

Beyond twitter

Exploring bluesky.social for digital disease detection and prototyping a data extraction pipeline for ILI surveillance

Heiner Atze, MSc, PhD

Digital Epidemiology 2025, Hasselt University

2025-04-10

Outlininglines I

- 1 Introduction
- 2 Exploration of bluesky data
- 3 Project
- 4 Methods
- 5 Data extraction
- 6 Results

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

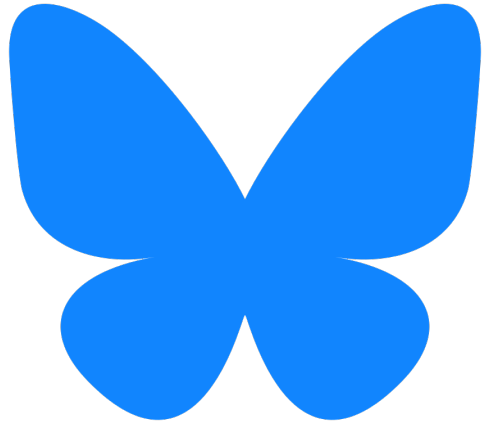
Data
extraction

Results

Introduction

bluesky: general aspects

- microblogging platform
- similar to twitter in user experience
- decentralized
- open source



Decentralization and Democratization of content algorithms ¹

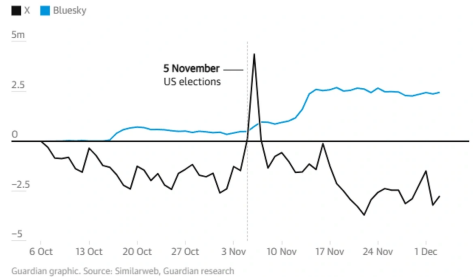
- Decentralized User Identifier (DID)
 - immutable, associated with human readable user handle
- Personal Data servers (PSDs)
- DIDs and affiliated contents are portable between PSDs
- Users can choose, prioritize and develop feed generators and content labelers

¹Balduf et al. (2024)

Development of user activity ²

- current estimate: ca. 33 Millions active users
- user base expanded in bursts after key events:
 - 2022: acquisition of twitter by Elon Musk
 - 2024: ban of X in Brazil, presidential election in the US

X has lost users since October while Bluesky has gained close to 2.5m
Change in active daily users since 6 October 2024



²Duarte, Balduf et al. (2024)

Literature addressing bluesky

Introduction

Exploration of bluesky data

Project

Methods

Data extraction

Results

- Google scholar search : “bluesky” AND “social” since 2022
- 43 articles
- main topics:
 - decentralized social network architecture
 - user migration from X to bluesky 2024
 - network structure and dynamics
- no results for
 - “bluesky” AND “disease”
 - “bluesky” AND “epidemiology”

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

Data
extraction

Results

Exploration of bluesky data

bluesky API

- publicly accessible for free
- extensive documentation at <https://docs.bsky.app/docs/category/http-reference>

searchPosts API method

- API documentation
- selected parameters:
 - q: search query
 - since, until: defining search period
- deterministic search
- allows exhaustive sampling

getProfiles

- allows to retrieve the author profile information
- for reference, not used in this project

Post metadata

- defined in the SDK documentation
- fields (selection):
 - `uri`: unique post identifier
 - `author`: contains `did` which allows to retrieve user profile
 - `record`: contains the text and time information of the message
 - `embedded`: any embedded media (images, other posts, etc ...)
- in contrary to former twitter post metadata, no geoinformation

User information

- Feedgens
- Labelers
- no geo information

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

Data
extraction

Results

Project

Introduction

Exploration of
bluesky data

Project

Methods

Data
extraction

Results

bluesky post data for digital disease surveillance

Introduction

Exploration of
bluesky data

Project

Methods

Data
extraction

Results

bluesky post data for digital disease surveillance

Implementation of a continuous surveillance pipeline

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

Data
extraction

Results

Methods

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

**Data
extraction**

Results

Data extraction

Symptom related message extraction

- focused on French bluesky posts (data volume constraint)
- extraction using list of keywords
 - grippe (flu, influenza)
 - rhume (common cold)
 - fièvre (fever)
 - courbature (muscle pain)
- extraction of
 - complete message data for further language processing
 -

Basal network activity

- probing of the basal network activity using keywords
 - travail (*work*)
 - demain (*tomorrow*)
 - voiture (*car*)
 - sommeil (*sleep*)
- post counts aggregated by day

Case data

- data downloaded from WHO Flumart = FluID: ILI case data
 - FluNet: virological data

Data processing for time series extraction

- Normalization of ILI post counts by basal network activity
-
- LLM
- ECDC case definition
 - LLM vs. random post selection

Beyond
twitter

Heiner Atze,
MSc, PhD

Introduction

Exploration of
bluesky data

Project

Methods

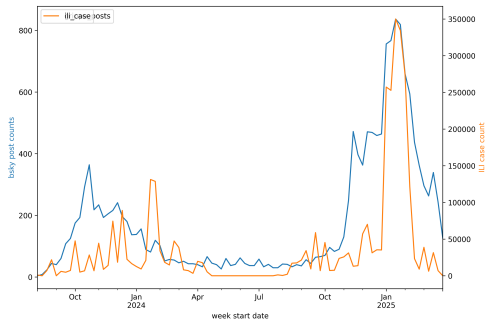
Data
extraction

Results

Results

Raw post counts

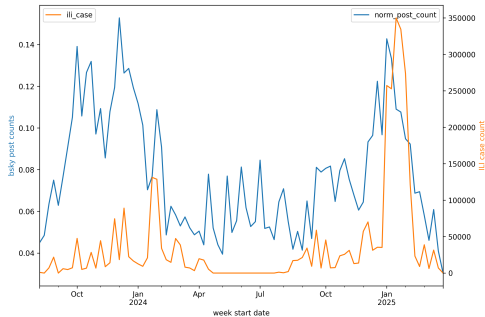
```
Text(0, 0.5, 'ILI case count')
```



Correlation

	grippe_posts	rest_posts
grippe_posts	1.000000	0.884120
rest_posts	0.884120	1.000000
ili_case	0.769662	0.552405

`Text(0, 0.5, 'ILI case count')`



It is not as simple as that :/

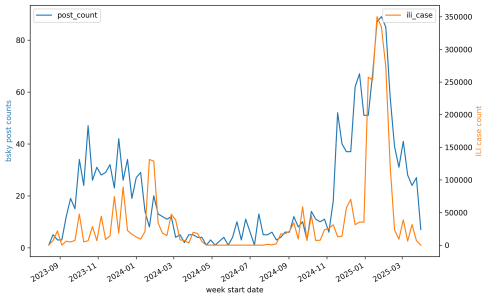
Normalized post counts

Correlation

	norm_post_count	re
norm_post_count	1.000	0.
rest_posts	0.197	1.
ili_case	0.436	0.

LLM annotated post counts, raw

Text(0, 0.5, 'ILI case count')



Correlation

	ili_case	post_count
ili_case	1.00	0.71
post_count	0.71	1.00

Bibliography

Balduf, Leonhard, Saidu Sokoto, Onur Ascigil, Gareth Tyson, Björn Scheuermann, Maciej Korczyński, Ignacio Castro, and Michał Król. 2024. “Looking at the Blue Skies of Bluesky.” In *Proceedings of the 2024 ACM on Internet Measurement Conference*, 76–91.

Duarte, Fabio. “Bluesky User Age, Gender, & Demographics (2025).” <https://explodingtopics.com/blog/bluesky-users>.