# Vluchtige tweets en honderden likes

Sociale media-archivering en het gebruik ervan door onderzoekers

**Anne Helmond** (Universiteit van Amsterdam).

'Webinar Sociale media-archivering?,' Netwerk Digitaal Erfgoed, 25 juni 2020.

# De onderzoeksmogelijkheden van archieven en platformen

# Introductie

- Archieven en archiveringspraktijken 'vormen' het corpus en beïnvloeden de soorten vragen die gesteld kunnen worden (zie Brügger, 2005; 2019).
- Sociale media platformen vormen ook het corpus door hun sterk gecontroleerde archiveringsmogelijkheden (zie Helmond en van der Vlist, 2019).
  - Technische controle: beperkte toegang tot een selectie van data
  - Contractuele controle: gebruik en hergebruik van data

# Vormen van webarchivering

1. Making an image
2. Making a screen movie
3. Downloading individual files
4. Web crawling
5. Collecting web material from a database via an API
6. Collecting the web that has been taken off-line and preserved unchanged
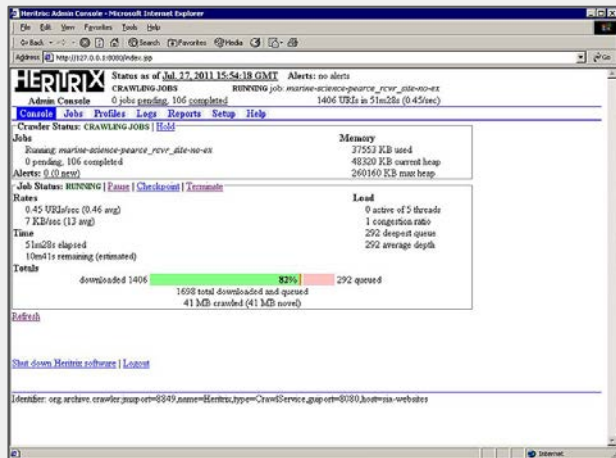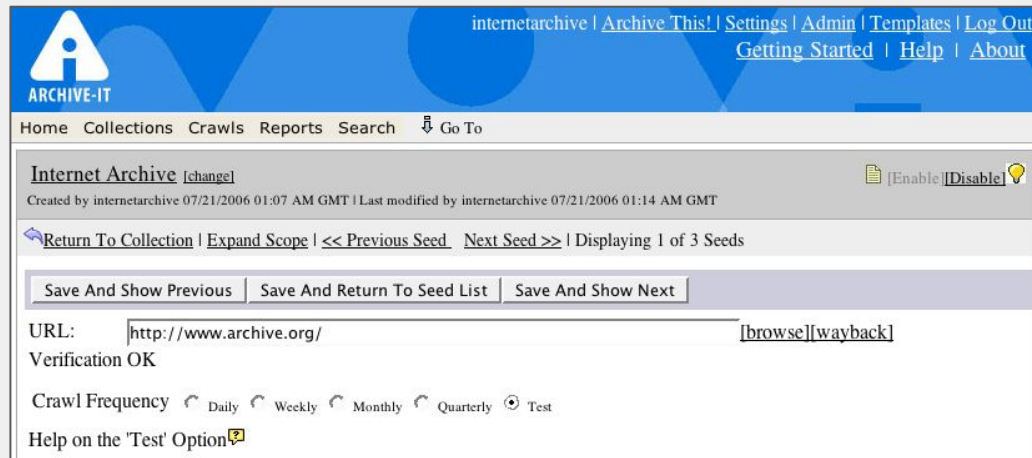7. Collecting the web as presented in other media types

(Brügger, 2018: 80).

# Web crawling

- **Crawl**: A web archiving (or "capture") operation that is conducted by an automated agent, called a crawler, a robot, or a spider. Crawls identify materials on the live web that belong in your collections, based upon your choice of seeds and scope. Crawl can also reference the archived content associated with the action.
- **Crawler**: Explores the web and collects data about its contents. A crawler can also be configured to capture web-based resources. It starts a capture process from a seed list of entry-point URLs (EPUs).
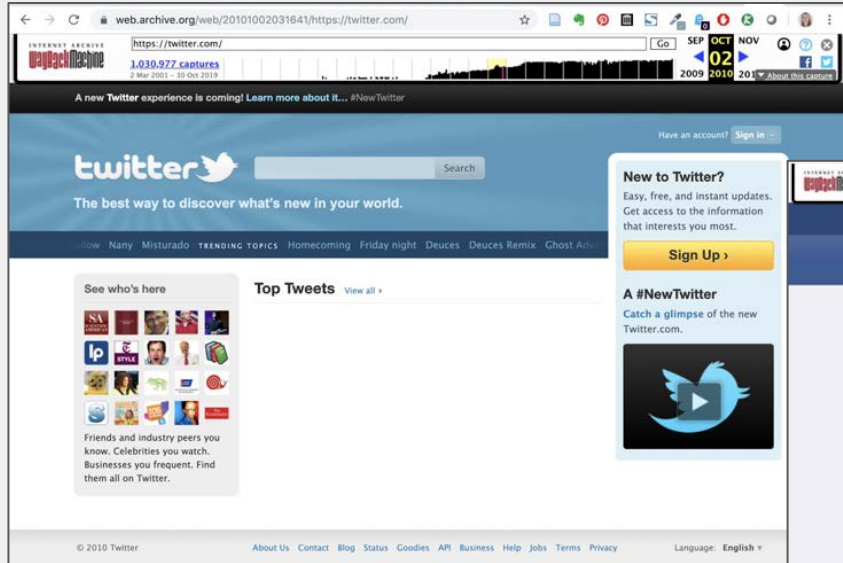
Glossary of Archive-It and Web Archiving Terms

# Het crawlen van sociale media



https://twitter.com/

https://facebook.com/

# Crawler profielen (bijv. Charlie Archivist)

# Tools voor sociale media archivering

Voor **kleine** individuele projecten:

Conifer to 'create high-fidelity, interactive captures'

- 👍Gratis account 5GB opslagruimte. Online interface. Opnemen en afspelen van websites. Standaard .warc file.
- 👎Files nemen veel ruimte in. 1 website tegelijk.

Webrecorder 'to capture and replay interactive websites'

- 👍Lokale versie van Conifer. Opnemen en afspelen van websites. Standaard .warc file.
- 👎Files nemen veel ruimte in. 1 website tegelijk.

# Sociale media archivering via crawlers

- De archivering van sociale media *via crawlers* levert een (leeg) **gebruikers blik** op het platform die verschillende vormen van **interface analyse** en **platform analyse** mogelijk maakt.
- De archivering van sociale media *via Conifer/Webrecorder* levert een **individuele gebruikers blik** die een **interface analyse** en **autobiografische analyse** mogelijk maakt.

# De vele gebruikers van sociale media



https://twitter.com/

http://business.twitter.com/

# Typen sociale media gebruikers en hun bronnen

**Table A-2.** User groups of social media platforms.

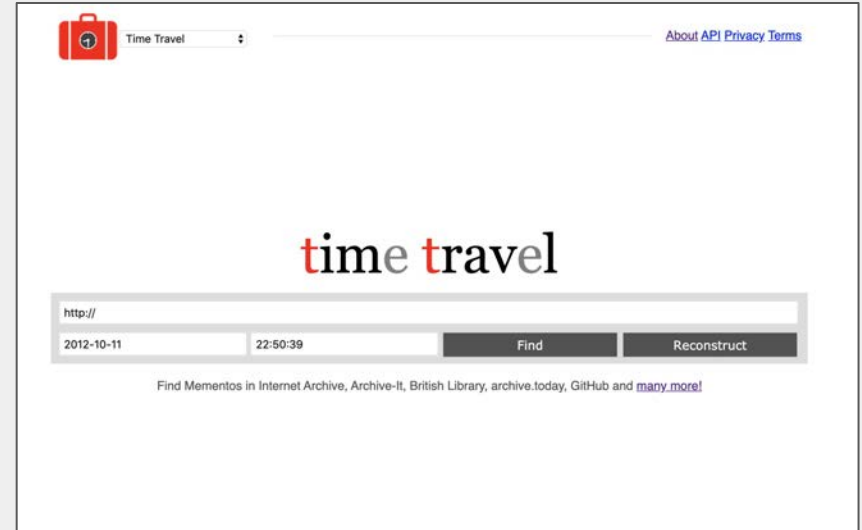| user group | archived resources (sample) | afforded histories (sample) |
|---|---|---|
| *end users* | <ul><li>graphical user interfaces and 'empty' frames</li><li>sign up and registration procedures</li><li>about and feature pages</li><li>data and privacy policies</li><li>terms of service and use</li><li>end-user license agreements (EULA)</li><li>account security pages</li><li>help pages</li><li>language support</li></ul> | <ul><li>self-description histories[17]</li><li>feature and practice histories[18]</li><li>'stakeholder politics' histories[19]</li><li>revenue model histories[20]</li><li>data and privacy policy histories</li><li>terms of service and use histories[21]</li></ul> |
| *developers* | <ul><li>tools and product pages</li><li>application programming interfaces (APIs) and endpoints</li><li>software development kits (SDKs)</li><li>integrated development environments (IDEs)</li><li>software and developer tools and frameworks</li><li>guides for app development, best practices, app review, and privacy and consent</li><li>online training courses for developers</li><li>developer support, help pages, and frequently asked questions (FAQs)</li><li>API reference documentation, version histories, and changelogs</li><li>developer news, blog posts, and blog archives</li><li>open source projects and code repositories</li><li>programming, query, and markup languages</li><li>bug reports</li><li>platform status</li><li>annual developer conferences</li><li>developers community groups, meetups and local developer communities</li><li>startup accelerator programmes</li><li>platform policies</li><li>careers</li><li>platform and privacy policies</li><li>cookies</li><li>terms of service and use</li></ul> | <ul><li>API-based data sharing histories[22]</li><li>data strategy and 'intraoperability' histories[23]</li><li>'programmability' and app development histories[24]</li><li>app ecosystem histories</li><li>tracking technology histories</li><li>standards and protocol histories</li><li>'datastructuring' histories[25]</li><li>platform architecture design and governance and control histories[26]</li><li>platform growth and embedding histories[27]</li><li>platform status, maintenance, repair, and 'issue' histories</li><li>platform and privacy policy histories[28]</li></ul> |
| *business* | <ul><li>ad creation and manag…</li><li>tools and product page…</li><li>business news, blog pages, and archives</li></ul> | <ul><li>ad creation and targeting fields histories</li><li>platform growth and embedding histories[29]</li></ul> |

Helmond A and van der Vlist FN (2019) Social Media and Platform Historiography: Challenges and Opportunities. *TMG – Journal for Media History* 22(1): 6–34. http://doi.org/10.18146/tmg.434
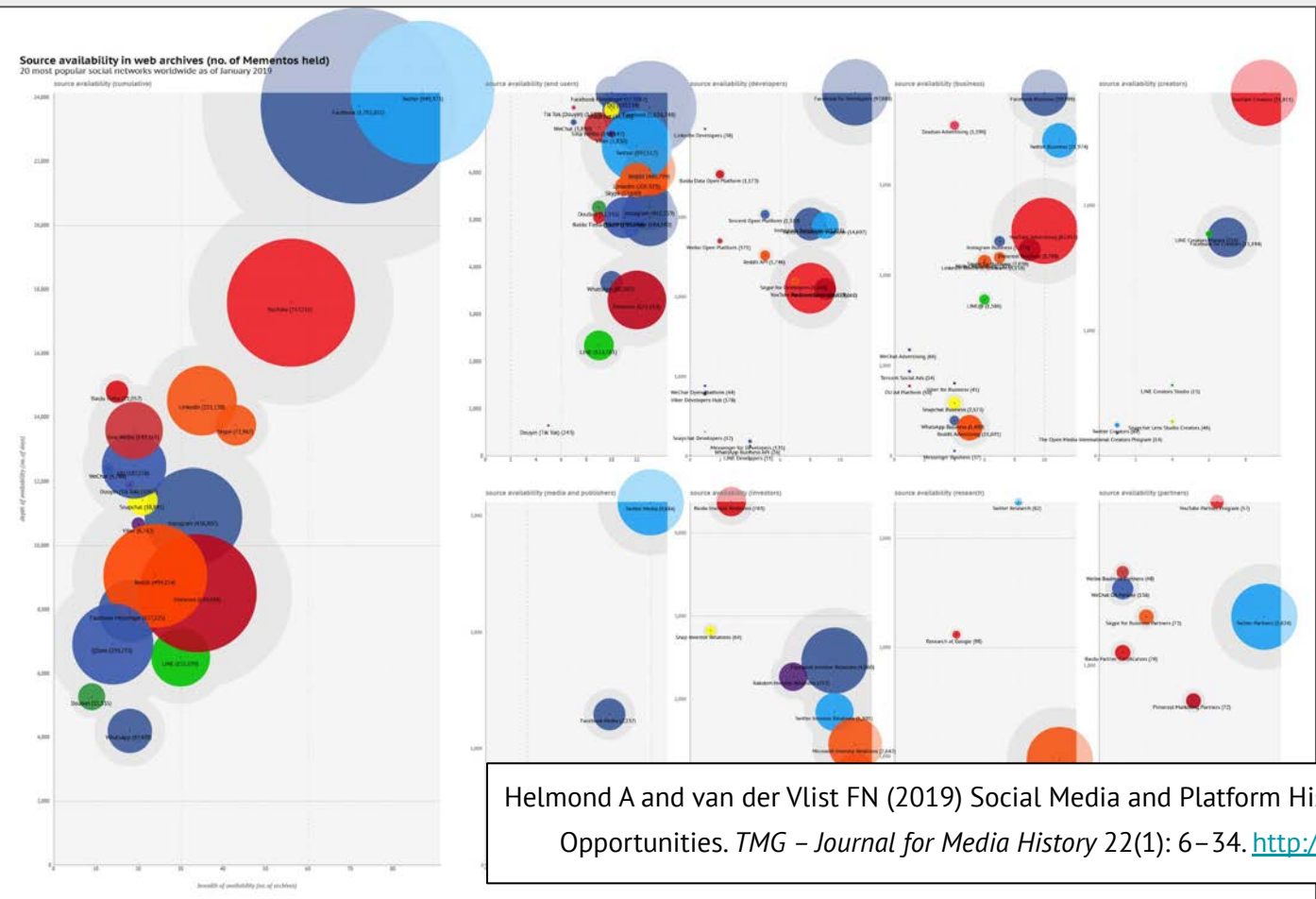
# Beschikbaarheid materialen: Memento & Memgator

- De beschikbaarheid van sociale media platform URLs voor alle gebruikersgroepen opvragen in 20+ web archieven
    - facebook.com
    - developers.facebook.com
    - facebook.com/business
    - facebook.com/politics
    - facebook.com/creators
    - etc
- Tools: Memento Time Travel API via Memgator (Alam & Nelson, 2016).



https://timetravel.mementoweb.org/

Helmond A and van der Vlist FN (2019) Social Media and Platform Historiography: Challenges and Opportunities. TMG – Journal for Media History 22(1): 6–34. http://doi.org/10.18146/tmg.434
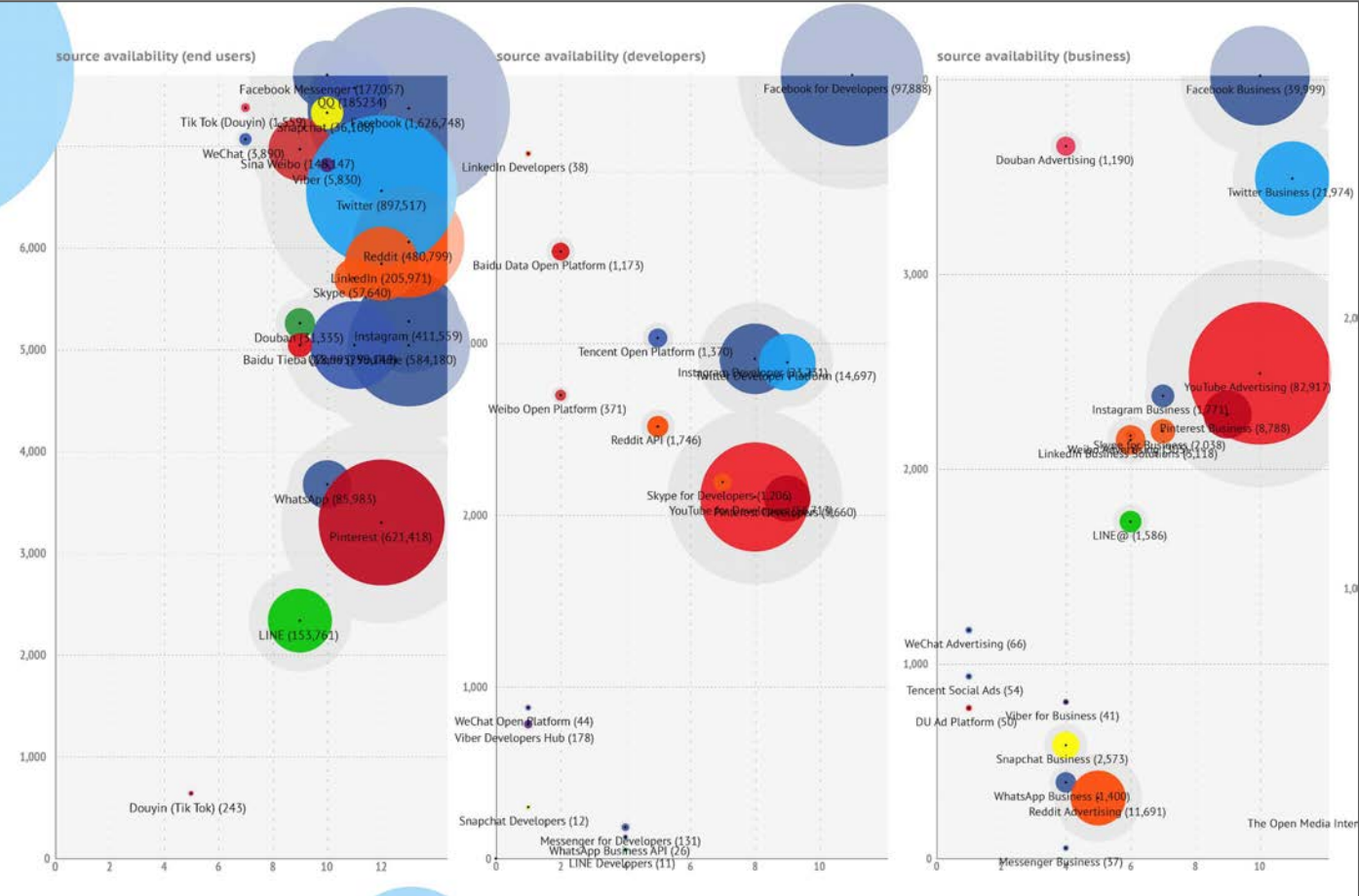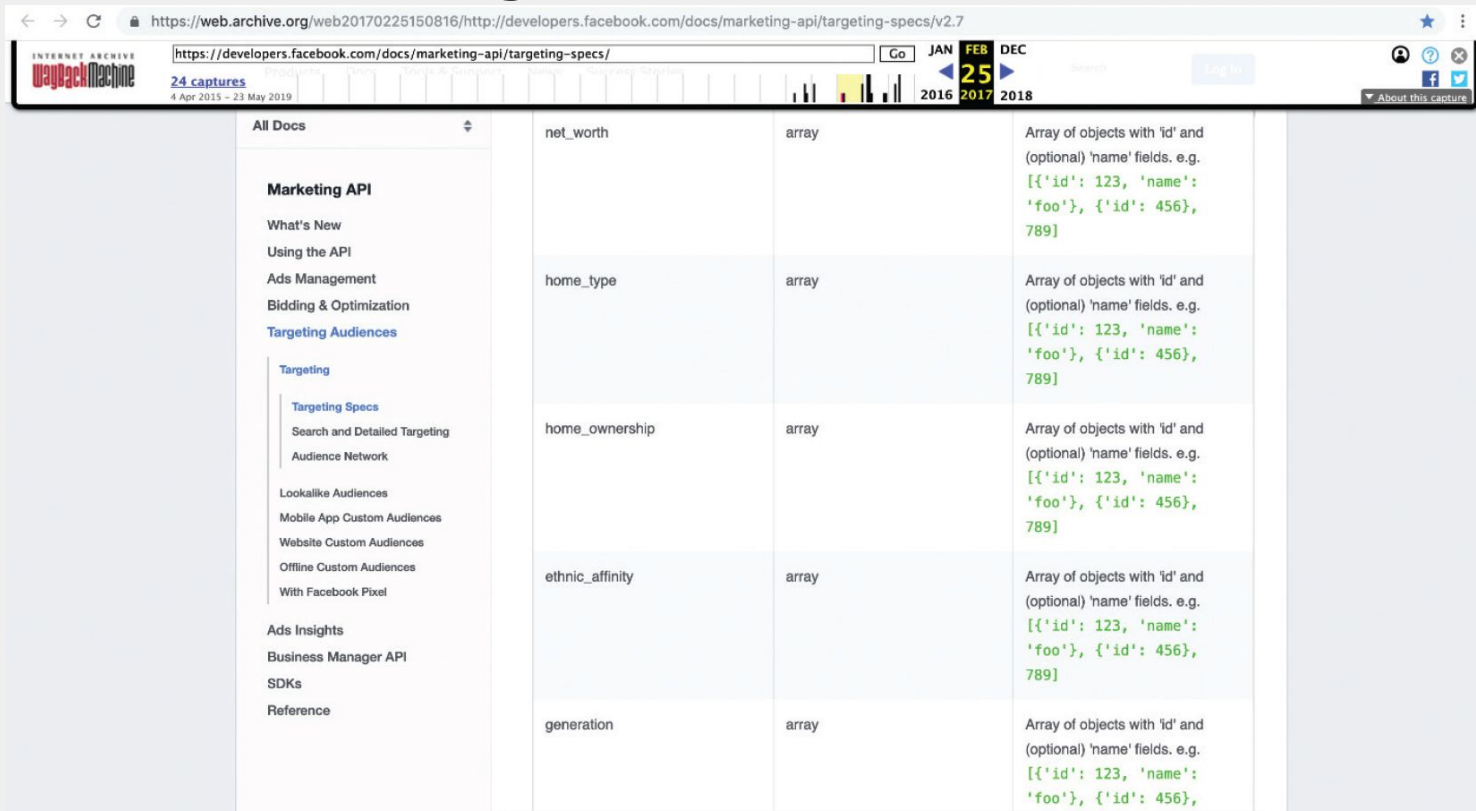
# Sociale media bronnen in web archieven



Helmond A and van der Vlist FN (2019) Social Media and Platform Historiography: Challenges and Opportunities. *TMG – Journal for Media History* 22(1): 6–34. http://doi.org/10.18146/tmg.434

# Sociale media bronnen in web archieven

# Advertentie doelgroepen van Facebook



Screen capture of Facebook Marketing API audience targeting specifications at developers.facebook.com/docs/marketing-api/targeting-specs/, including the 'ethnic_affinity' field (now disabled). *Internet Archive Wayback Machine, February 25, 2017* (Helmond & van der Vlist, 2019).

# De partners van Facebook

# Officiële Facebook partners, 2010–2017

bit.ly/2oYvpI5

Zoom: 25%



| 2010 Dec 05 | 2011 Aug 13 | 2012 Jan 03 | 2012 May 12 | 2013 Jan 18 | 2014 Feb 08 | 2014 Nov 05 | 2016 Mar 12 | 2017 Jan 30 | 2017 Apr 07 |
|---|---|---|---|---|---|---|---|---|---|
| Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks | Brand Networks |
| Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems | Hearsay Systems |
| Komfo | Komfo | Komfo | Komfo | Komfo | Komfo | Komfo | Komfo | Komfo | Komfo |
| Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers | Socialbakers |
| Zibaba | Zibaba | Zibaba | Zibaba | Zibaba | Zibaba | Zibaba | Zibaba | Zibaba | Zibaba |
| Adobe | Conversocial | Conversocial | Conversocial | Conversocial | Conversocial | Conversocial | Conversocial | Conversocial | Conversocial |
| KRDS | Fanbooster | Fanbooster | Fanbooster | Fanbooster | Fanbooster | Fanbooster | Fanbooster | Fanbooster | Fanbooster |
| Kenshoo | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media | Hootsuite Media |
| Kremsa Digital | Shoutlet | Shoutlet | Shoutlet | Shoutlet | Shoutlet | Shoutlet | Shoutlet | Shoutlet | Shoutlet |
| Marin Software | StitcherAds | StitcherAds | StitcherAds | StitcherAds | StitcherAds | StitcherAds | StitcherAds | StitcherAds | StitcherAds |
| Sprout | KRDS | KRDS | Adobe | Adobe | Adobe | Adobe | Adobe | Adobe | Adobe |
| Alphabet | Kremsa Digital | Kremsa Digital | KRDS | KRDS | KRDS | KRDS | KRDS | Kenshoo | Kenshoo |
| Blueye Creative | Sprout | Alphabet | Kenshoo | Kenshoo | Kenshoo | Kenshoo | Kenshoo | Marin Software | Marin Software |
| Carrot Creative | Alphabet | Blueye Creative | Kremsa Digital | Kremsa Digital | Kremsa Digital | Kremsa Digital | Kremsa Digital | Sprout | Sprout |
| Dachis Group | Blueye Creative | Carrot Creative | Marin Software | Marin Software | Marin Software | Marin Software | Marin Software | 4C Insights | 4C Insights |
| Experian | Carrot Creative | Dachis Group | 4C Insights | Sprout | Sprout | Sprout | Sprout | AdParlor | AdParlor |
| Fan Appz | Dachis Group | Dentsu | AdParlor | 4C Insights | 4C Insights | 4C Insights | 4C Insights | Adaptly | Adaptly |
| Fluid | Dentsu | Fan Appz | Adaptly | AdParlor | AdParlor | AdParlor | AdParlor | Computerology | Computerology |
| Friend2Friend | Fan Appz | Fluid | Alphabet | Adaptly | Adaptly | Adaptly | Adaptly | Die Socialisten | Die Socialisten |
| Promoqube | Fluid | Friend2Friend | Blueye Creative | Alphabet | Alphabet | Alphabet | Computerology | Experian | Experian |
| SocialAmp | Friend2Friend | MakeMeReach | Carrot Creative | Blueye Creative | Blueye Creative | Blueye Creative | Dentsu | HYFN | HYFN |
| Thuzi | MakeMeReach | Promoqube | Computerology | Carrot Creative | Carrot Creative | Carrot Creative | Die Socialisten | Innobirds Media | Innobirds Media |
| 77Agency | Promoqube | SocialAmp | Dachis Group | Computerology | Computerology | Computerology | Experian | Kinetic Social | Kinetic Social |
| Gamaroff Digital | SocialAmp | Thuzi | Dentsu | Dachis Group | Dachis Group | Dachis Group | HYFN | Nanigans | MakeMeReach |
| Glow Digital Medi | Thuzi | Tigerlily | Die Socialisten | Dentsu | Dentsu | Dentsu | Innobirds Media | Napoleon | Nanigans |
| Syncapse | Tigerlily | Gamaroff Digital | Fan Appz | Die Socialisten | Die Socialisten | Die Socialisten | Kinetic Social | Qwaya | Napoleon |
| BLiNQ M | Fission | | MakeMeReach | | | | | | |
| Likeabl | | | | | | | | | |
| Sociabl | | | | | | | | | |
| Votigo | | | | | | | | | |
| Buddy Media | Strutta | Syncapse | MakeMeReach | Innobirds Media | Innobirds Media | Innobirds Media | Spredfast | AdRoll | Unified |

# Sociale media platforms: GUI versus API

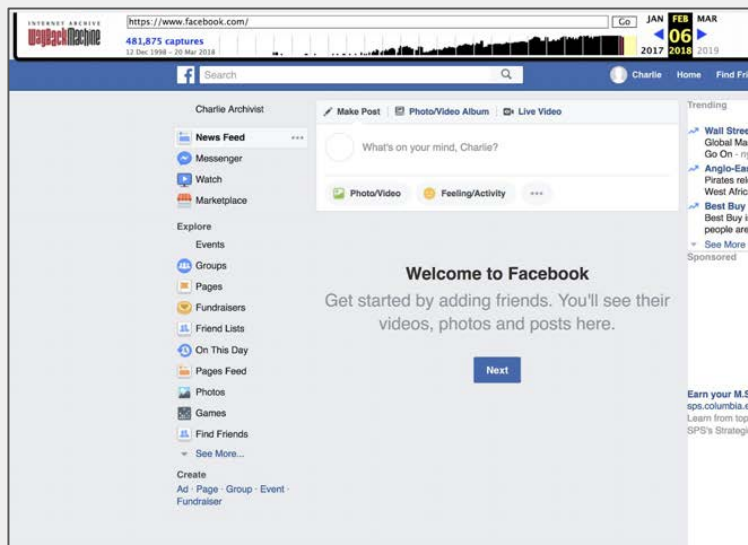Graphical User Interface (GUI):

- Archiveringsmethode: crawler
- Toegang: via login
- De "voorkant" van sociale media, toegang tot de gebruikersinterface.
- Gebruikersgroepen: eindgebruikers, maar ook politici, adverteerders, etc.
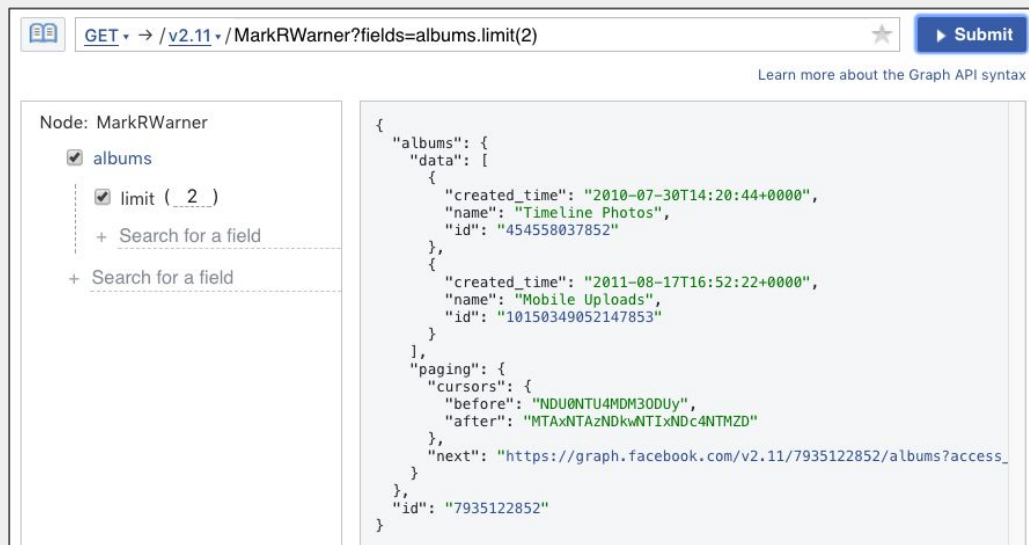
Application Programming Interface (API):

- Archiveringsmethode: API
- Toegang: via developer account
- De "achterkant" van sociale media, toegang tot de data in de database.
- Gebruikersgroepen: (app) ontwikkelaars, adverteerders en marketing bedrijven.

# Sociale media archivering: GUI/API

**Crawler-methode**

**API-methode**



**GUI ≠ API: Niet alle data die je in de GUI ziet kun je via de API krijgen!**

# Sociale media archivering via APIs

- De archivering van sociale media via APIs levert een **data-blik** op het platform die verschillende vormen van **data-analyse** mogelijk maakt.

- Onderzoekers maken zelf datasets of gebruiken bestaande datasets.
  - Bestaande datasets gebruiken:
    - 👍Iemand anders heeft al waardevolle historische data verzameld
    - 👎Lastig te vinden en niet altijd duidelijk gedocumenteerd
  - Zelf datasets maken:
    - 👍De onderzoeksvraag bepaalt welke data je nodig hebt
    - 👎Tools, skills en (veel) opslagruimte nodig

# Bestaande datasets gebruiken

- Opslag en ontsluiting van sociale media data is technisch lastig en kan zeer veel ruimte kosten.
  - 'The Twitter Archive at the Library of Congress: Challenges for information practice and information policy' (Zimmer, 2015).
  - 'The Library of Congress Twitter Archive: A Failure of Historic Proportions' (Bruns, 2018).
- Zeer sterke restricties voor het maken, delen en gebruiken van datasets.
  - Veranderende terms of service.
  - Platformen stellen steeds minder data beschikbaar 'APIcalypse' (Bruns, 2019).
  - Twitter: datasets herpubliceren alleen op basis van tweet-ID.
  - Facebook: Social Science One (partnerovereenkomst)
  - Historische data als verdienmodel (Twitter).

# Bestaande Twitter datasets ([DocNow Catalog](#))
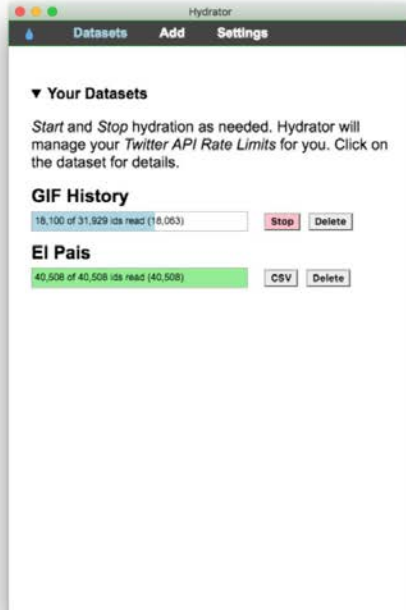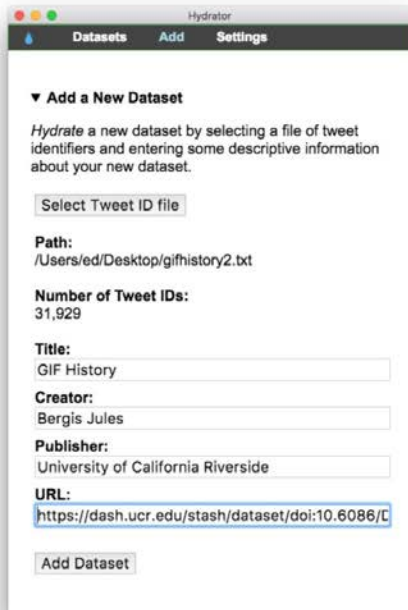
# Van tweet id naar volledige tweet

`"id": 1050118621198921728`

```
{
  "created_at": "Wed Oct 10 20:19:24 +0000 2018",
  "id": 1050118621198921728,
  "id_str": "1050118621198921728",
  "text": "To make room for more expression, we will now count all emojis as equal—including those w
and skin t… https://t.co/MkGjXf9aXm",
  "truncated": true,
  "entities": {
    "hashtags": [],
    "symbols": [],
    "user_mentions": [],
    "urls": [
      {
        "url": "https://t.co/MkGjXf9aXm",
        "expanded_url": "https://twitter.com/i/web/status/1050118621198921728",
        "display_url": "twitter.com/i/web/status/1…",
        "indices": [
          117,
          140
        ]
      }
    ]
  },
  "source": "<a href=\"http://twitter.com\" rel=\"nofollow\">Twitter Web Client</a>",
  "in_reply_to_status_id": null,
  "in_reply_to_status_id_str": null,
  "in_reply_to_user_id": null,
  "in_reply_to_user_id_str": null,
  "in_reply_to_screen_name": null,
```

```
  "user": {
    "id": 6253282,
    "id_str": "6253282",
    "name": "Twitter API",
    "screen_name": "TwitterAPI",
    "location": "San Francisco, CA",
    "description": "The Real Twitter API. Tweets about API changes, service issues and our Developer Platform.
Don't get an answer? It's on my website.",
    "url": "https://t.co/8IkCzCDr19",
    "entities": {
      "url": {
        "urls": [
          {
            "url": "https://t.co/8IkCzCDr19",
            "expanded_url": "https://developer.twitter.com",
            "display_url": "developer.twitter.com",
            "indices": [
              0,
              23
            ]
          }
        ]
      },
      "description": {
        "urls": []
      }
    },
    "protected": false,
    "followers_count": 6128663,
    "friends_count": 12,
    "listed_count": 12900,
    "created_at": "Wed May 23 06:01:13 +0000 2007",
    "favourites_count": 32,
    "utc_offset": null,
    "time_zone": null,
    "geo_enabled": null,
    "verified": true,
    "statuses_count": 3659,
```

# DocNow Twitter [Hydrator](#)



Hydrator is an Electron based desktop application for hydrating Twitter ID datasets. Twitter's Terms of Service do not allow the full JSON for datasets of tweets to be distributed to third parties. However they do allow datasets of tweet IDs to be shared. Hydrator helps you turn these tweet IDs back into JSON and also CSV from the comfort of your desktop.

# Zelf data verzamelen: Twitter Search API

## Feature summary

| Category | Product name | Supported history | Query capability | Counts endpoint | Data fidelity |
|---|---|---|---|---|---|
| Standard | Standard Search API | 7 days | Standard operators | Not available | Incomplete |
| Premium | Search Tweets: 30-day endpoint | 30 days | Premium operators | Available | Full |
| Premium | Search Tweets: Full-archive endpoint | Tweets from as early as 2006 | Premium operators | Available | Full |
| Enterprise | 30-day Search API | 30 days | Premium operators | Included | Full |
| Enterprise | Full-archive Search API | Tweets from as early as 2006 | Premium operators | Included | Full |

https://developer.twitter.com/en/docs/tweets/search/overview

# Tools voor sociale media archivering (Twitter 1)

Twitter Archiving Google Sheet (TAGS): [https://tags.hawksey.info/](https://tags.hawksey.info/)

- 👍Gratis, veel mogelijkheden om Twitter data te archiveren via de Twitter API. Gekoppeld aan Google Sheets. Enkele ingebouwde analyse mogelijkheden voor onderzoekers.
- 👎Alleen tweets van de laatste 7 dagen (limitatie Twitter API)

# Tools voor sociale media archivering (Twitter 2)

The Digital Methods Initiative Twitter Capture and Analysis Toolset (DMI-TCAT) allows one to retrieve and collect tweets from Twitter and to analyze them in various ways: https://github.com/digitalmethodsinitiative/dmi-tcat

- 👍Gratis, open source. Zeer veel mogelijkheden om Twitter data te archiveren via de Twitter API. Zeer veel ingebouwde analyse mogelijkheden voor onderzoekers (zie instructie video).
- 👎Installeren op eigen server (ook een voordeel).

## Analyses

- **#hashtags**: sociale en culturele issues, bijvoorbeeld #BlackLivesMatterNL, #coronamaatregelen
- **mentions**: identificeren van expertise, discussie netwerken
- **following**: volgnetwerken
- **retweets**: resonantie, viraliteit
- **URL**: distributie van content

# De rol van web archieven en instanties

# Tools voor archivarissen

[DocNow](#) is an **appraisal tool** for the web.

"DocNow allows archivists to tap into conversations in Twitter to help them discover what web resources are in need of archiving. Its goal is to help ensure **ethical practices** in web archiving by building conversations between archivists and the communities they are documenting."

# DocNow ([demo](#))



72 #sayhername

65 #justiceforbreonnataylor

31 #corneliusfredericks

22 #elijahmcclain

21 #blm

14 #antifia

14 #breonnataylor

12 #justiceforelijah

8 #retweet

7 #racism

## HASHTAGS

19 VIDEOS

VIEW VIDEO INSIGHTS →

959 USERS

BLACK MARKET

I ♥ USA

IN DEFENSE BLACK LIFE

VIEW USER INSIGHTS →

# Doel van de collectie

- Keuzes over selectie bepalen later welke vragen onderzoekers kunnen stellen.
- Wie krijgt er toegang tot de collectie en hoe en waar?
  - Onderzoekers werken het liefst met 'rauwe data' op hun eigen computer.
  - Maar: privacy en ethiekvraagstukken en platform regels.

# Collectiebeschrijving en metadata

- Essentieel voor onderzoekers is informatie over:
  - Hoe is de collectie gemaakt?
  - Welke beslissingen zijn er genomen?
  - Wat zijn de selectiecriteria?
  - Welke data zit er in de set?
  - Hoe is deze data verzameld (en met welke API)?
- "Web archives **provenance**: what users need to know about how a collection was made as they use, analyze, and make inferences from these aggregations" (Maemura et al., 2018).
- Een collectiebeschrijving: Metadata over het gearchiveerde corpus.
- Rauwe data, bijvoorbeeld in .csv of .json (.warc is lastig om mee te werken)

# Voorbeeld collectiebeschrijving (1)

- Bode, P. de, Lin, K., Teszelszky, K. (2019) Chinese Netherlands web collection. KB Lab: The Hague. https://lab.kb.nl/dataset/web-collection-chinese-netherlands

Information about the collection and its heritage value can be found in a 📥 collection description (in English) a collection description in 📥 Traditional Chinese and one in 📥 Simplified Chinese. The collection can be studied on the terminals in the reading room of KB with a valid library card. Researches can also use the 📥 dataset with URL's (note: UTF-8 encoding) and a 📥 link analysis.

# Voorbeeld collectiebeschrijving (2)

● Internet Archive Wayback Machine 'About this capture'

# Voorbeeld collectiebeschrijving (3)

- The DocNow Catalog



**Catalog**

**Title:**
Twitter Historical Dataset
**Repository:**
Zenodo
**Repository URL:**
https://doi.org/10.5281/zenodo.3833781
**Creator(s):**

Daniel Gayo-Avello

**Subjects:**

Social Media
History
Twitter
Politics
Languages

**Dates:**

2006-03-21T17:04:15.796Z - 2009-07-31T16:04:15.808Z

**Number of Tweets:**
1,499,896,115

**Description:**

This dataset is distributed by Daniel Gayo-Avello, an associate professor at the Department of Computer Science in the University of Oviedo, for the sole purpose of non-commercial research and it just includes tweet ids.

The dataset contains tweet IDs for all the published tweets (in any language) between March 21, 2006 and July 31, 2009 thus comprising the first whole three years of Twitter from its creation, that is, about 1.5 billion tweets (see file *Twitter-historical-20060321-20090731.zip*).

It covers several defining issues in Twitter, such as the invention of hashtags, retweets and trending topics, and it includes tweets related to the 2008 US Presidential Elections, the first Obama's inauguration speech or the 2009 Iran Election protests (one of the so-called Twitter Revolutions).

Finally, it does contain tweets in many major languages (mainly English, Portuguese, Japanese, Spanish, German and French) so it should be possible—at least in theory—to analyze international events from different cultural perspectives.

The dataset was completed in November 2016 and, therefore, the tweet IDs it contains were publicly available at that moment. This means that there could be tweets public during that period that do not appear in the dataset and also that a substantial part of tweets in the dataset has been deleted (or locked) since 2016.

To make easier to understand the decay of tweet IDs in the dataset a number of representative samples (99% confidence level and 0.5 confidence interval) are provided.

# Conclusie

- Archieven, archiveringspraktijken en sociale media platformen 'vormen' het corpus en beïnvloeden de soorten vragen die gesteld kunnen worden.
- De onderzoeksvraag is leidend voor welke archief of welke data je nodig hebt.
- Het nauwkeurig documenteren van een collectie is essentieel voor onderzoekers om gebruik te maken van de gearchiveerde sociale media data.
- Voor web archieven, instellingen en onderzoekers is API kennis belangrijk.

# Bedankt! Vragen?

Helmond A and van der Vlist FN (2019) Social Media and Platform Historiography: Challenges and Opportunities. *TMG – Journal for Media History* 22(1): 6–34. doi:10.18146/tmg.434

Helmond A, Nieborg DB and van der Vlist FN (2019) Facebook's evolution: Development of a platform-as-infrastructure. *Internet Histories: Digital Technology, Culture and Society*, 3(2), 123–146. doi: 10.1145/3097286.3097324.

Anne Helmond (Universiteit van Amsterdam)
**a.helmond@uva.nl**

# Bronnen en links

- Brügger N (2018a) *The Archived Web: Doing History in the Digital Age*. Cambridge, MA: The MIT Press.
- Brügger N (2018b) *Web History and Social Media*. In: Burgess J, Poell T, and Marwick A (eds) The SAGE Handbook of Social Media. London: SAGE Publications, pp. 196–212.
- Bruns A (2019) After the 'APIcalypse': social media platforms and their fight against critical scholarly research. *Information, Communication & Society* 22(11). Routledge: 1544–1566. DOI: 10.1080/1369118X.2019.1637447.
- Helmond A and van der Vlist FN (2019) Social Media and Platform Historiography: Challenges and Opportunities. *TMG – Journal for Media History* 22(1): 6–34. http://www.tmgonline.nl/articles/434/
- Helmond A, Nieborg DB and van der Vlist FN (2019) Facebook's evolution: development of a platform-as-infrastructure. *Internet Histories* 3(2): 123–146. DOI: 10.1080/24701475.2019.1593667.
- https://tags.hawksey.info/
- https://github.com/digitalmethodsinitiative/dmi-tcat
- https://www.docnow.io/
- https://catalog.docnow.io/
- https://webrecorder.net/
- https://conifer.rhizome.org/
- https://timetravel.mementoweb.org/
- https://socialscience.one/
- https://github.com/emilymae/web-archives-bib/blob/master/research-by-theme.md#5-twitter-and-social-media