# gesis

Leibniz Institute
for the Social Sciences

## Automatic Sampling and Analysis of YouTube Data

### Tools for collecting YouTube data

*Johannes Breuer, Annika Deubel, & M. Rohangis Mohseni*

*February 14th, 2023*

Leibniz
Association

# How to Collect *YouTube* Data

There are many different ways in which data from *YouTube* and other social media can be collected (see Breuer et al., 2020):

- Manually (e.g., via copy & paste and manual content analysis)

- Using existing data, such as *YouNiverse: Large-Scale Channel and Video Metadata from English YouTube* (also see the accompanying preprint by Ribeiro & West, 2021)

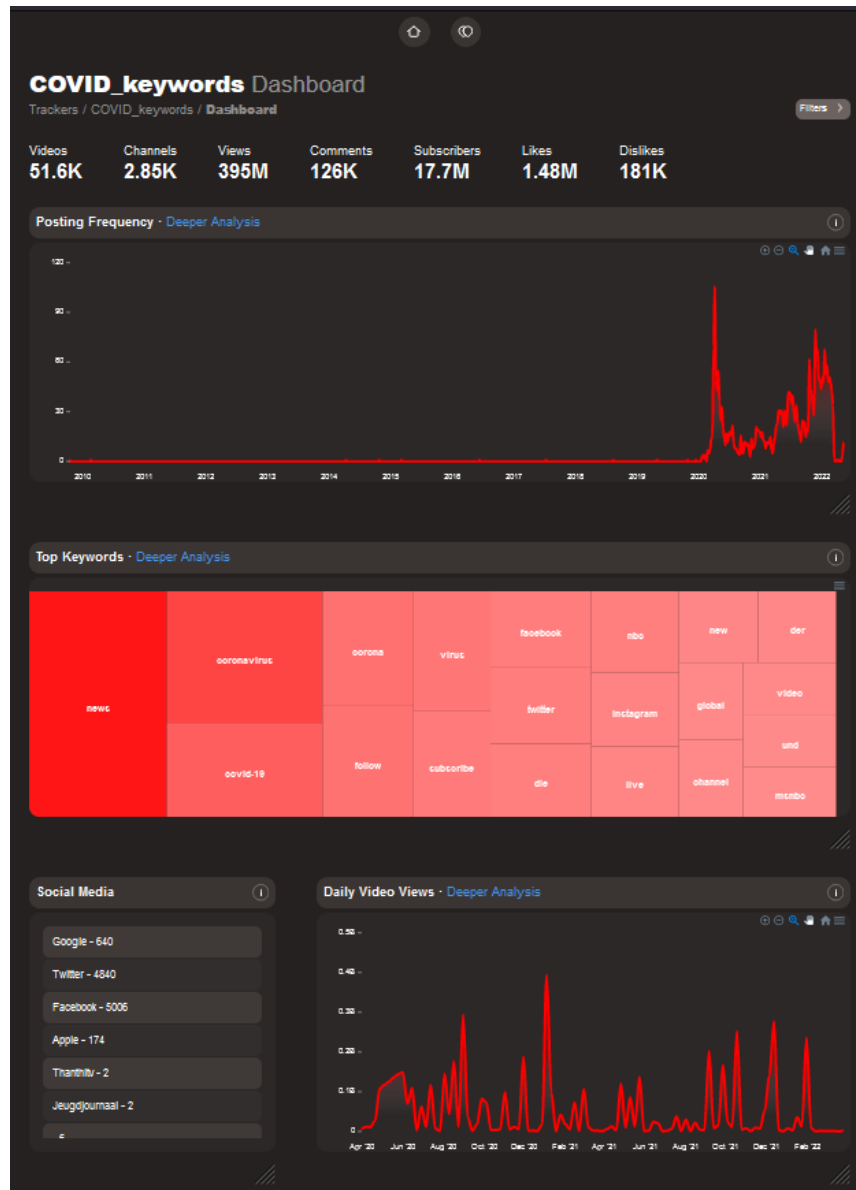- Automatically via the *YouTube* API or web scraping

# Identifying Relevant Channels or Videos

If new data is collected, it is necessary to identify relevant channels and videos for the sample.

- VTracker
- Socialblade
- YouTube Channel Crawler

# VTracker

- Search for and tracking of videos
- Low-key analysis such as engagement, keyword trends, influence detection
- Creation of Dashboard for different metrics
- Data can't be collected for further analysis
- Still a bit buggy

# Socialblade

- Ranked lists of channels
- Useful if there are no content-related criteria for channel selection

# YouTube Channel Crawler

- Search for channels with the help of filters (e.g. language, likes)
- Useful if there are no content-related criteria for channel selection

# Excluding Problematic Channels

- **YouTube Wiki**
  - Social background information on channels (only in German)
  - Useful to identify reasons for exclusion (e.g., fight between channels)

If the relevant channels are identified and potentially problematic channels are excluded, the next step would be to sample the comments.

Some of the comment sampling tools also offer search functions that can be used in addition to or instead of the tools mentioned above.

# Comparisons of Approaches for Collecting *YouTube* Data

| Software | Type | Can collect | Comment Scope | Needs API Key |
|---|---|---|---|---|
| YouTube Data Tools 1.22 | Website | Channel Info, Video Info, Comments | x top-level or all | No |
| Webometric 4.3 | Standalone app | Channel Info, Video Info, Comments, Video Search | 100 most recent or all | Yes |
| Tuber 0.9.9 | R package | Channel Info, Video Info, Comments, Subtitles, All searches | 20-100 most recent or all | Yes |
| vosonSML 0.29.13 | R package | Video IDs, Comments | 1-x top-level | Yes |
| youtubecaption 1.0.0 | R package | Subtitles | n/a | No |

gesis  Leibniz Institute for the Social Sciences

# YouTube Data Tools

## YouTube Data Tools

**YouTube Data Tools**

blog  software  research  DMI  about

Home | Channel Info | Channel List | Channel Network | Video List | Video Network | Video Comments | FAQ

**Video Info and Comments Module**

This module starts from a video id and retrieves basic info for the video in question and provides a number of analyses of the comment section. Comments are retrieved via the commentThreads/list API endpoint.

The number of comments the script is able to retrieve can vary wildly. In some cases, only a relatively small percentage is made available, while in others well over 100.000 comments have been successfully retrieved. This seems to be mainly related to the age of the video in question.

The module creates the following outputs:

- a tabular file containing basic info and statistics about the video;
- a tabular file containing all retrievable comments, both top level and replies;
- a tabular file containing comment authors and their comment count;
- a network file (gdf format) that maps interactions between users in the comment section;

The first three elements can be shown directly in the browser by enabling HTML output.

**Parameters**

**Video selection and comment cutoff:**

Video id:  _____  (video ids can be found in URLs, e.g. https://www.youtube.com/watch?v=**aXnaHh40xnM**)

Limit to:  _____  top level comments (ranked by relevance, leave empty for no limit)

**Output option:**

HTML output:  ☐ (displays HTML result tables in addition to file exports)

File format:  csv ⦿ / tab ○

# Webometric

## Webometric 4.3

# Exemplary Comparison of the Different Tools

| Software | Ease of Use | Disadvantages | No. of Comments |
|---|---|---|---|
| YouTube Data Tools 1.30 | High | Lacking flexibility, less information | 54,850 |
| Webometric 4.1 | Low | Only first 5 follow-up comments, no error feedback, undetectable time-outs | 51,095 |
| Tuber 0.9.9 | Low | Only first 5 follow-up comments | 51,084 |
| vosonSML 0.29.13 | Low | Lacking flexibility, only comments | 52,679 |

Example data source: Dayum Video

# A Note on Using FOSS

The tools listed are free and open source software (FOSS). Using FOSS has many advantages (availability, adaptability, etc.). However, one risk associated with using FOSS is that tools are not maintained anymore and cease to function. After all, people create and maintain these tools in their spare time or as side projects and this work is often not recognized enough (esp. within academia). For this reason it is important to acknowledge the work that goes into these tools by properly citing them.

```r
citation("tuber")
```

```
##
## To cite package 'tuber' in publications use:
##
##   Gaurav Sood (2020). tuber: Access YouTube from R. R package version 0.9.9.
##
## Ein BibTeX-Eintrag für LaTeX-Benutzer ist
##
##   @Manual{,
##     title = {tuber: Access YouTube from R},
##     author = {Gaurav SOod},
##     year = {2020},
##     note = {R package version 0.9.9},
##   }
```

# Any questions so far?