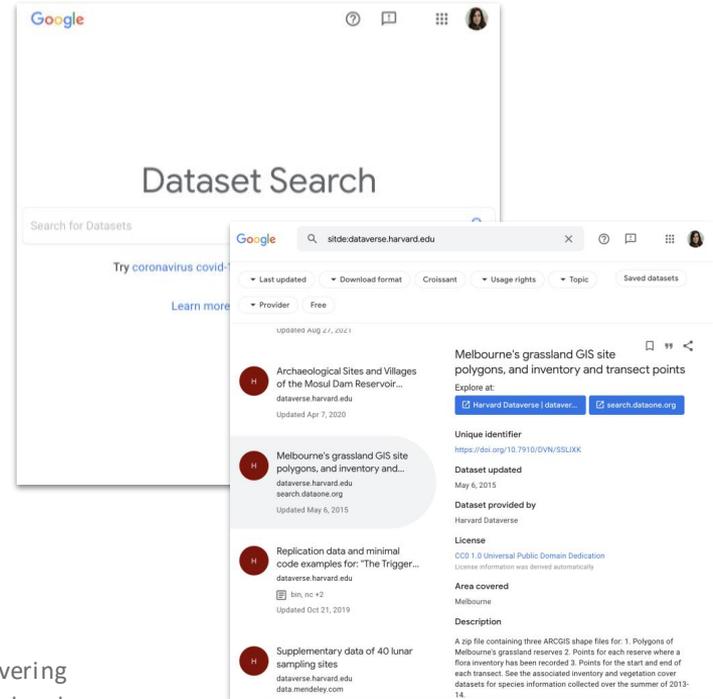


Google

# Limitations

• Dataset Search is the sole tool in its category, hindering broader comparison with similar tools.

• Participants varied in dataset discovery experience, shaping their interview approaches. However, common findings among participants enhanced the validity of our results.

• Participants used different queries rather than perform the same specific tasks.

• All participants were based in the US and used English-language queries.

• We did not have any participants who were regular users of Dataset Search.

Google

# Thank You

## Any Questions?

Sostek, K., Russell, D. M., Goyal, N., Alrashed, T., Dugall, S., & Noy, N. (2024). Discovering Datasets on the Web Scale: Challenges and Recommendations for Google Dataset Search. *Harvard Data Science Review*, (Special Issue 4). https://doi.org/10.1162/99608f92.4c3e11ca

datasetsearch.research.google.com

Google