# Pardon Our Dust
## public data infrastructure under construction

👷🏾 **Rahul Bhargava, Data Culture Group** Ⓜ rahulbot@vis.social 🐦 rahulbot

👷 **Hayden DelCiello, Lazer Lab**

👷 **Zhen Guo, Lazer Lab** 🐦 VeraGuo5

👷 **Alyssa Smith, Lazer Lab** 🐦 cetacean_needed Ⓜ cetacean_needed@mas.to

# Background

the need for public data infrastructure

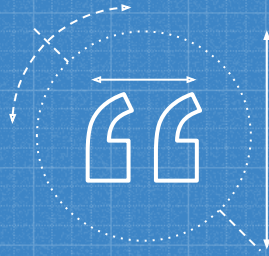# Why do we need public data infrastructure?

- Online media data mostly serves private interests, with very high barriers to access

- This data is critical for understanding changes in online media patterns and impacts

- There is a *public good* here here that is poorly understood and lacks tooling

# Why are we "under construction"?

The landscape, builders, and tools all keep changing.

We're trying to make the process of constructing this infrastructure transparent so it seems more achievable.

We need adaptable, functional infrastructure – but that doesn't come out of thin air!

# Existing Tools

# Big Picture

An evolving ecosystem of online social/media archives

# Building Infrastructure Together

The growing need to understand the societal impacts of online media mean we need to work together across sectors to build new shared digital public infrastructure.

Technically: APIs are how we can interoperate and together create more than our individual sum of parts.

Socially: sharing methods and approaches, governance structures, and convenings can help us build on each other's learnings.

There are a large set of partners we already know about, build on, and work with.

# Media Cloud

**Access**: free

**Data Source**: custom open news scraping

**Interface**: web-based search UI and API

**Metadata**: basic

**Concerns**: current data instability



https://search.mediacloud.org

# Stanford TV News Analyzer

**Access**: free

**Data Source**: large-scale ingestion of Reddit content recordings of US-based cable news networks; indexed on closed captions

**Interface**: web-based search and analysis tool with CSV downloads

**Metadata**: time, channel, person entities, duration

**Concerns**: depends on cloud service support grants



https://tvnews.stanford.edu/

# Pushshift.io

**Access**: free

**Data Source**: large-scale custom ingestion of Reddit submissions and comments

**Interface**: API

**Metadata**: time, subreddit, score, etc.

**Concerns:** handoff to NCRI leaves dataset in flux, and aligned with defense interests; take-down requests



https://pushshift.io

# CrowdTangle

**Access**: free-ish

**Data Source**: Facebook, Instagram, Reddit

**Interface**: web-based monitors, API

**Metadata**: various

**Concerns**: owned by Meta; most of staff laid off; unreliable statements about content included and data; data only includes subsets of platform



https://crowdtangle.com/

# newscatcher

**Access**: freemium

**Data Source**: open news on the web

**Interface**: API

**Metadata**: date, title, authors, country, and others

**Concerns**: startup that might go out of business

</newscatcher>

News API    Pricing    »Developers    Live Demo    Get AP

# News Data: structured, relevant, real-time

Search multi-language worldwide news articles published online with NewsCatcher's News API
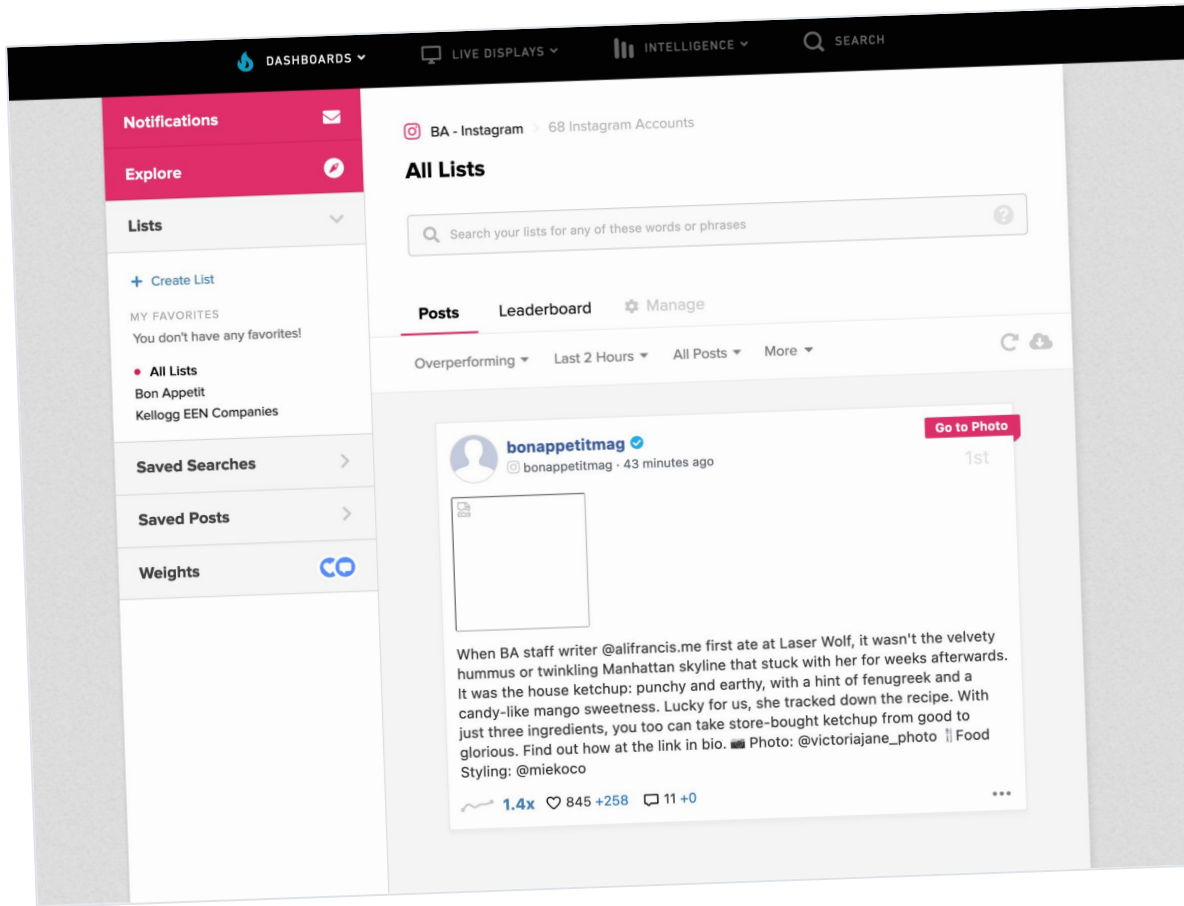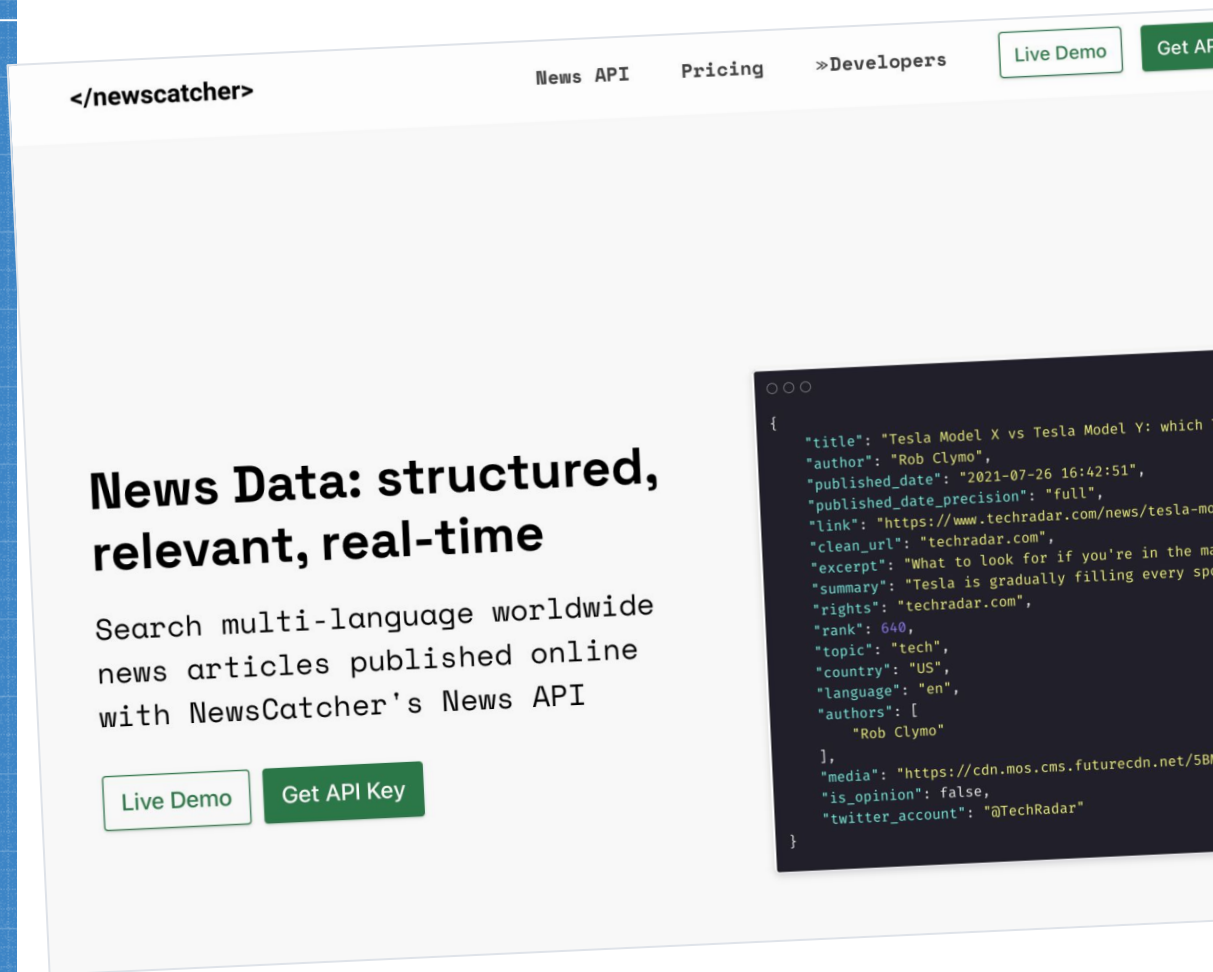
Live Demo    Get API Key

```
{
    "title": "Tesla Model X vs Tesla Model Y: which
    "author": "Rob Clymo",
    "published_date": "2021-07-26 16:42:51",
    "published_date_precision": "full",
    "link": "https://www.techradar.com/news/tesla-mc
    "clean_url": "techradar.com",
    "excerpt": "What to look for if you're in the ma
    "summary": "Tesla is gradually filling every spc
    "rights": "techradar.com",
    "rank": 640,
    "topic": "tech",
    "country": "US",
    "language": "en",
    "authors": [
        "Rob Clymo"
    ],
    "media": "https://cdn.mos.cms.futurecdn.net/5BB
    "is_opinion": false,
    "twitter_account": "@TechRadar"
}
```

ttps://newscatcherapi.com
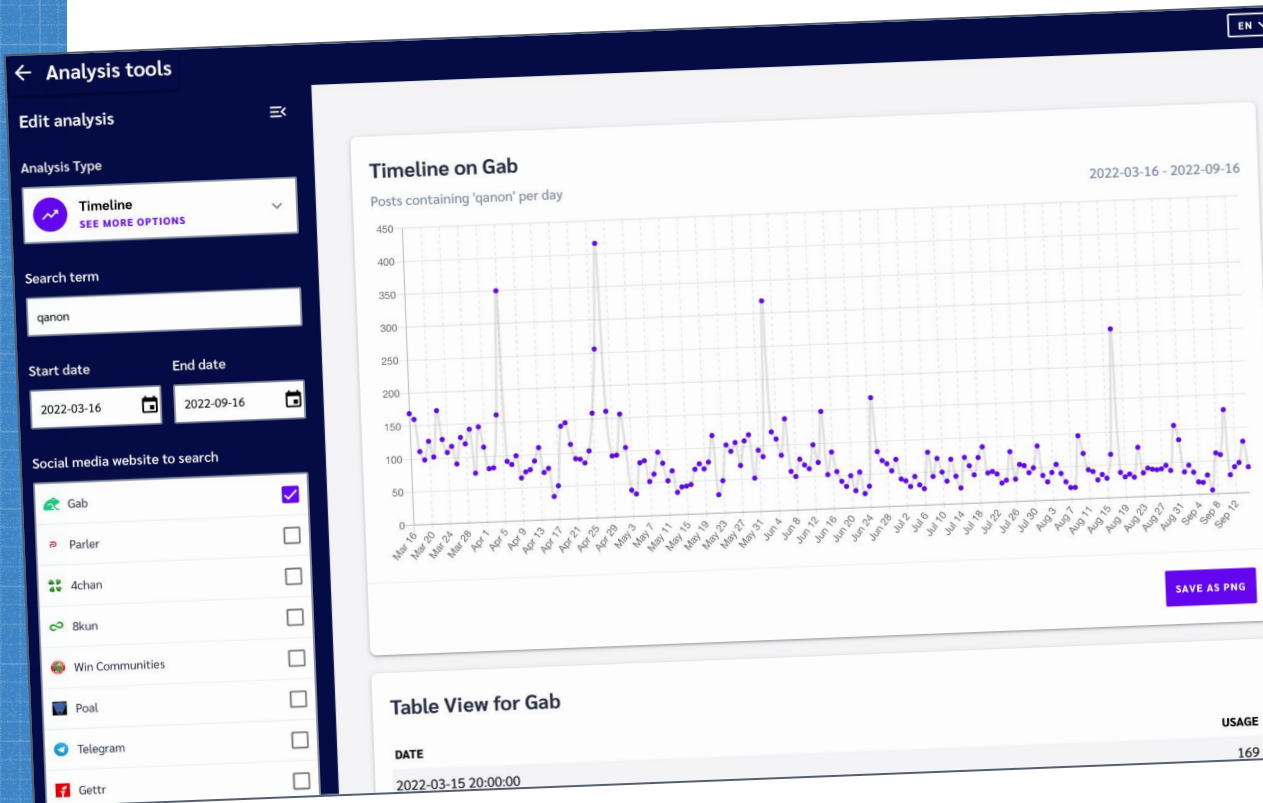
# SMAT

**Access**: freemium

**Data Source**: Parler, Telegram, 8kun, 4chan, Gab, and other fringe sites

**Interface**: web-based search UI and API

**Metadata**: basic

**Concerns**: longer term sustainability; sustainability of antagonistic data collection



https://www.smat-app.com

# Platform APIs

**Access**: freemium

**Data Source**: internal

**Interface**: API calls

**Metadata**: most of it

**Concerns**: cost; lack of transparency; real usage limits; changing policies

# Social Analytics Companies

We've used Brandwatch (used to be called Crimson Hexagon).

**Access**: costs real money

**Data Source**: paying social media platforms for access

**Interface**: web-based dashboards and API calls

**Metadata**: loads

**Concerns**: little transparency; cost; unclear methods for things like sampling

# Media Cloud

an searchable public archive of
global online news stories

# Media Cloud is...

The most comprehensive **database of digital news** in the world available to researchers.

A set of online analysis and visualization **tools and methods.**

A team of **technologists and researchers**.

A **research service** for a community of academics, journalists, foundations, and nonprofits that want to understand the online media ecosystem.

# The open news archive: what is/isn't it?

## Featured Collections

**ONLINE NEWS**

### U.S. Top Digital Native Sources 2018

Top U.S. digital native sources of 2018, based on research from the Pew Research Center published in Aug. 2019.

**ONLINE NEWS**

### U.S. Top Newspapers 2018

Top U.S. newspapers of 2018, based on research from the Pew Research Center published in Aug. 2019.

**ONLINE NEWS**

### U.S. Top Sources 2018

Top U.S. newspapers and digital native sources of 2018, based on research from the Pew Research Center published in Aug. 2019.

**ONLINE NEWS**

### Tweeted Somewhat More by Followers of Liberal Politicians 2019 (US Center Left 2019)

Media for which url sharing on twitter is aligned with the U.S. partisan center left

**ONLINE NEWS**

### Tweeted Mostly by Followers of Conservative Politicians 2019 (US Right 2019)

Media for which url sharing on twitter is aligned with the U.S. partisan right

**ONLINE NEWS**

### Tweeted Evenly by Followers of Conservative & Liberal Politicians 2019 (US Center 2019)

Media for which url sharing on twitter is aligned with the U.S. partisan center

**ONLINE NEWS**

### Tweeted Mostly by Followers of Liberal Politicians 2019 (US Left 2019)

Media for which url sharing on twitter is aligned with the U.S. partisan left

**ONLINE NEWS**

### Tweeted Somewhat More by Followers of Conservative Politicians 2019 (US Center Right 2019)

Media for which url sharing on twitter is aligned with the U.S. partisan center right

**ONLINE NEWS**

### India - National

Media is largely about India

# Search interface demonstration

# Existing Applications: What are people doing now?

- Researching misinformation
- Monitoring media
- Using news as data

# Technical concerns

**Stability** – large software and data infrastructures are hard to maintain without a team of technical PhD students or staff

**Grants** – few private funders want to support infrastructure and maintenance, but public funding is competitive too
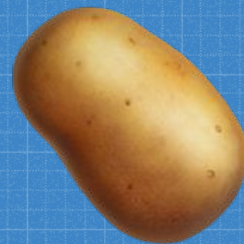
# Techno-social concerns

**Copyright –** 2 billions documents is a *lot* of copyright liability

**Compute Power –** processor cycles add up to become expensive very quickly

**Governance –** shared resources require shared governance structures, which need trust, shared interests, etc.

# What we're starting with: the Twitter panel dataset

- Matches state-level voter records + demographic info to social media accounts.
- We have ~1.6 million users' tweets from the inception of Twitter to around October 2022.
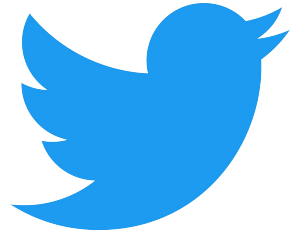
# Privacy Concerns

- **<u>This is sensitive data!</u>** Our data usage agreement doesn't allow us to share the voter files or the linkages from voter file to Twitter accounts.
- We are able to report aggregate statistics about the panel (e.g. "What is the gender breakdown of people tweeting about lizards?")
- Even when we report aggregate statistics, if the buckets become small, we worry about identifiability.

# Our Solution (a very brief overview)

- POTATO (Panel-Based Open Term-Level Aggregate Twitter Observatory) is a searchable version of the Twitter panel.
- It returns search results in aggregate along with data visualizations of the demographic/geographic information.
- Users can also download a JSON blob of the aggregate data for further analysis.
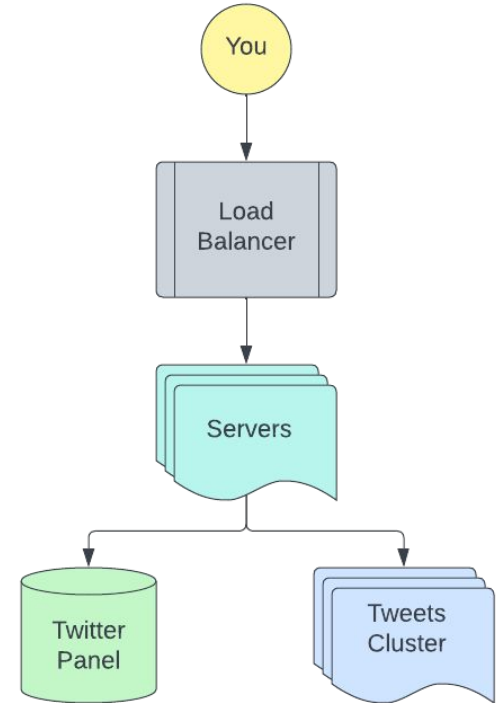
# Searchability

- You search Tweet text -> We tell you who Tweeted it

- We store Tweets from the Twitter Panel

- Searches not hindered by Twitter API limits

# Scalability

- On-demand aggregation is expensive

- Vertically scalable in hardware
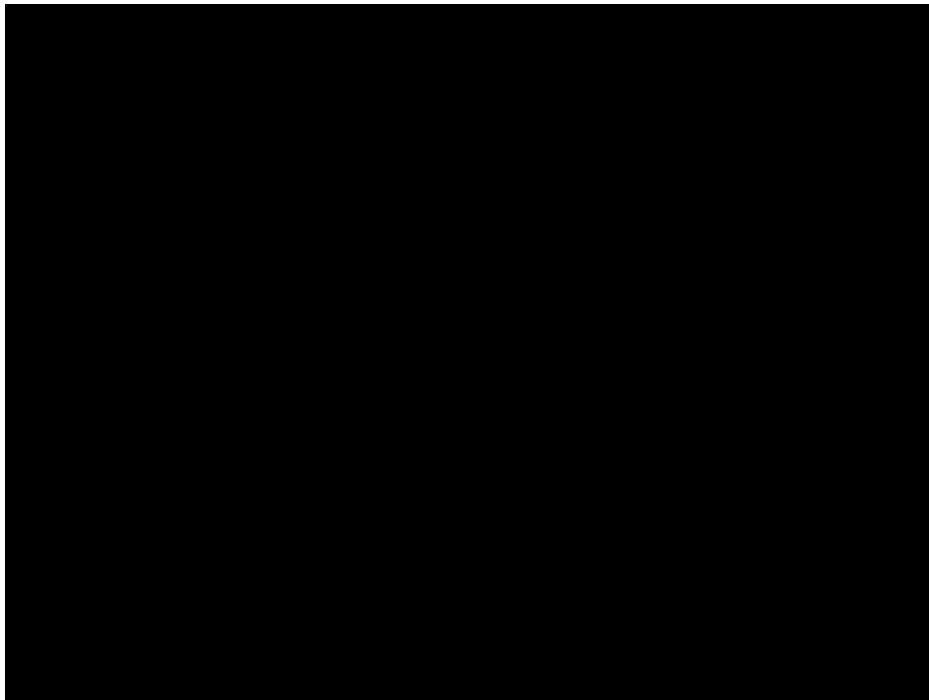
- Horizontally scalable in software

# Visualization Tool & Website UI

- Streamlit
- Open-source
- Pure python
- Beautiful visualization
- Easy and quick to learn
- Support all the features we need

# A Quick Video of POTATO In Action

# Visualization Tool & Website UI
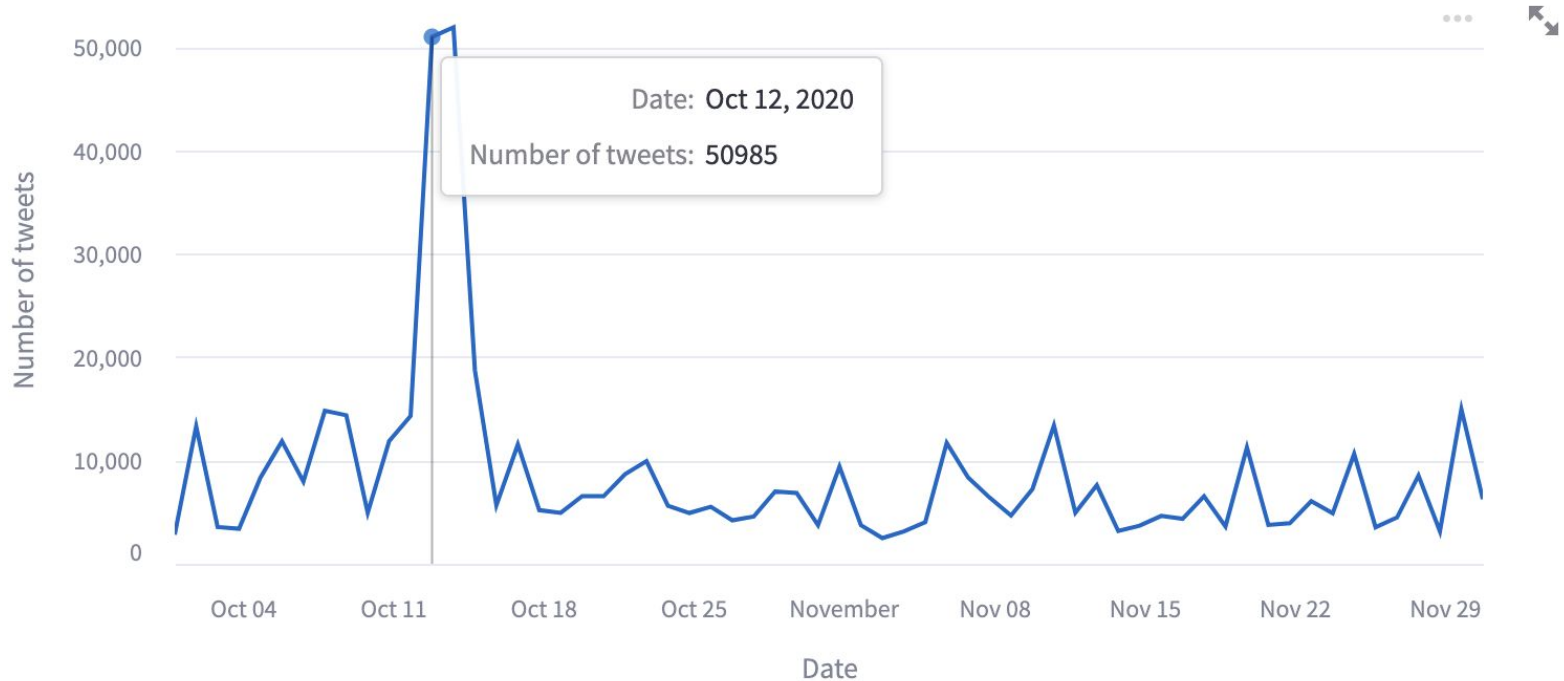
# Twitter Panel Dashboard

Search for keyword

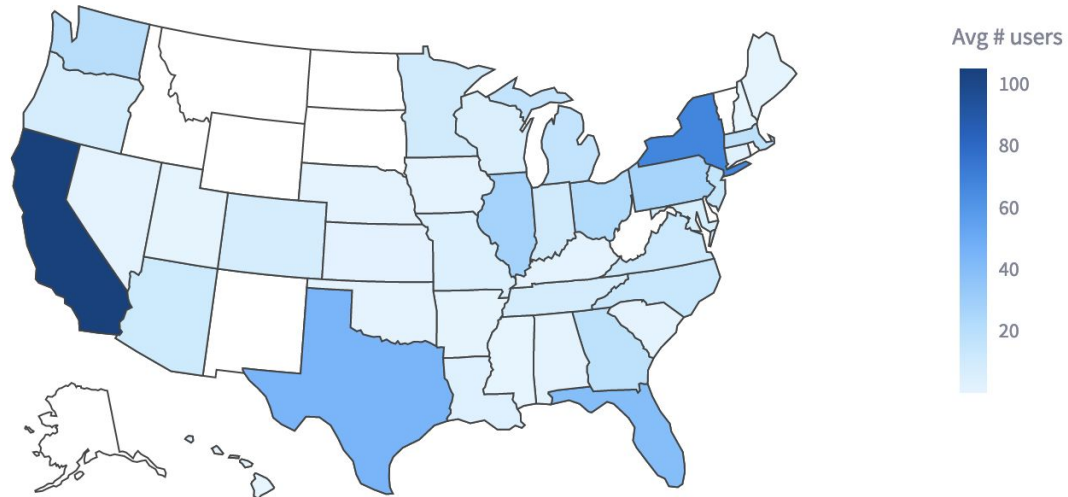optimus prime

Aggregate time based on:

day ▾

# Visualization Tool & Website UI

## Number of tweets per day

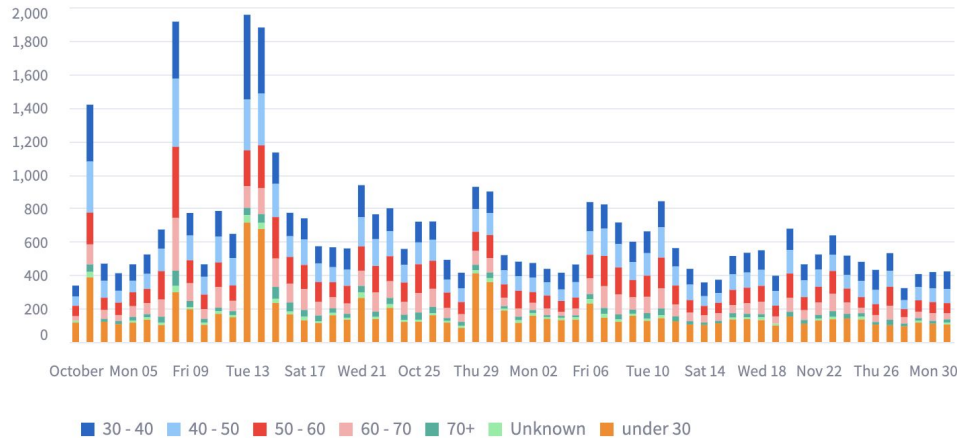# Visualization Tool & Website UI

**Average number of unique users per day**

# Visualization Tool & Website UI
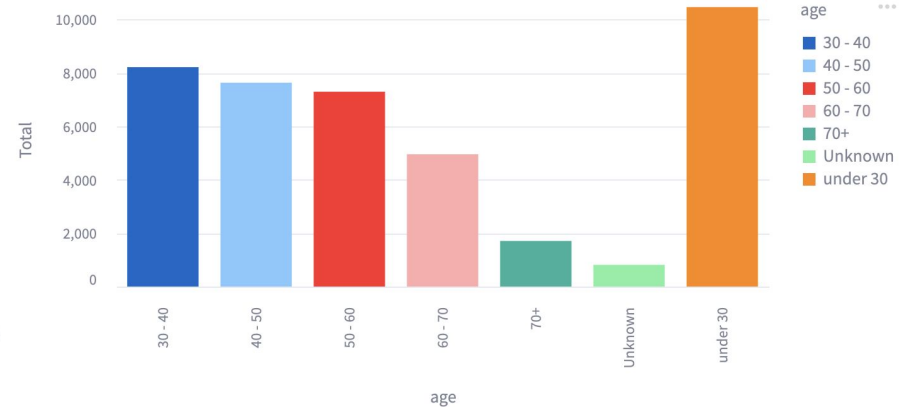
Select variables to visualize:

age ▾

## Users by age category per day



October  Mon 05  Fri 09  Tue 13  Sat 17  Wed 21  Oct 25  Thu 29  Mon 02  Fri 06  Tue 10  Sat 14  Wed 18  Nov 22  Thu 26  Mon 30

■ 30 - 40   ■ 40 - 50   ■ 50 - 60   ■ 60 - 70   ■ 70+   ■ Unknown   ■ under 30

## Total users by age category*



age
■ 30 - 40
■ 40 - 50
■ 50 - 60
■ 60 - 70
■ 70+
■ Unknown
■ under 30

# Visualization Tool & Website UI

Select variables to visualize:

gender ▾

## Users by gender category per day



Female  Male  Unknown

## Total users by gender category*



gender
- Female
- Male
- Unknown

# Visualization Tool & Website UI

Select variables to visualize:

race|  ▼

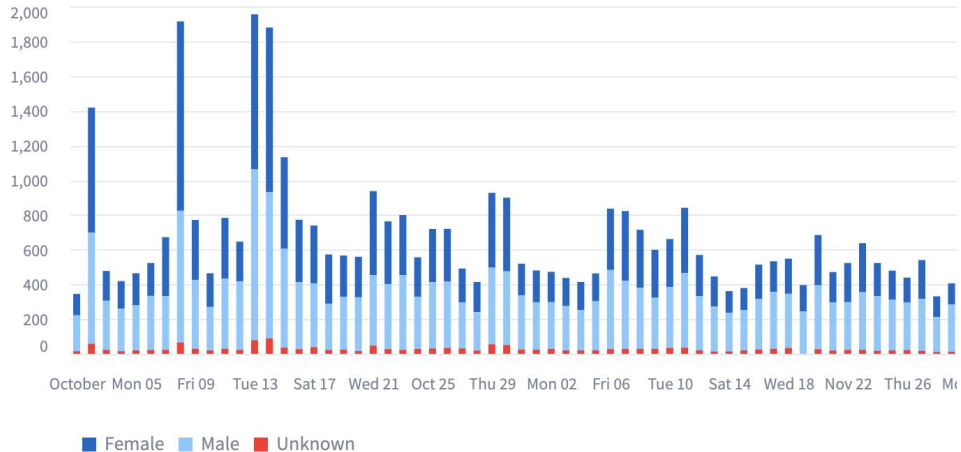## Users by race category per day
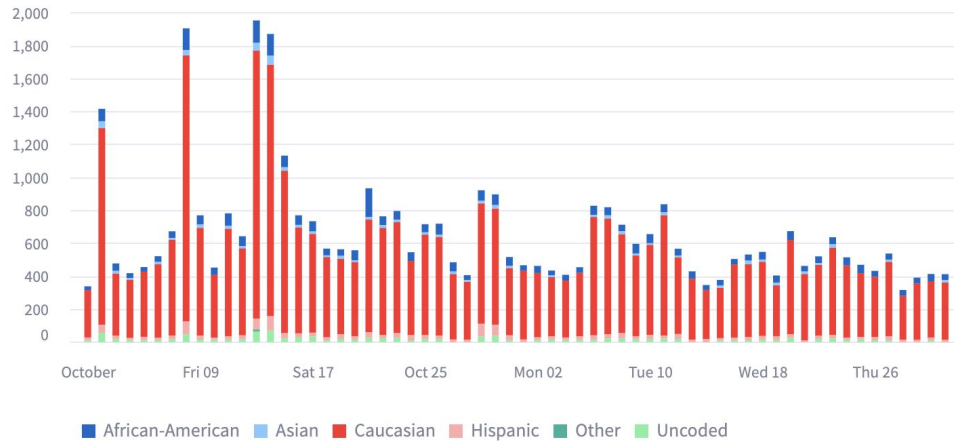


## Total users by race category*



race

- African-American
- Asian
- Caucasian
- Hispanic
- Other
- Uncoded

# Visualization Tool & Website UI
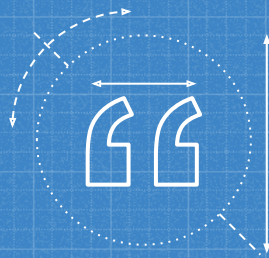
# What Now?

a call to action

# Calls to Action

- How would you, or are you, using these types of tools now?
- How can we leverage collective infrastructure to open access to more archives?
- What approaches to creating sustainable technologies like this are working?

# General Discussion

# General Discussion (with audience)

- What else exists in this space that we can support & work with?
- What are people working on at the moment?
- What needs do you, in particular, have?
- What are you hoping to build? How can we help you build it?

Thank You