

Epicosm: A framework for linking social media data in large-scale birth cohorts

Dr Oliver Davis

Associate Professor and Turing Fellow

University of Bristol and the Alan Turing Institute

bristol.ac.uk





- Engage participants and cohort leaders to assess acceptable approaches
- Develop easy-to-use open-source software for securely linking Twitter data
- Proof-of-principle linking of Twitter data in the ALSPAC cohort
- Dissemination of the software and findings across CLOSER cohorts, funders, participants and policy makers

- Roll out the software to other cohorts
- Develop software to link other social media platforms

The Alan Turing Institute

- Build the framework for algorithm developers to validate their approaches against “ground truth” in cohorts



What do
participants
think?



What do
participants
think?

“Children of the Nineties is fine
because I know you’re not going
to sell it”

“We’re always sharing our data with loads of people all the time who are using our data for advertising and selling it on. At least with this we would have given our consent and knew it was for a good cause.”



What do
participants
think?

“Children of the Nineties is fine
because I know you’re not going
to sell it”

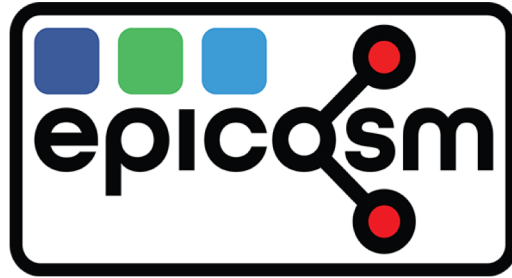
“We’re always sharing our data with loads of people all the time who are using our data for advertising and selling it on. At least with this we would have given our consent and knew it was for a good cause.”



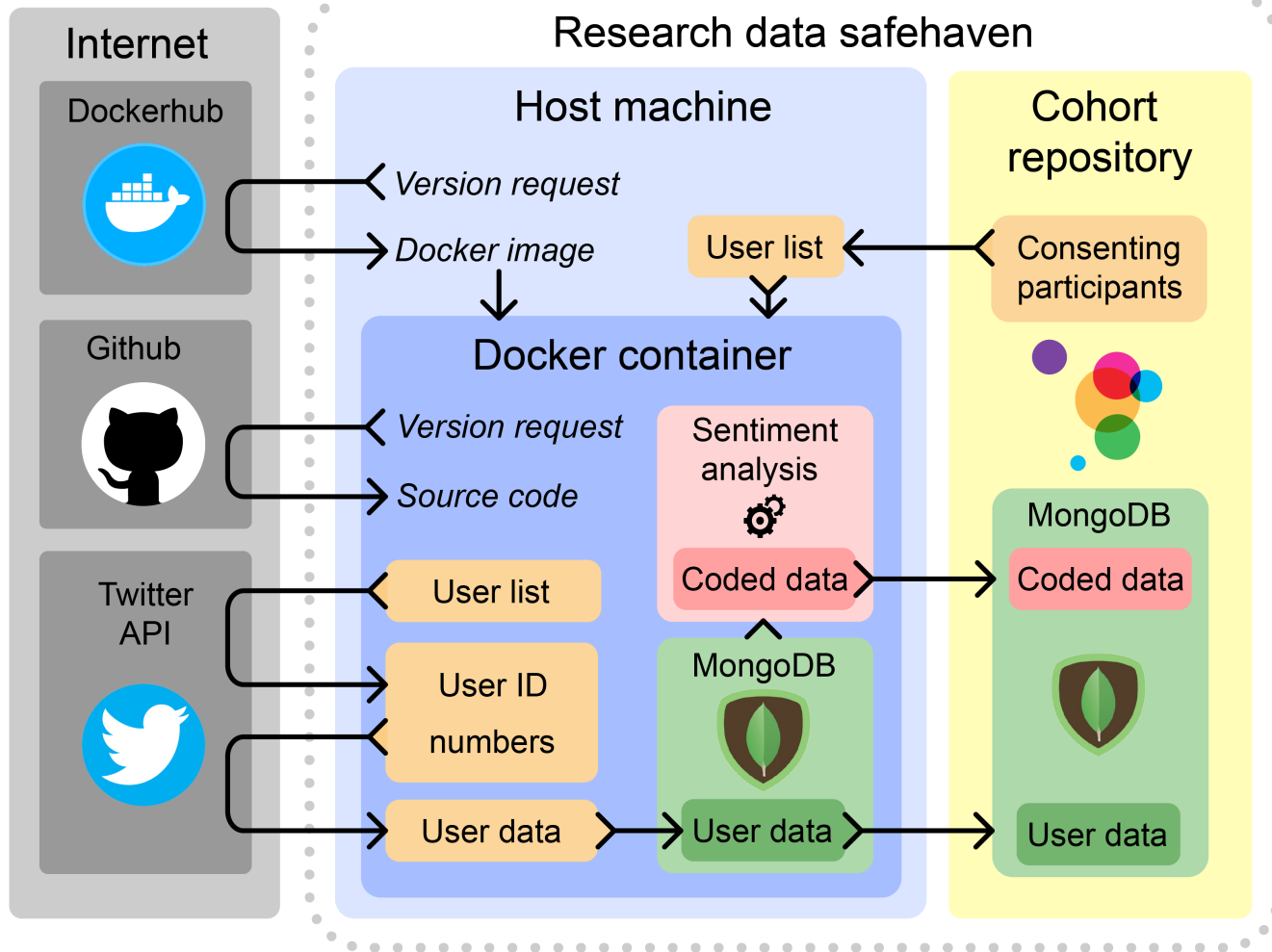
What do
participants
think?

“Children of the Nineties is fine
because I know you’re not going
to sell it”

“I think it’s important. Because to get the
fullest roundest picture you need to do
that anyway, don’t you.”



Epidemiological Cohort Online Social Media



If you're a #twin, go follow @TedsProject.
@oliviaabrownn22 and I have helped them with
their research for years. It's v. interesting

```

{
  "id": "1d93eaa8cc0cacce6787778a51307564e",
  "rev": "1-69c0e71bd5550d3f59f445fdac3f1f94",
  "contributors": null,
  "truncated": false,
  "text": "RT @RealMattieBrown: If you're a #twin, go follow @TedsProject. @oliviaabrown22 and I have helped them with their research for years. ...",
  "in_reply_to_status_id": null,
  "id": 313624475394506750,
  "favorite_count": 0,
  "source": "web",
  "retweeted": false,
  "coordinates": null,
  "entities": {
    "user_mentions": [
      {
        "id": 112801662,
        "indices": [
          3,
          19
        ],
        "id_str": "112801662",
        "screen_name": "RealMattieBrown",
        "name": "Matt"
      }
    ]
  }
}

```

```

{
  "_id": "1d93eaa8cc0cacce6787778a51307564e",
  "_rev": "1-69c0e71bd5550d3f59f445fdac3f1f94",
  "contributors": null,
  "truncated": false,
  "text": "RT @RealMattieBrown: If you're a #twin, go follow @TedsProject. @oliviaabrown22 and I have helped them with their research for years. ...",
  "in_reply_to_status_id": null,
  "id": 313624475394506750,
  "favorite_count": 0,
  "source": "web",
  "retweeted": false,
  "coordinates": null,
  "entities": {
    "user_mentions": [
      {
        "id": 112801662,
        "indices": [
          3,
          19
        ],
        "id_str": "112801662",
        "screen_name": "RealMattieBrown",
        "name": "Matt"
      }
    ]
  }
}

```



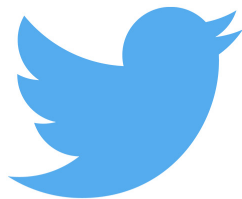

participants



researcher



harvester code



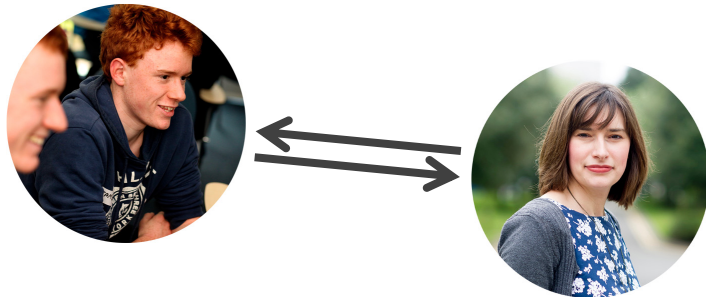
Twitter API



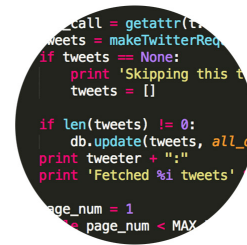
Tweet-scoring
code



mongoDB.

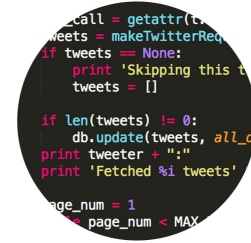
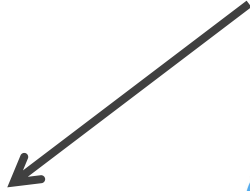


The researcher asks
participants for consent
and user names



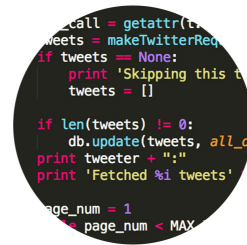


The researcher provides
the user names to the
Python Tweet-harvesting
code





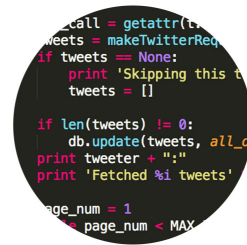
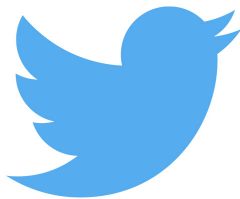
The harvester contacts the Twitter API to convert the participants' user names to persistent identifiers



mongoDB.



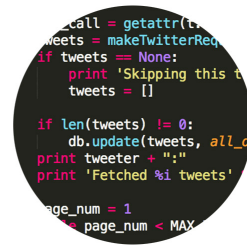
Every week, the harvester contacts the database to find the ID of the most recent Tweet collected for each participant



mongoDB.



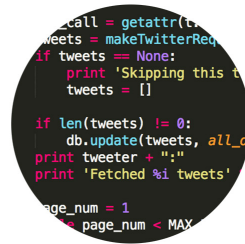
Then the harvester
contacts the Twitter API
to download all the
Tweets from the past
week for each
participant (as JSON
objects)



mongoDB.



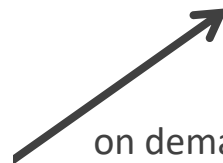
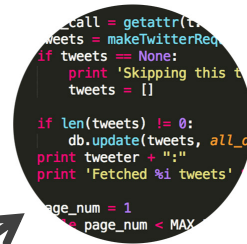
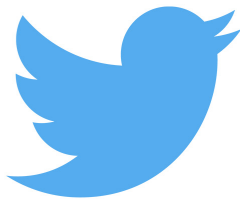
The harvester stores the
Tweets in the database
natively as JSON
documents



mongoDB.



When requested, the database exports the Tweet data to a text file and sends the file to the Tweet-scoring code



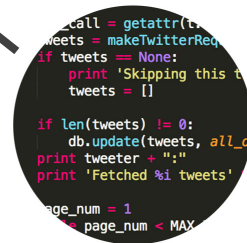
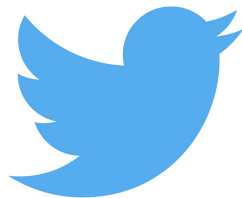
on demand



mongoDB.



The Tweet-scoring code
analyses the text of
every Tweet and
provides scores to the
researcher

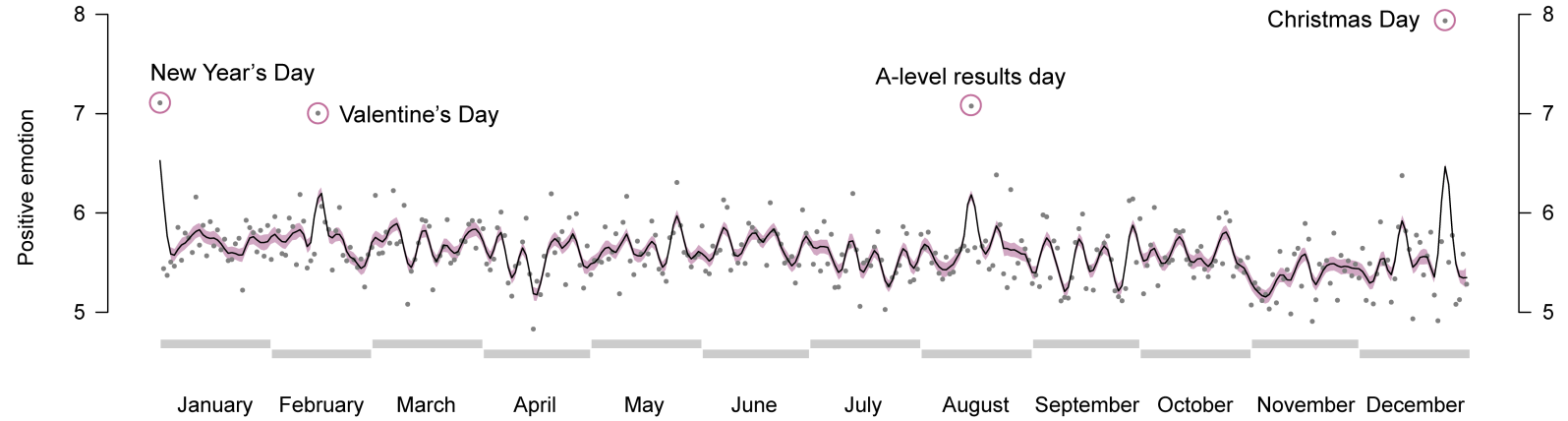


mongoDB.

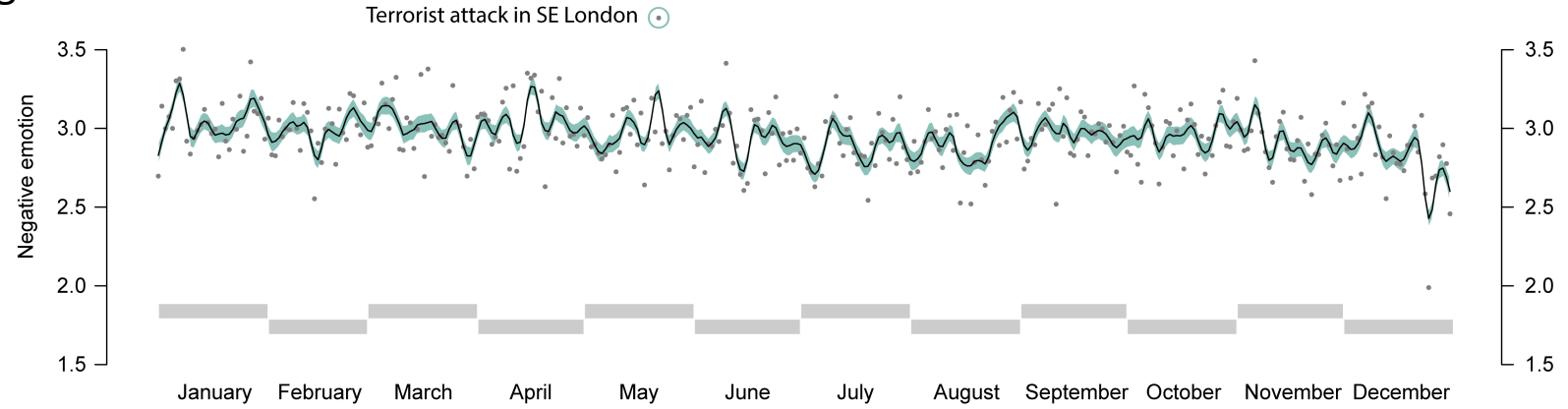
If you're a #twin, go follow @TedsProject.
@oliviaabrownn22 and I have helped them with
their research for years. It's v. interesting

if	4.66		
you're	5.30		
a	5.24		
twin	6.14	6.14	
go	5.54		
follow	5.66		
and	5.22		$6.14 + 7.28 + 6.46 + 6.12 + 7.52$
i	5.92		
have	5.82		$= 33.52$
helped	7.28	7.28	
them	4.92		
with	5.72		$33.52 / 5$
their	5.16		
research	6.46	6.46	$= 6.704$
for	5.22		
years	5.28		
it's	4.88		
very	6.12	6.12	
interesting	7.52	7.52	

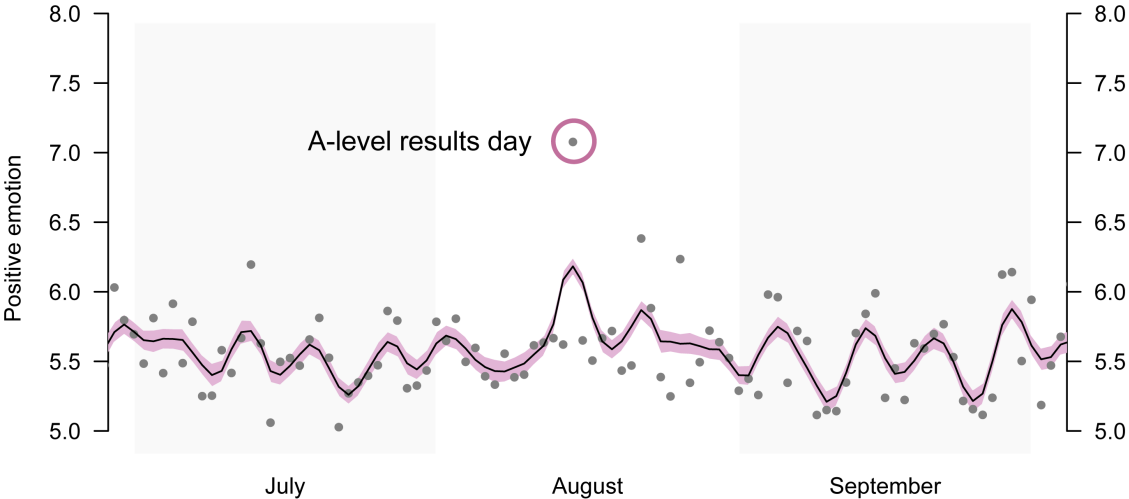
Positive emotion



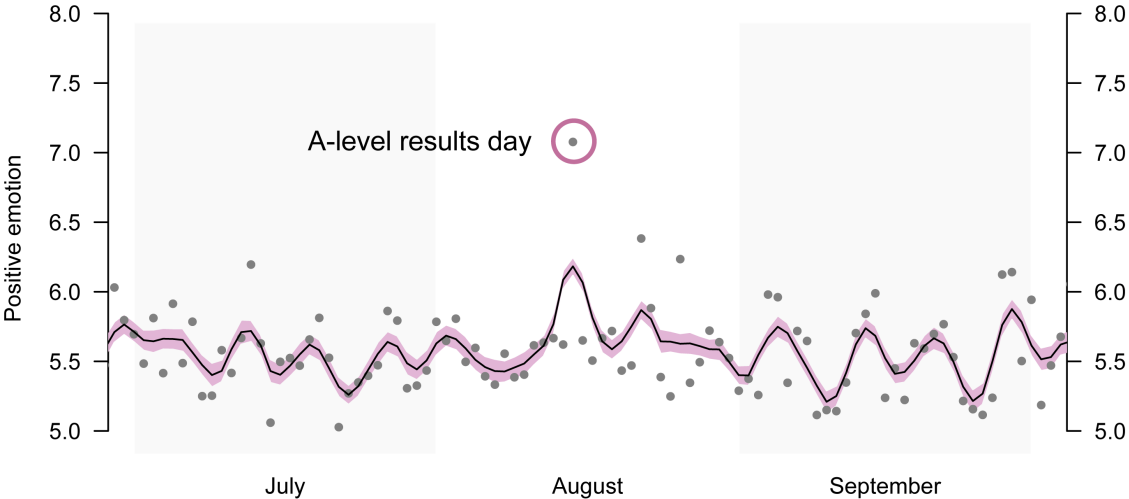
Negative emotion



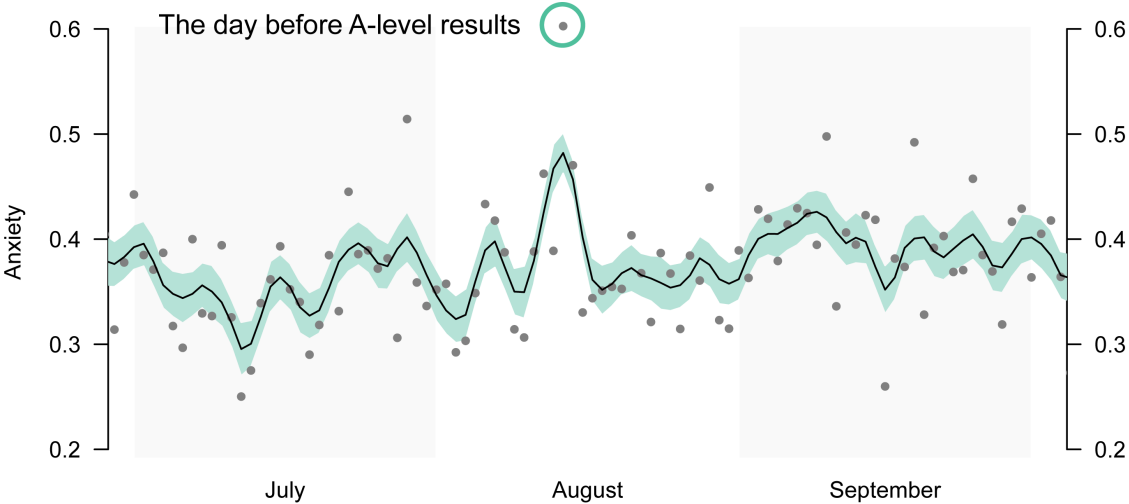
Positive emotion



Positive emotion



Anxiety



Acknowledgements

University of Bristol

Claire Haworth

Al Tanner

Nina Di Cara

Andy Boyd

Claire Bowring

Richard Thomas

Lynn Molloy



King's College London

Robert Plomin

Stephanie Kliskey

Andy McMillan



UCL

Lisa Calderwood



Cardiff University

Luke Sloan



University of Essex

Tarek Al Baghal

