

On the use of Applied Machine Learning and Digital Infrastructure to leverage Social Media Data in Health and Epidemiology

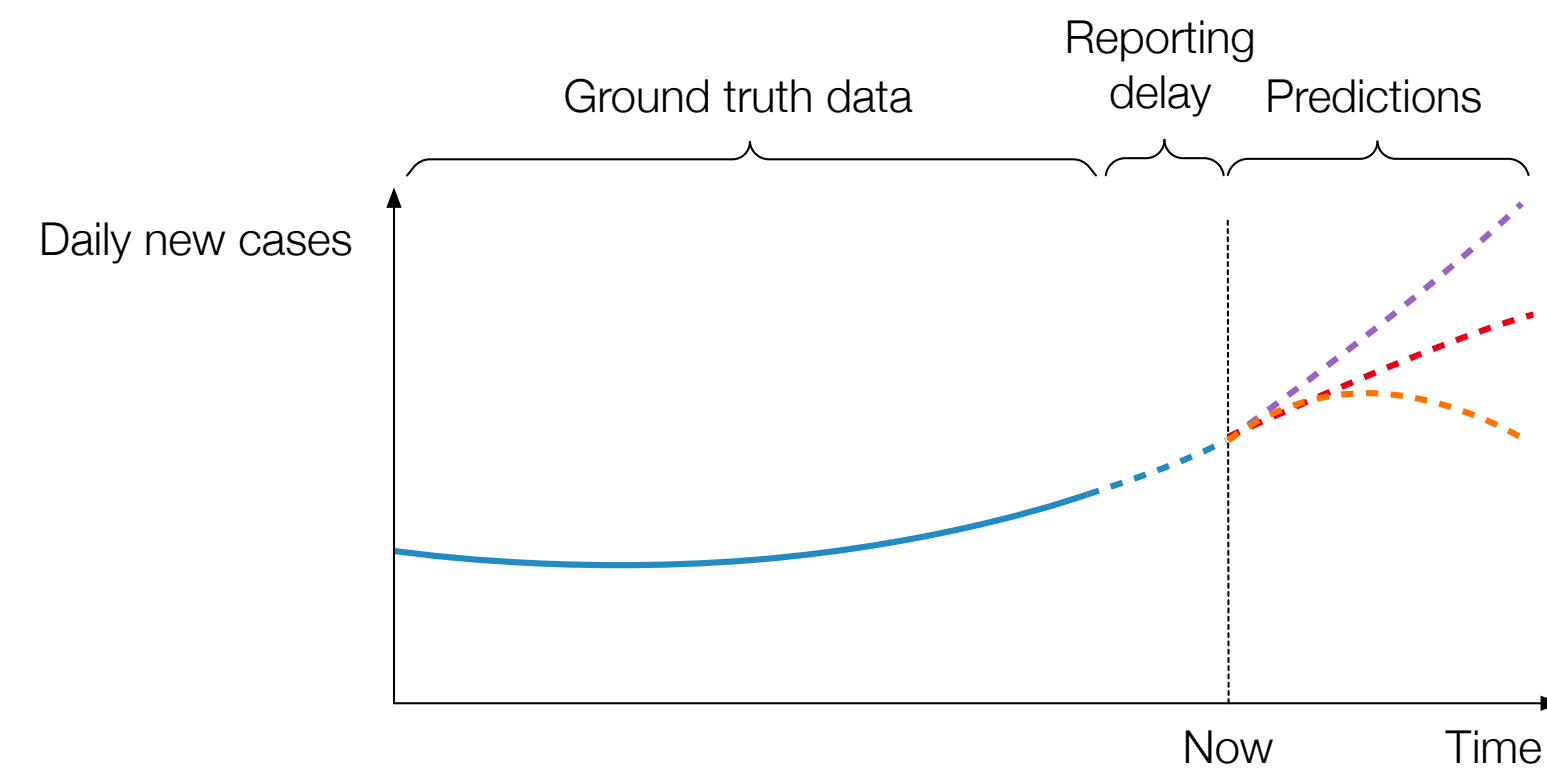


Martin Müller

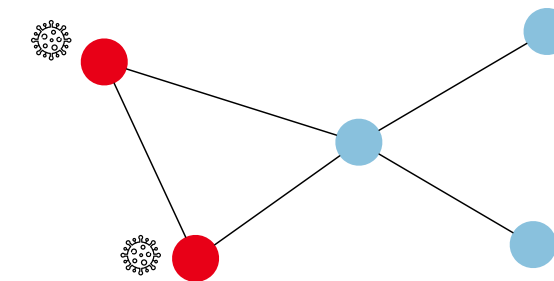
Digital Epidemiology Lab, EPFL

Summary of thesis defense (January 2021)

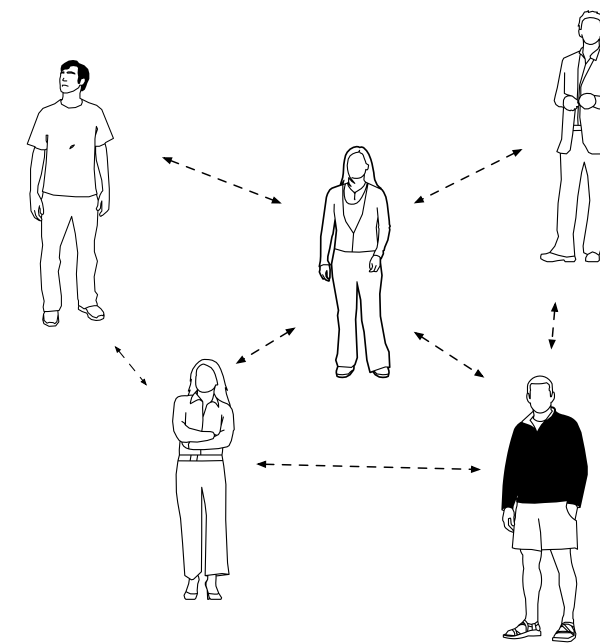
Human behavior has major impacts on disease spread



Biological contagion
on a contact network



Social contagion
network



Key characteristics

- Transmissibility
- Contact structure
- Immunisation status
- Adherence to control measures
- Vaccine hesitancy
- Health behaviors

... but how can we quantify these health behaviors?

Digital Epidemiology

- Physical activity trackers
- Weblogs
- Search engine data
- Mobile communication data
- Social media

and

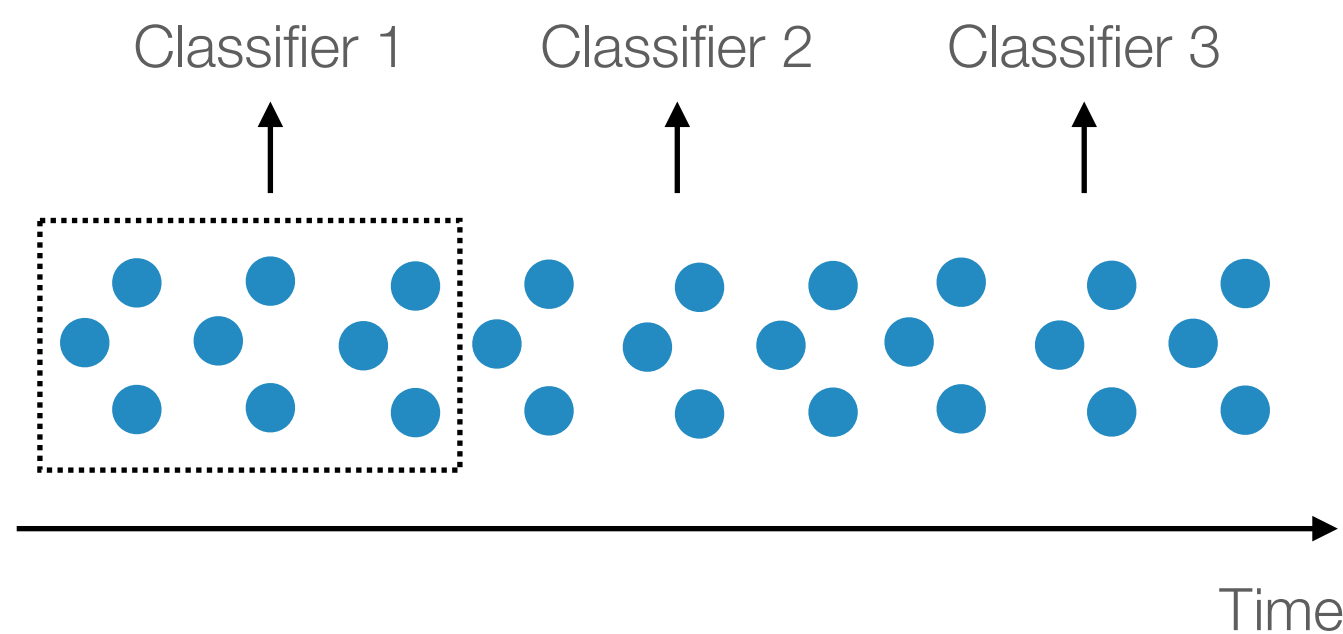
Traditional methods

- Surveys / questionnaires
- Clinical encounters

... however, traditional methods have several drawbacks

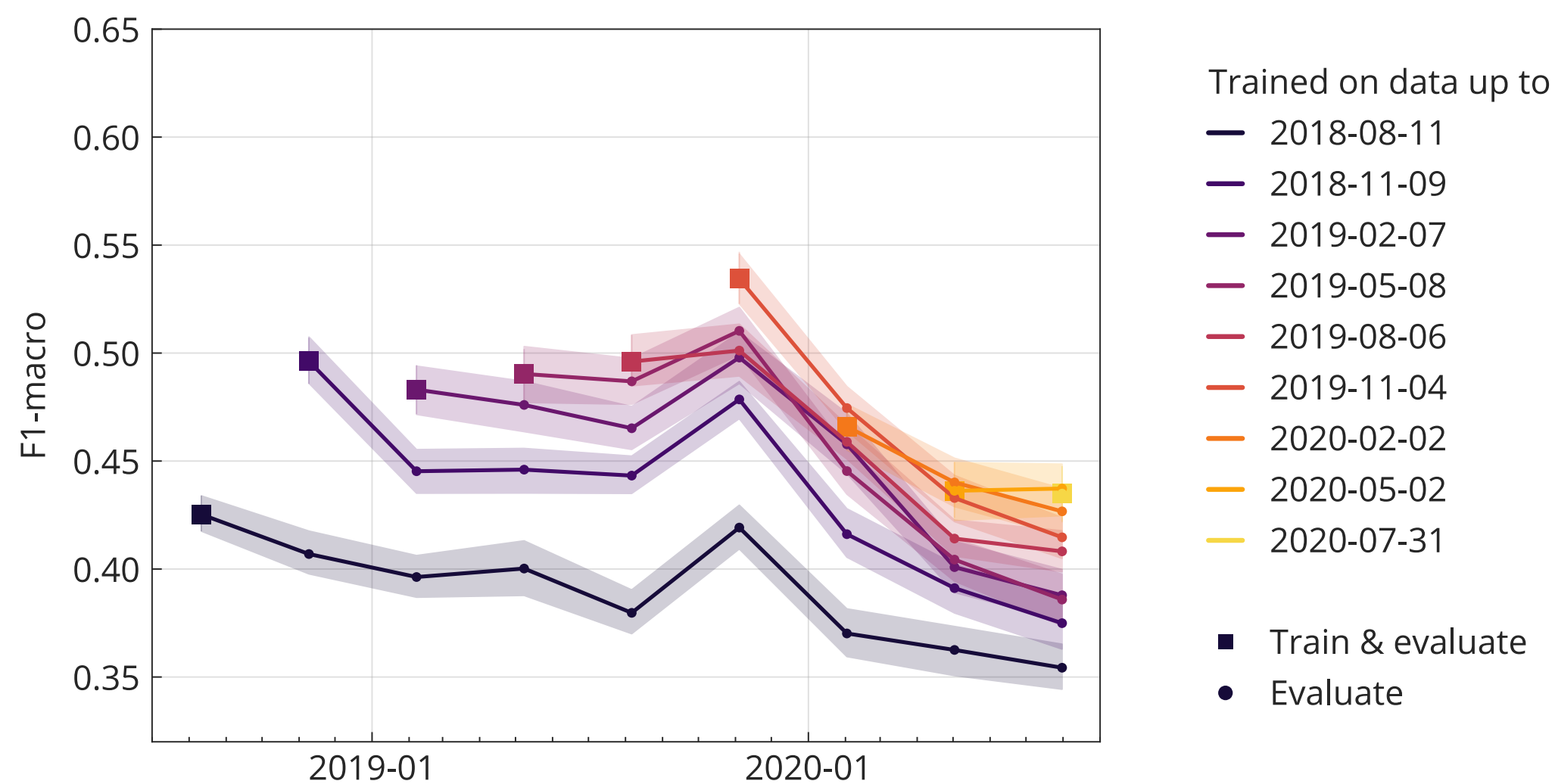
1. Lack of temporal dimension
2. Reporting delay
3. Suitable for reportable diseases only
4. Expensive/time-consuming to conduct
5. Response bias

Machine Learning concept drift: A key challenge in Digital Epidemiology



Question: How strong is the effect of concept drift of text classifiers trained on social media data?

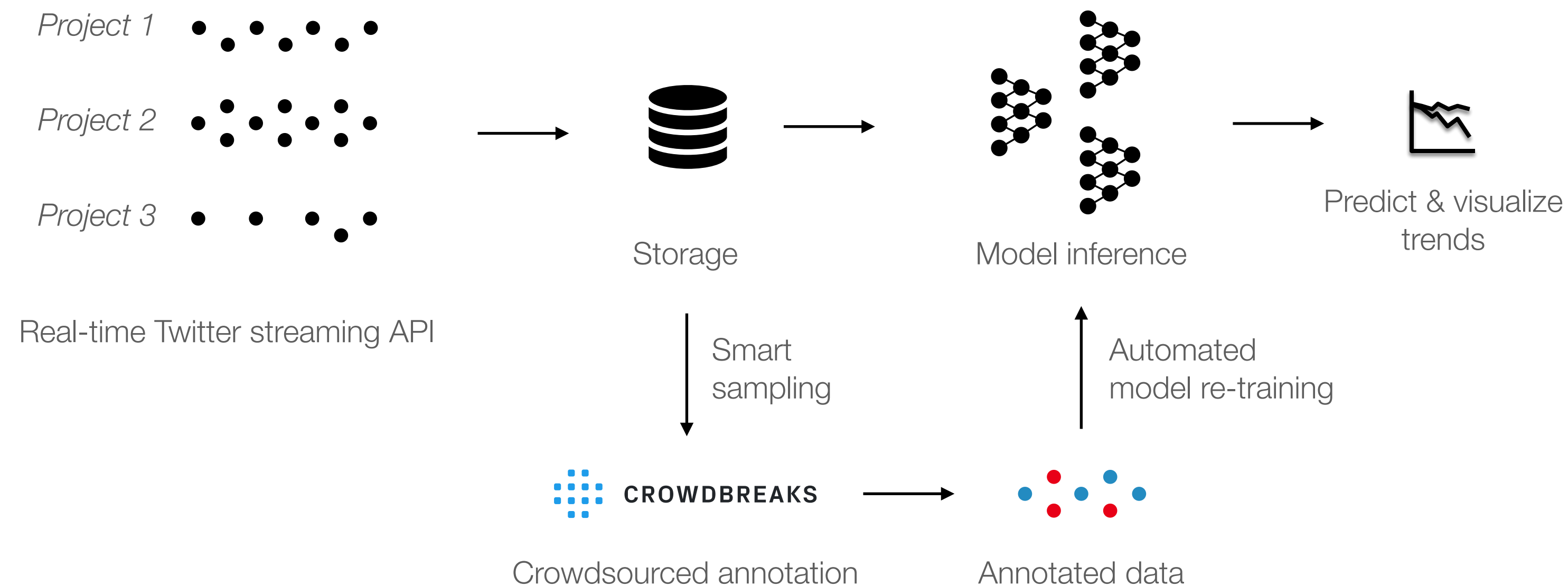
Method: We train text classifiers to predict vaccine hesitancy from Twitter data of various time windows and evaluate it on future/unseen test data



Model performance drops by up to 20% over the course of 6 months.

This highlights the need for **continuous model retraining** for Machine Learning based Public Health surveillance projects

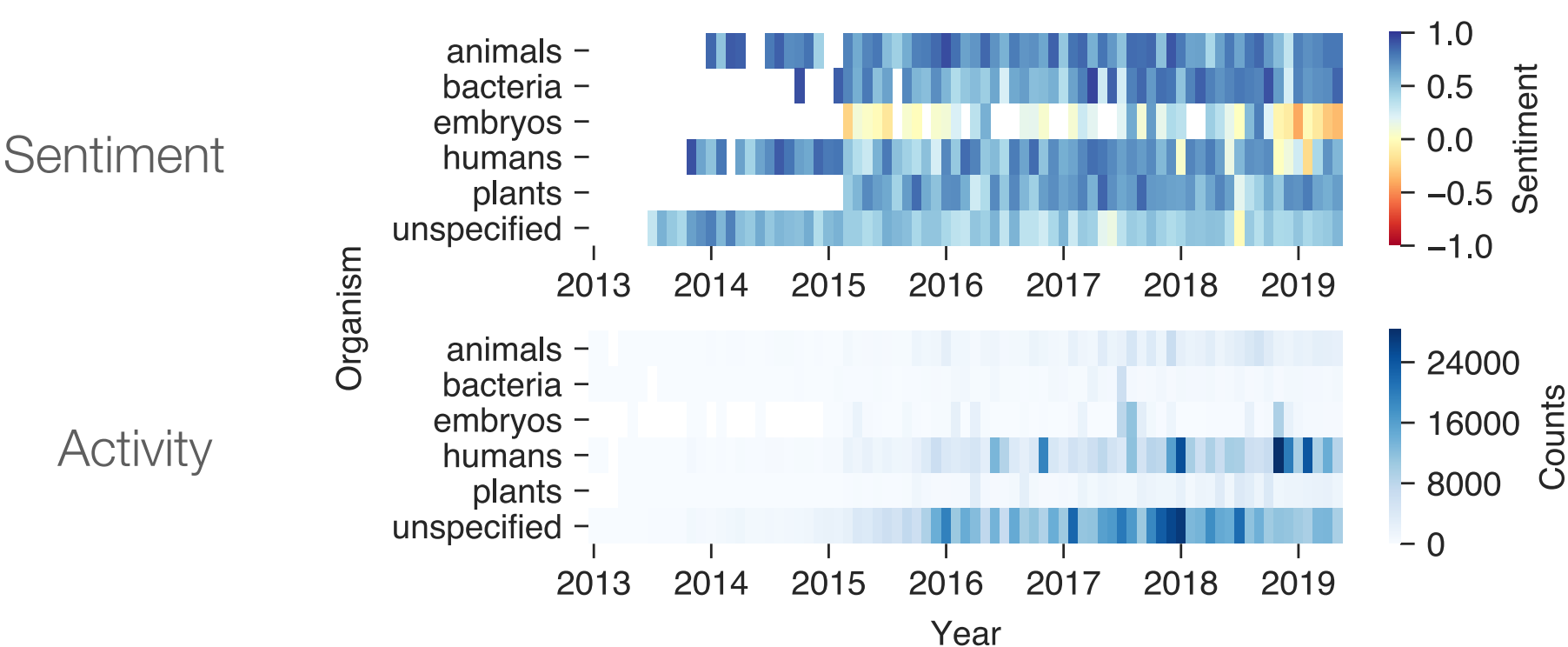
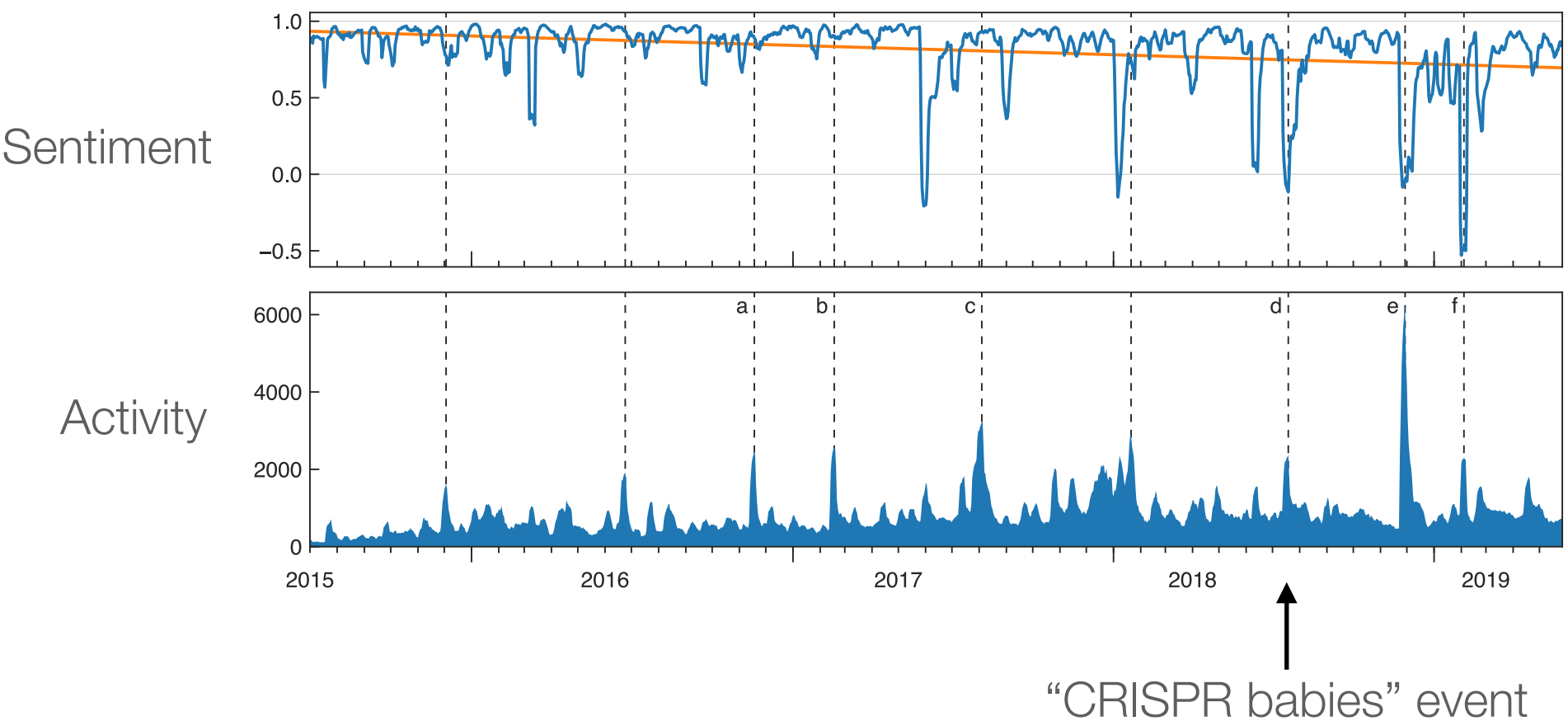
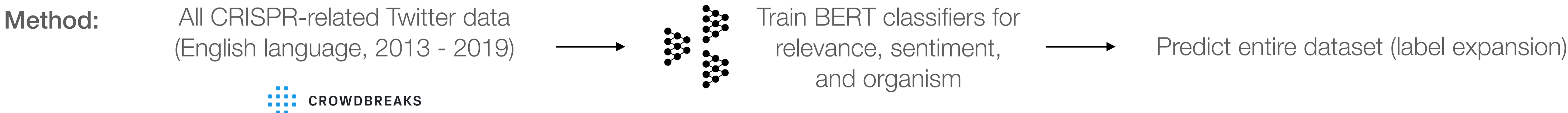
Crowdbreaks: An Open Source platform to track Public Health trends



- Ingestion of real-time signals from social media enables to follow **long-term** health trends
- Automated re-training as a means to combat Machine Learning **concept drift**
- Improved **automation and reproducibility** for research studies using social media data

Public opinion towards CRISPR/Cas9

Question: How is the novel gene technology CRISPR/Cas9 discussed online?



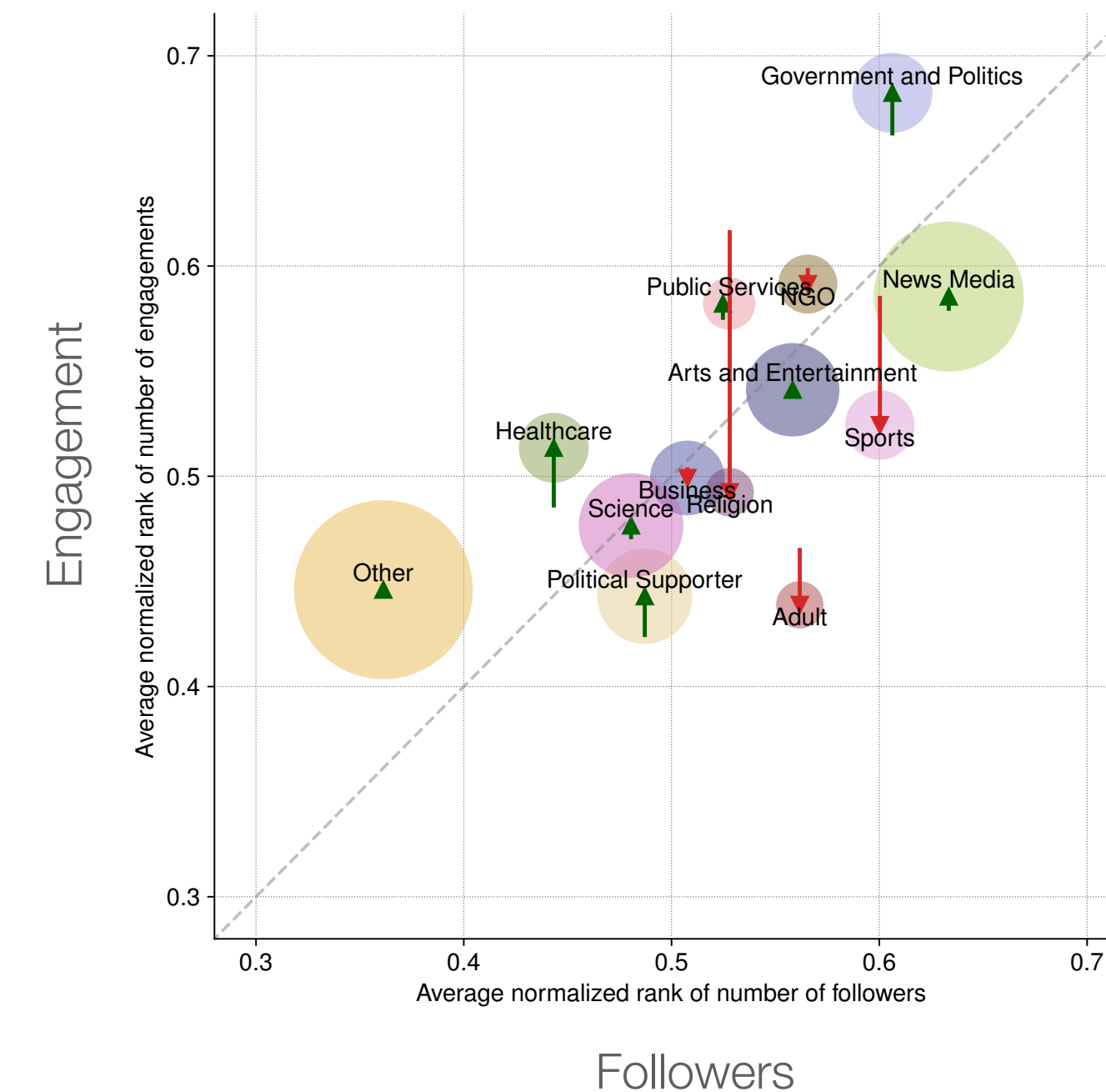
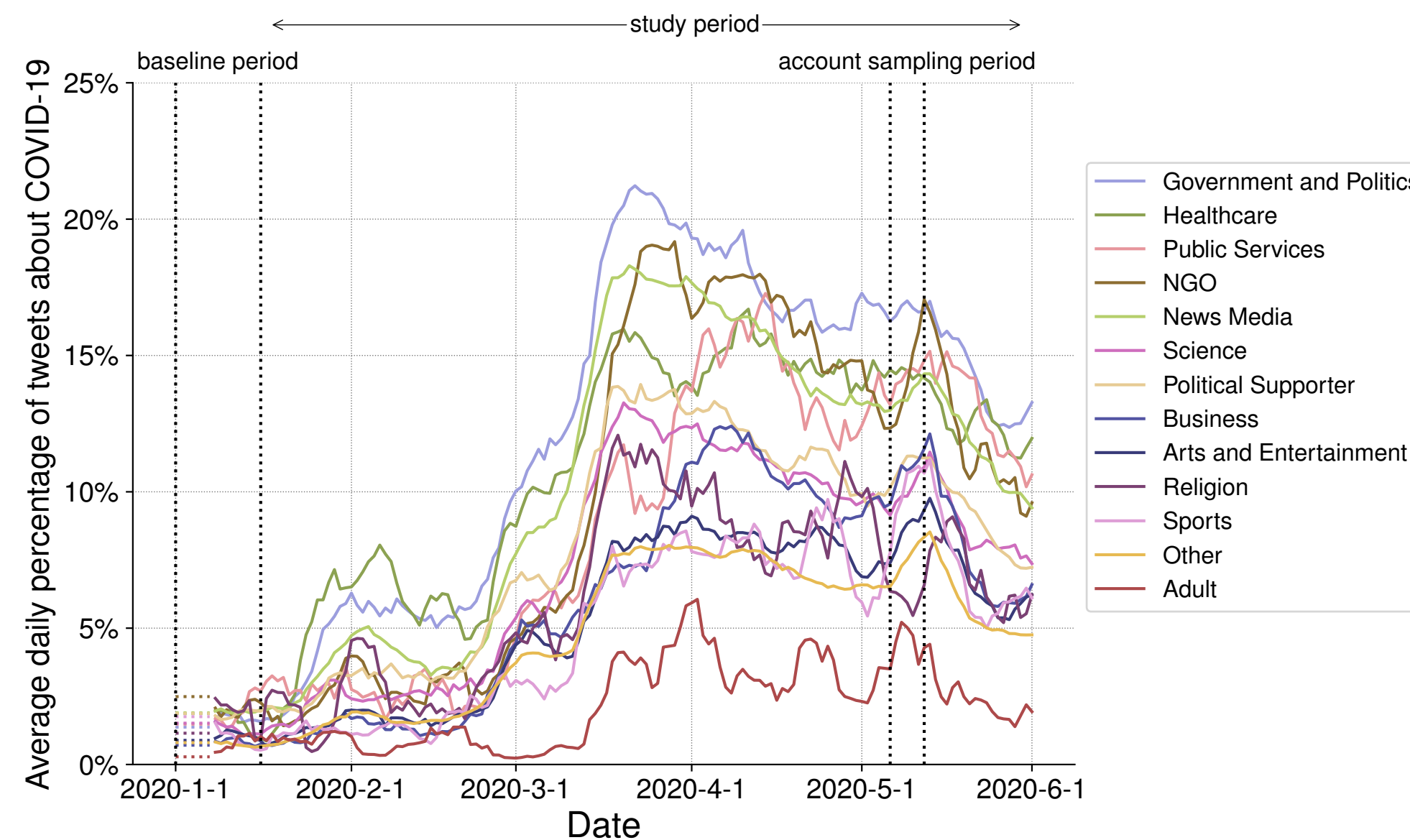
- Overall CRISPR is discussed in positive light, despite multiple scandals
- The use of CRISPR in the context of embryos is viewed increasingly negative
- Results and general trends are in line with several past surveys

Attention mechanisms during early phase of COVID-19

Question: Who was speaking, and who was being heard in the beginning of the COVID-19 pandemic?

Method:

- 1) Sample users from complete COVID-19 Twitter stream in all languages
- 2) Large-scale annotation to categorize user descriptions
- 3) Compare engagement on users' activity between January and June, 2020



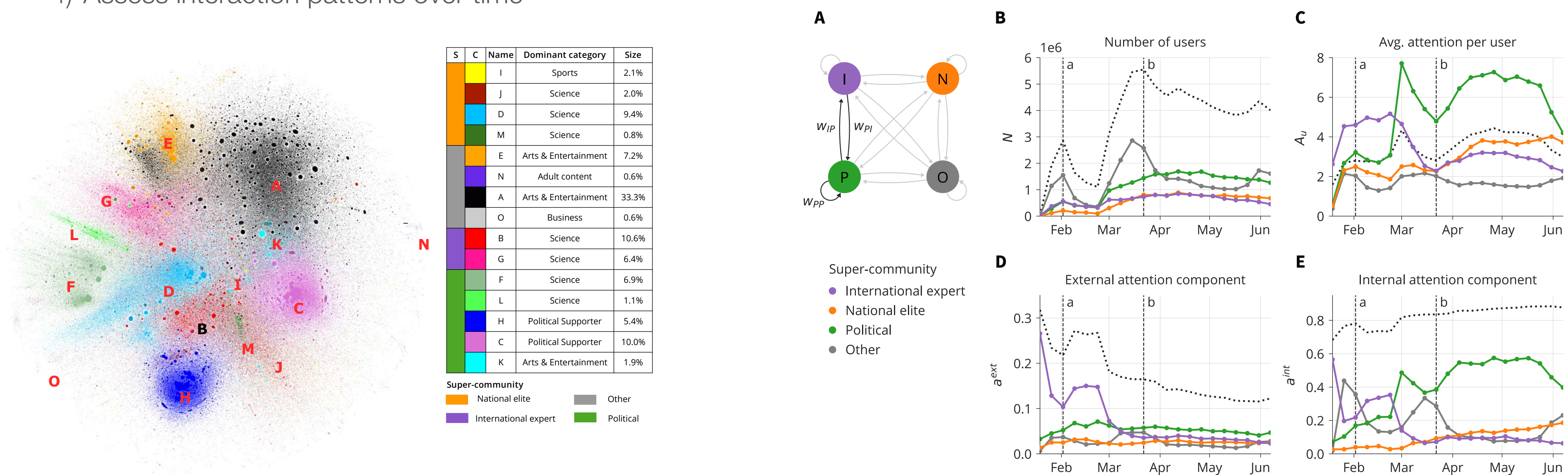
- Accounts linked to Healthcare, Science and Government & politics received strong boost
- Experts received disproportionately more attention than news media

Network dynamics of interactions during COVID-19

Question: How did COVID-19 change the nature of interactions online?

Method:

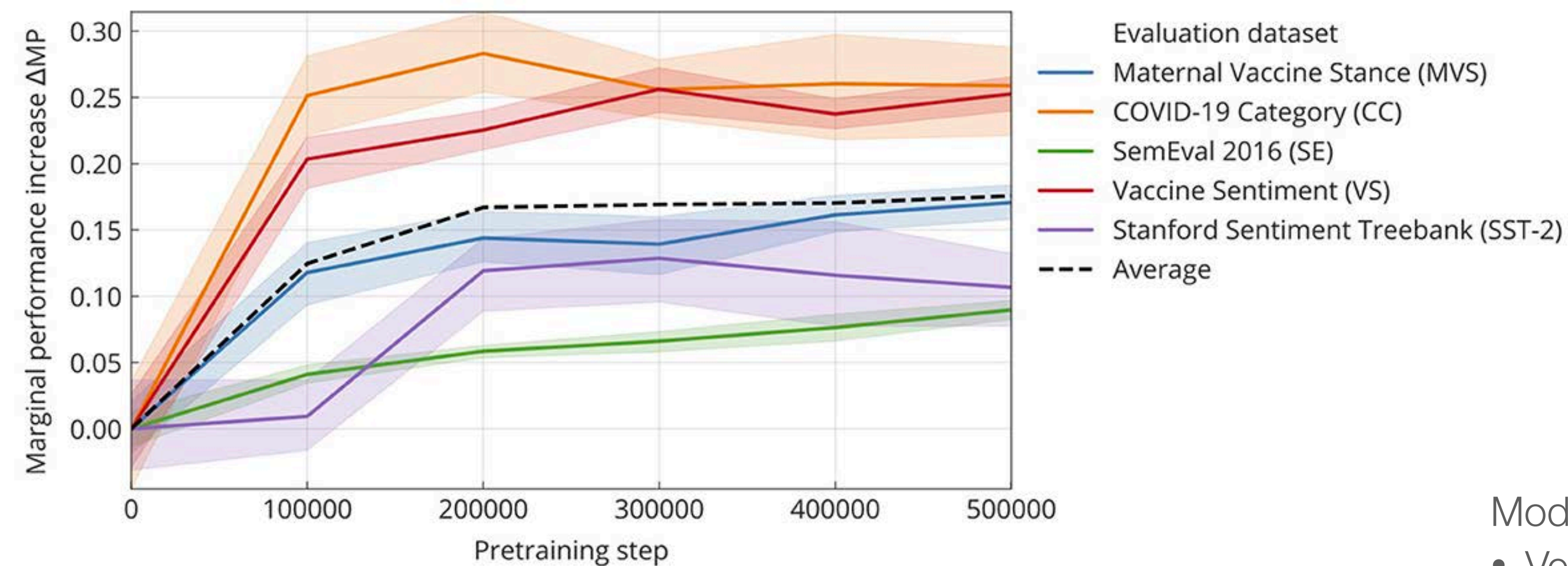
- 1) Build retweet network of 350M tweets, 26M users
- 2) Run a community detection algorithm
- 3) Characterize communities and merge into 4 super-communities: International expert, national elite, political, and other
- 4) Assess interaction patterns over time



- Initially international experts receive a lot attention and have broad reach
- Throughout the pandemic experts become increasingly isolated
- Overall, a growing politicization and polarization of the discourse can be measured

CT-BERT: Domain-adaptation of BERT models

Domain-specific pretraining of BERT-Large on 160M tweets related to COVID-19



(Official CT-BERT logo)

Domain-specific pretraining of language models, such as BERT, can yield up to a 30% relative increase in downstream text classification performance

Models available on the HuggingFace Hub 🤗:

- Version 1: `digitalepidemiologylab/covid-twitter-bert`
- Version 2: `digitalepidemiologylab/covid-twitter-bert-v2`
- Zero-shot model: `digitalepidemiologylab/covid-twitter-bert-v2-mnli`

<https://github.com/digitalepidemiologylab/covid-twitter-bert/>

Acknowledgements

Supervisor

Marcel Salathé

Collaborators

Per Egil Kummervold
Manuel Schneider
Effy Vayena
Francesco Durazzi
Kristina Gligorić
Manoel Horta Ribeiro
Robert West
Daniel Remondini
Michael Edelstein
Elaine Okanyene Nsoesie
Rafael Ruiz de Castañeda

Developers

Yannis Jaquet
Sean Carroll

MSc students

Tigist Menkir
Axel Uran
David da Rocha Rodrigues

Digital Epidemiology Lab

Marina Secat
Gianrocco Lazzari
Sharada Prasanna Mohanty
Djilani Kebaili
Chloé Alléman
Talía Salzmänn
Sylvain Bernard
Olesia Altunina
François Pichard

EPFL Extension School

Harry Anderson
Xavier Adam
Florian Laurent

Partner

Burcu Tepekule Müller