



Demographic Inference and Representative Population Estimates from Multilingual Social Media Data

Zijian
Wang



Scott
Hale



David
Adelani



Przemyslaw
Grabowicz



Timo
Hartmann



Fabian
Flöck



David
Jurgens

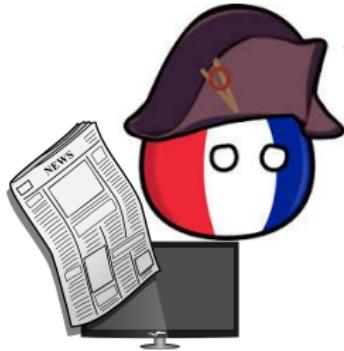


“What is the most important issue
facing the country today?”



Unemployment!

"What is the most important issue facing the country today?"



L'environnement!

Refugee crisis!





“What is the most important issue
facing the ~~country~~ today?”
European Union



Surveys often used to measure issue importance

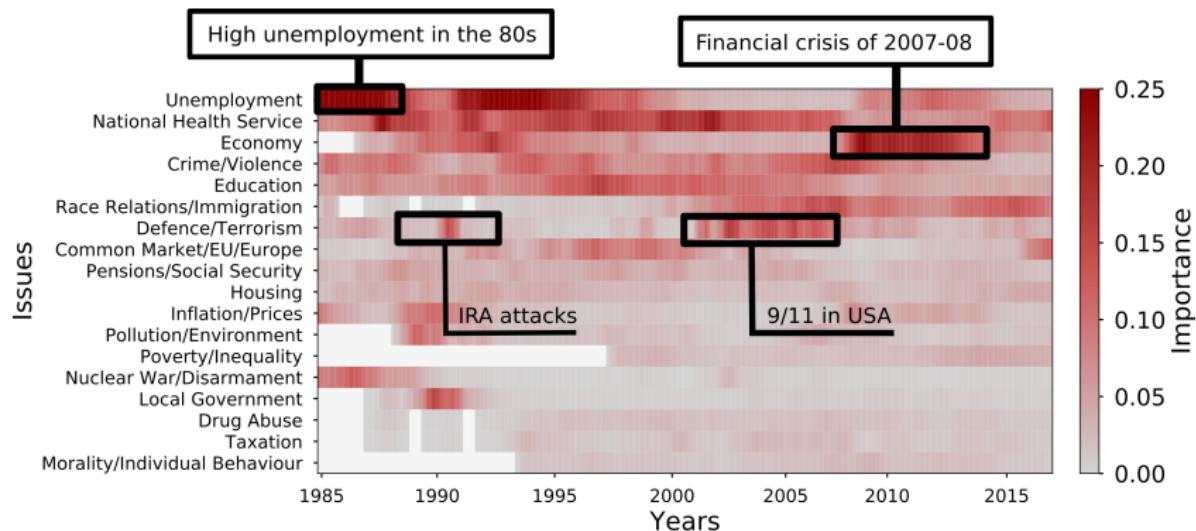


Figure Public attention to policy issues by month from 1985 to 2016 as reported by representative surveys of the UK population collected by Ipsos MORI.

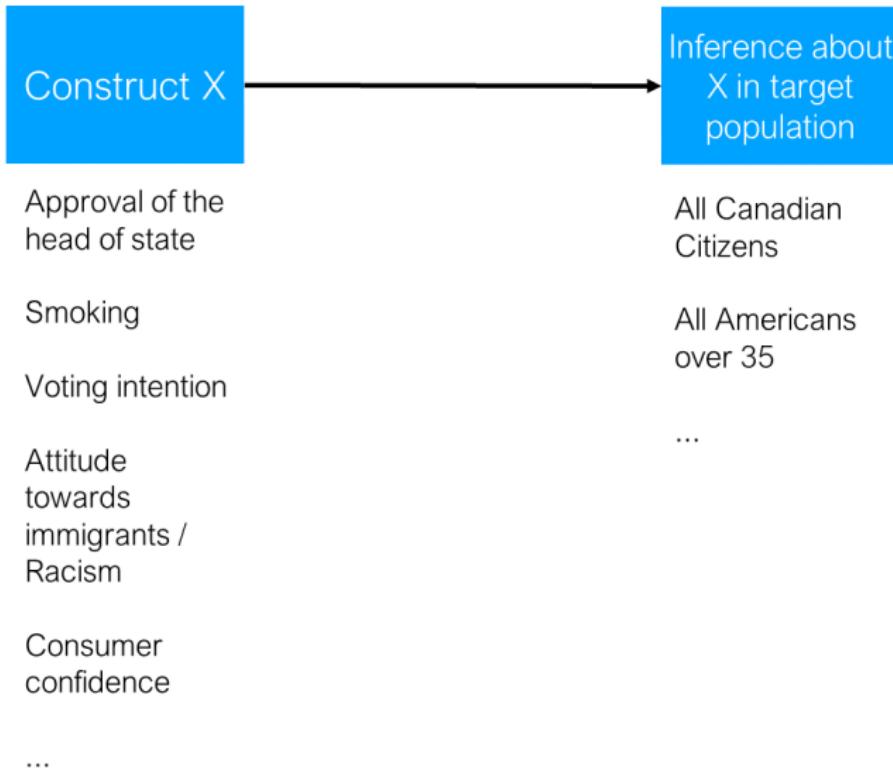
Could social media measure issue importance?

Surveying has been the gold standard, but observation on social media is cheap and easy...

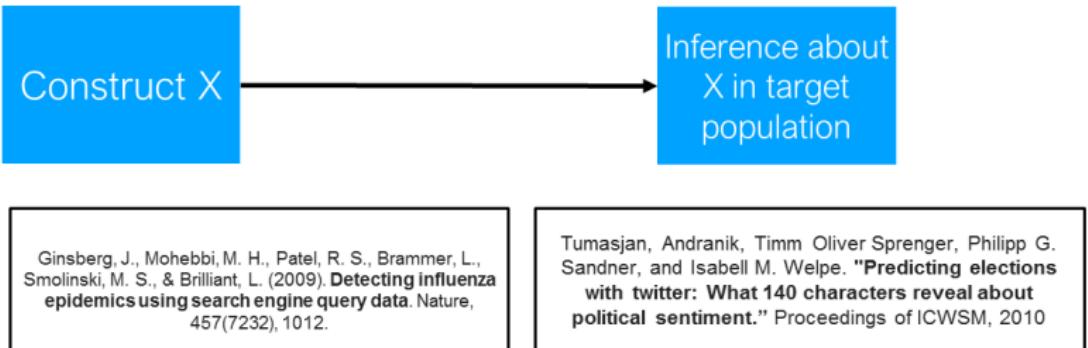
	Survey	Social Media
Purposefully sampled	✓	✗
Known demographics	✓	✗
Easy to get	✗	✓
Cheap	✗	✓
Quick reaction to events	✗	✓
Available retrospectively	✗	✓
Social desirability	✓	✓

Can we use social media data to measure public attention to societal issues?

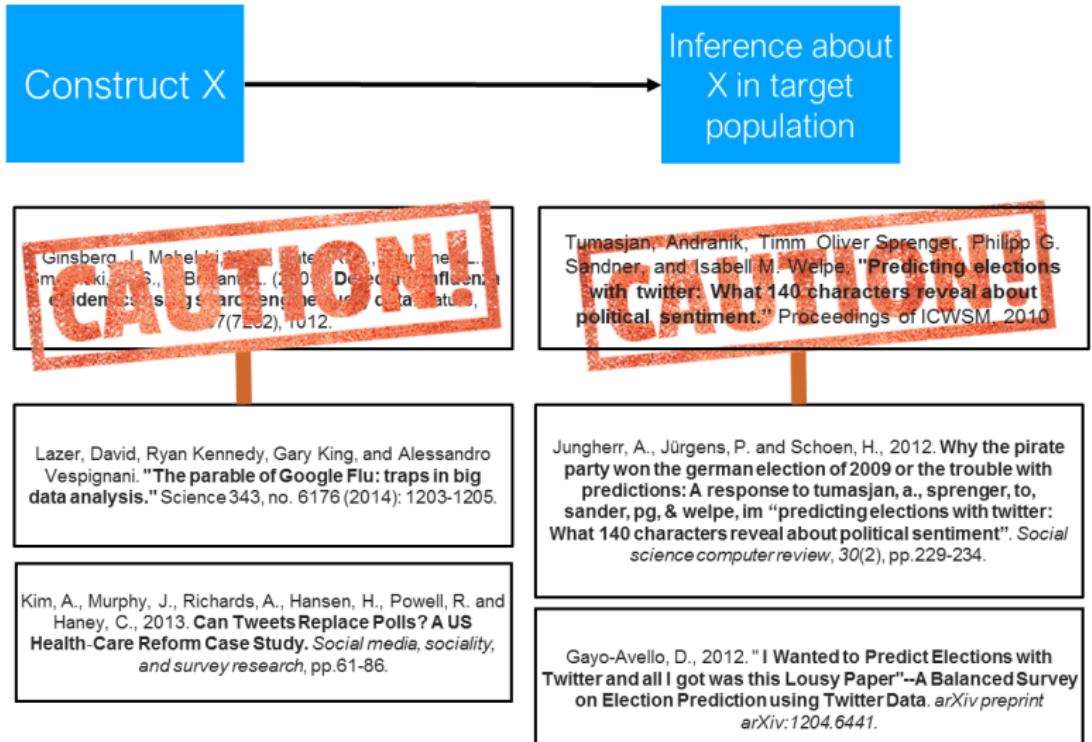
Typical survey research design

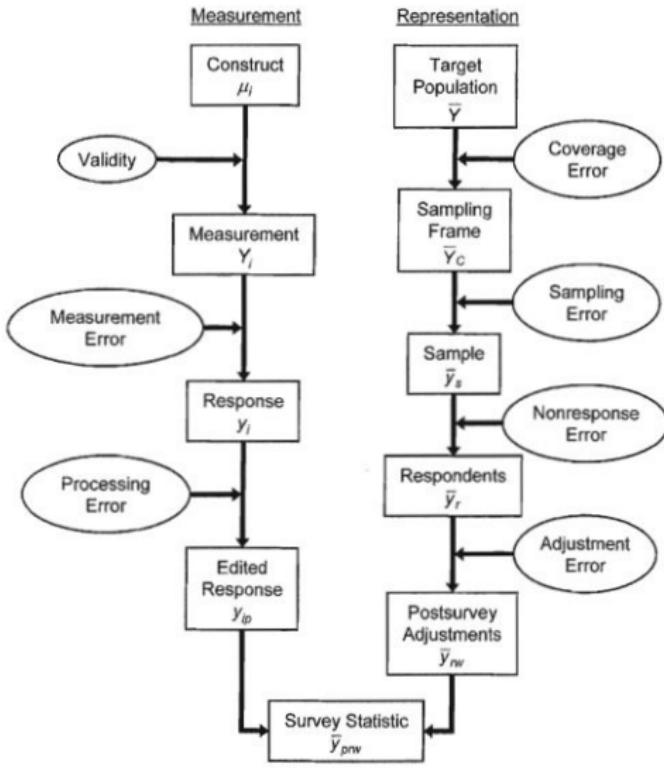


Using social media to measure X

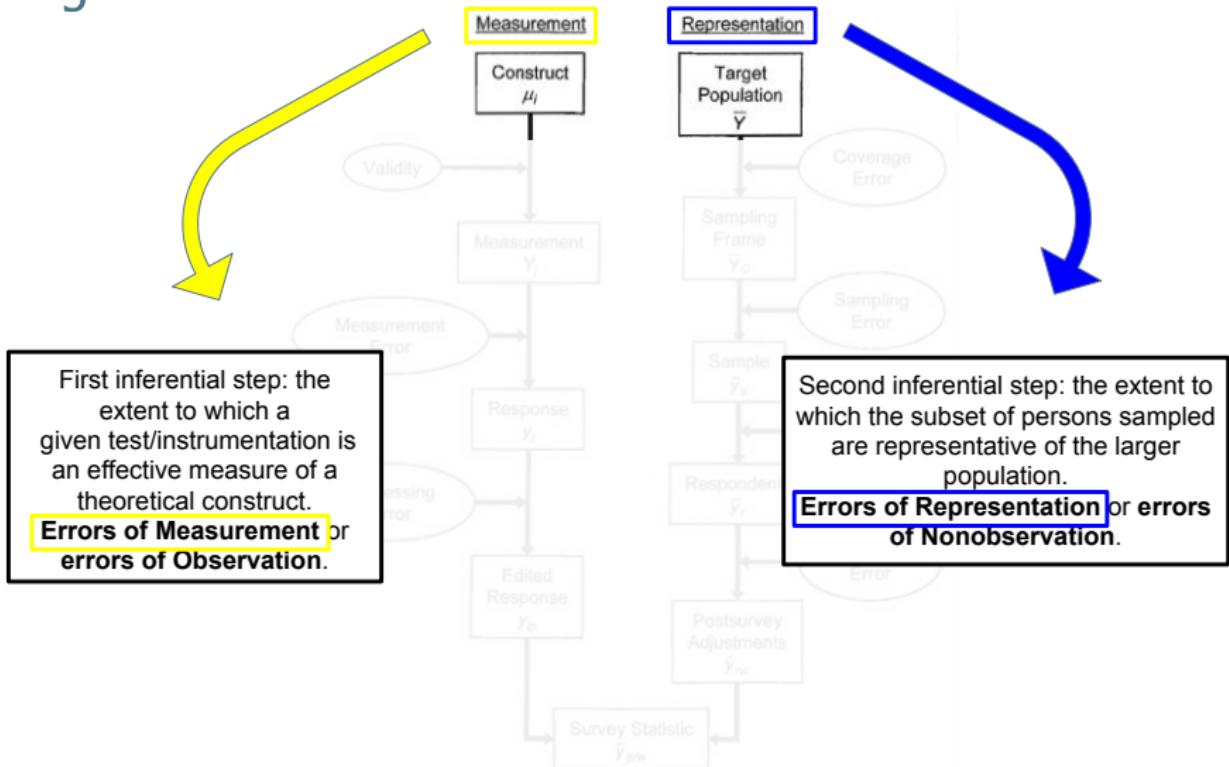


Using social media to measure X



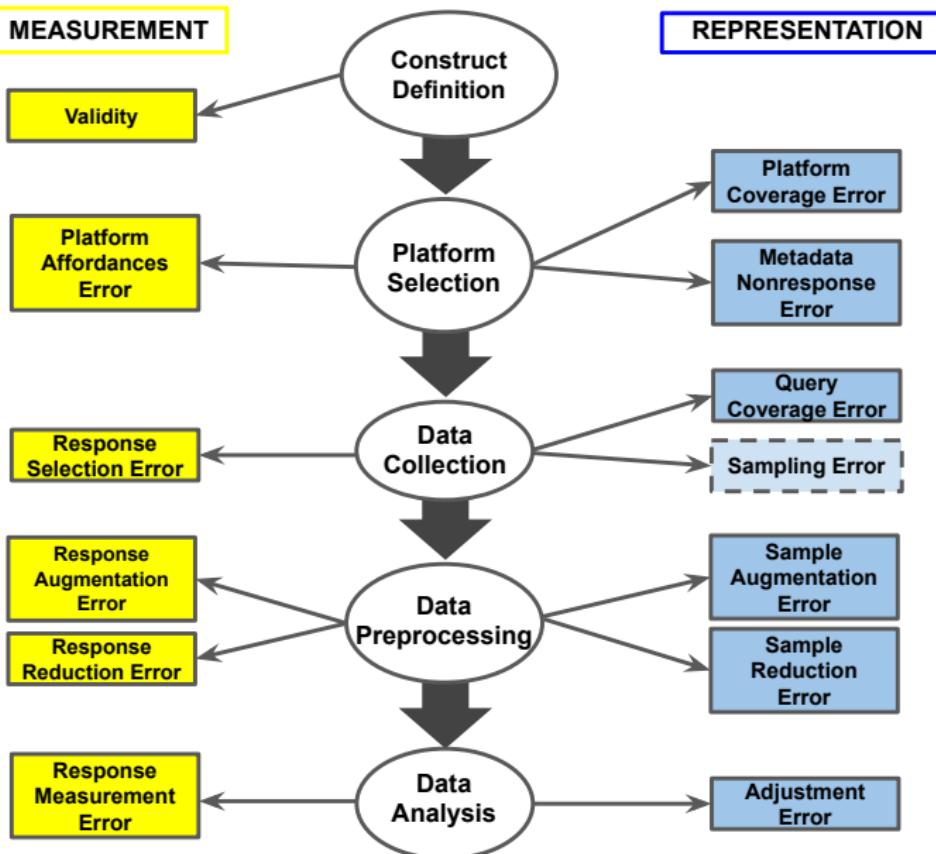


Total Survey Error Framework: Survey Lifestyle
from a Quality perspective. Adapted from Groves et
al.(2011) P. 48



Total Survey Error Framework: Survey Lifestyle
from a Quality perspective. Adapted from Groves et
al.(2011) P. 48

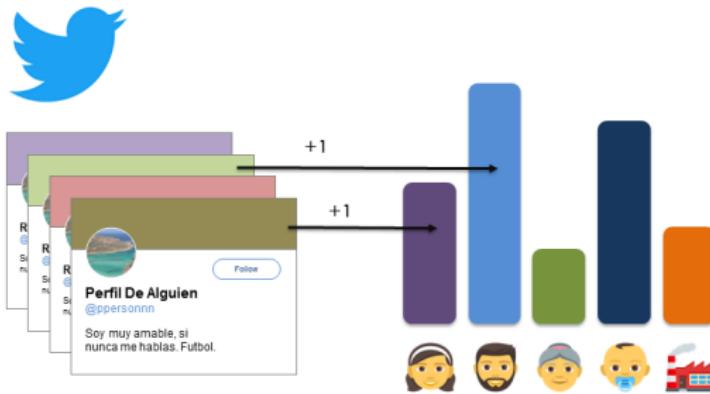
The Digital Trace Data Error Framework



How might debiasing work in practice in social media?

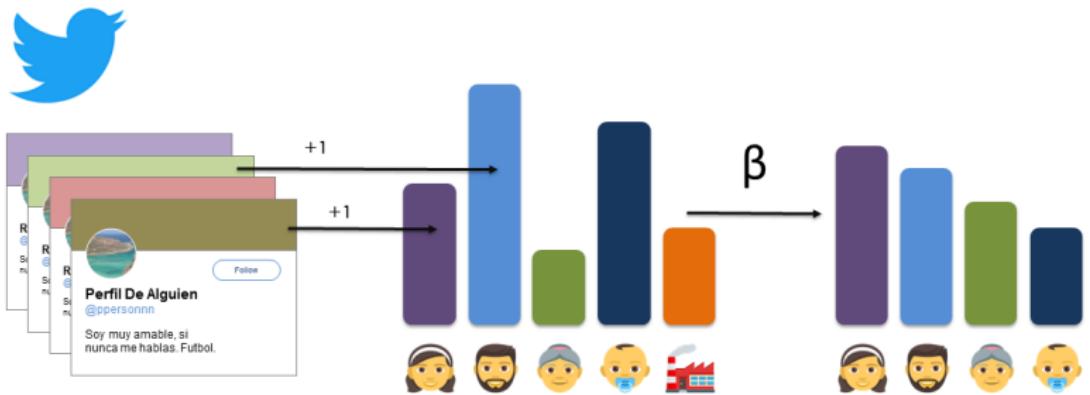


How might debiasing work in practice in social media?



Twitter population,
non-representative

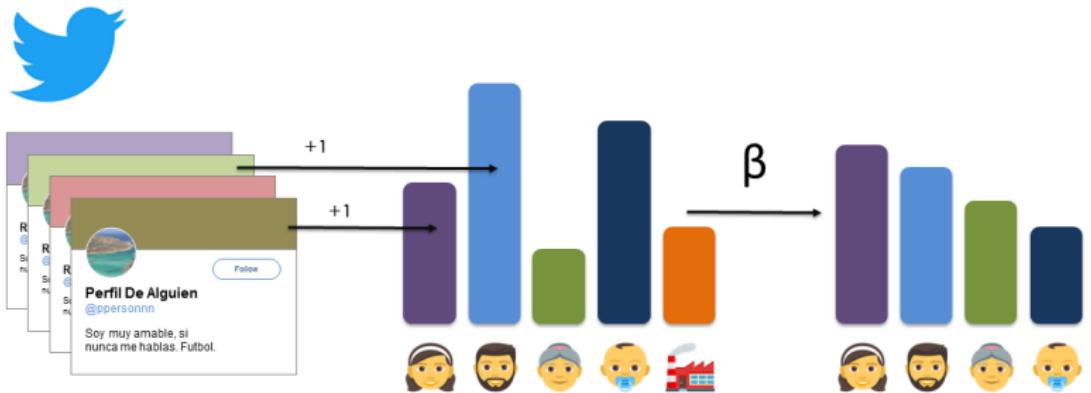
How might debiasing work in practice in social media?



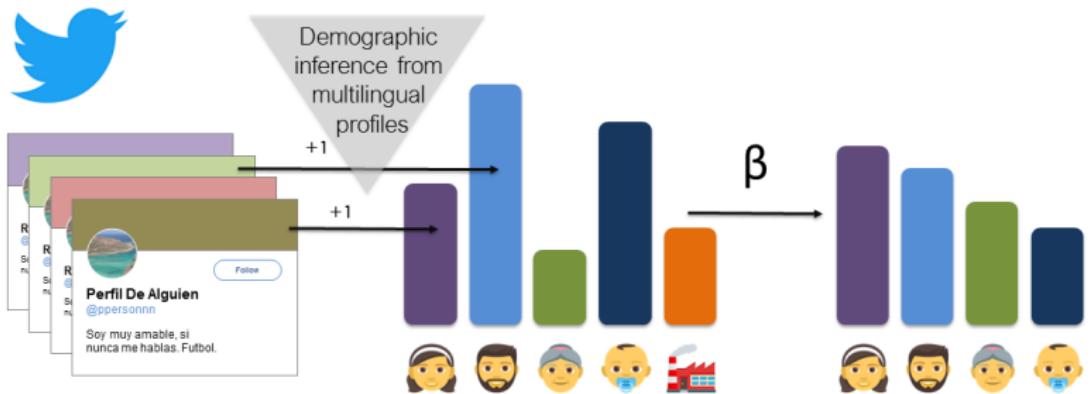
Twitter population,
non-representative

More representative
sample for geographic
region

Poststratification requires the ability to first stratify!



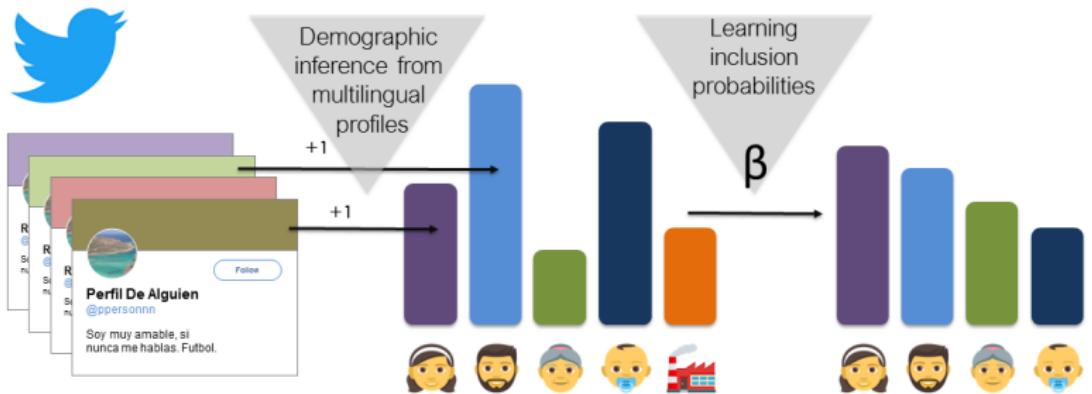
Poststratification requires the ability to first stratify!



Twitter population,
non-representative

More representative
sample for geographic
region

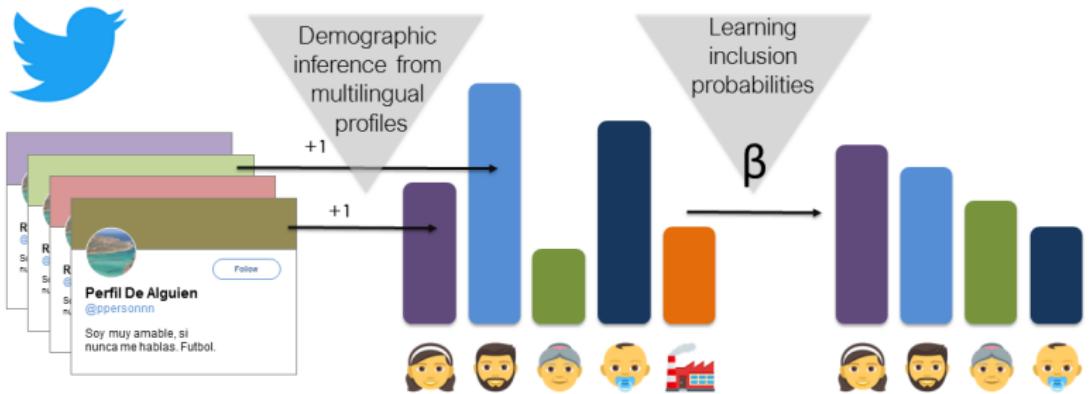
Poststratification requires the ability to first stratify!



Twitter population,
non-representative

More representative
sample for geographic
region

Poststratification requires the ability to first stratify!



For 32 languages
(EU-based)

Twitter population,
non-representative
binary gender,
4-category age,
binary is-org status

More representative
sample for geographic
region

Aside: Not all users are human

Which of the following is not from a personal account?

This is it – it's all come down to the last few minutes to change our country #VoteLabour

Unemployment sucks but it sucks more when it affects your friendships. We've all been there bit.ly/1IfGDpW

In swine flu season, sanitizer for the hand of God – fluviruses.com/?p=72702 #Flu #viruses #virus #disease

Depression is not feeling sad once in a while. It's so much more complicated

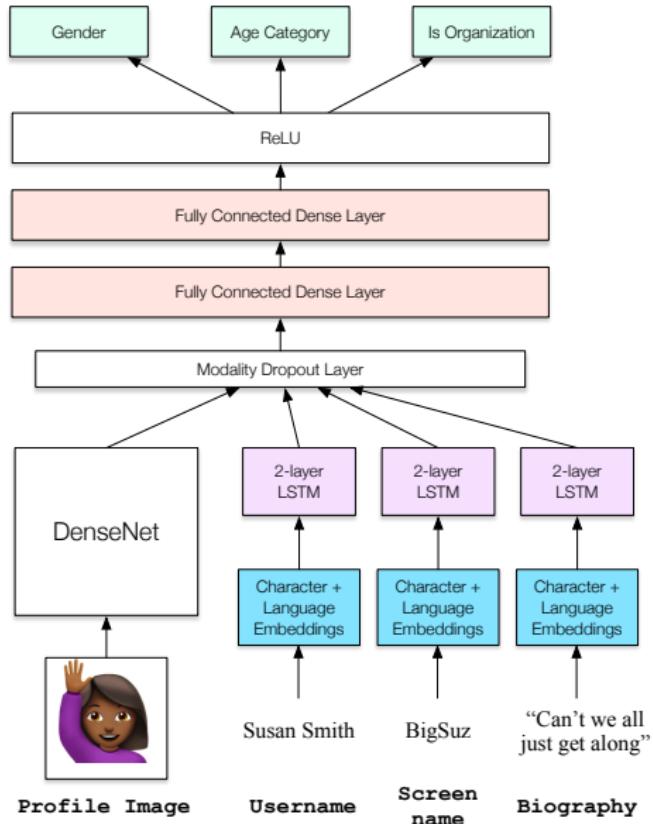
Happy Mother's Day to all the amazing Hockey Moms out there!

Aside: Not all users are human

Nearly 10% of all Twitter accounts belong to organisations. They can play a large role in agenda setting, but cannot participate in human phenomena.

Organizations are Users Too: Characterizing and Detecting the Presence of Organizations on Twitter. James McCorriston, David Jurgens, and Derek Ruths. Proceedings of the 9th International AAAI Conference on Web and Social Media (ICWSM). 2015.

Demographic inference (m3)



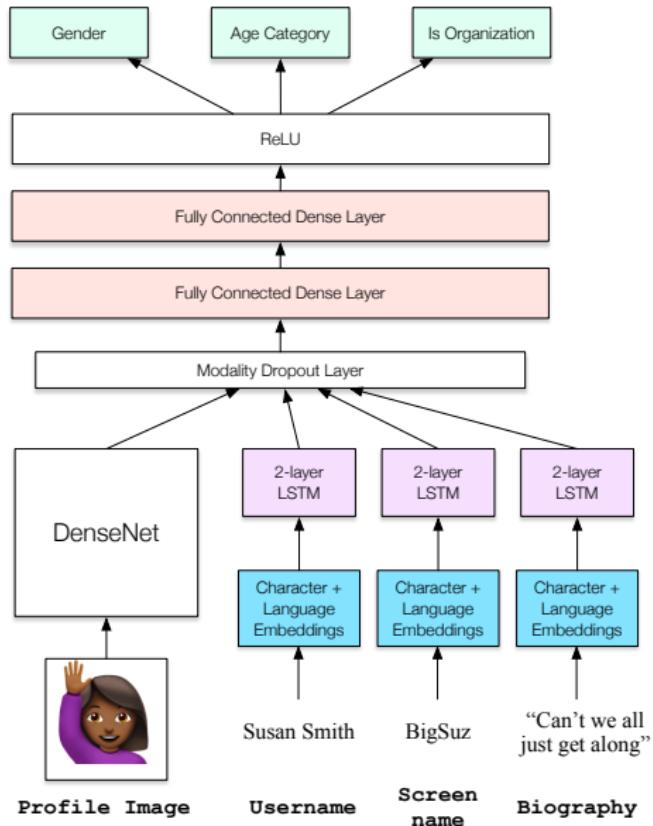
m3 inference:

multimodal Images and text

multiatribute Age, Gender, and
is-organization

multilingual For 32 languages spoken
in Europe

Demographic inference (m3)



Data:

- Twitter (random sample)
- Twitter (hand-curated lists of organisations)
- IMDB/Wikipedia
- Crowdsourcing (evaluation only)

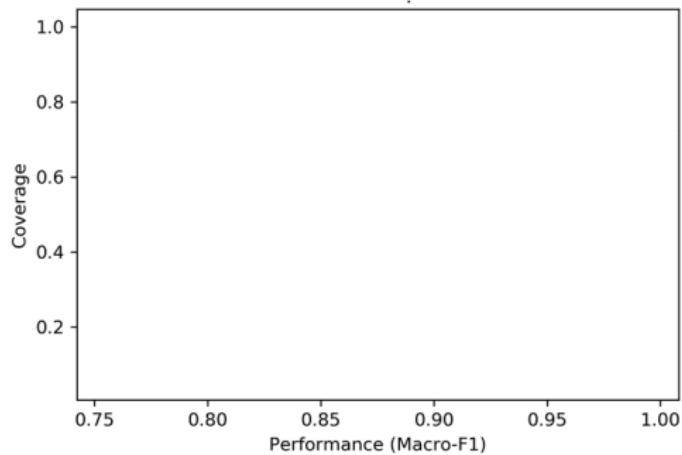
Training and Evaluation Data

- 10% sample of Twitter from 2014 to 2017
 - Heuristically labeled for age/gender and organizations
 - ~40% English
- Headshots of actors in IMDB and biographies in Wikipedia: age and gender
- Cross-language augmentation from EN
- Co-training → expanding data set on large unlabeled twitter corpus

	Gender	Age	Organization
Initial	3.98M	1.20M	59.92K
Total after augment	14.53M	2.61M	23.86M
Held-out Evaluation	0.38M	0.36M	6.91K

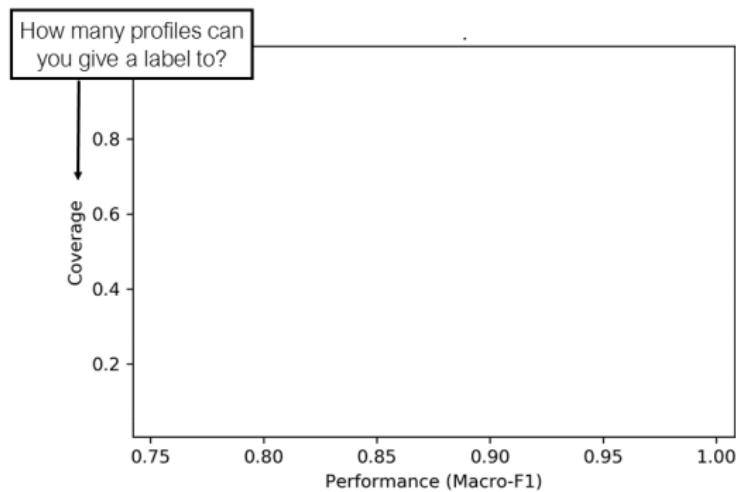
State of the art for gender inference

Gender from Images



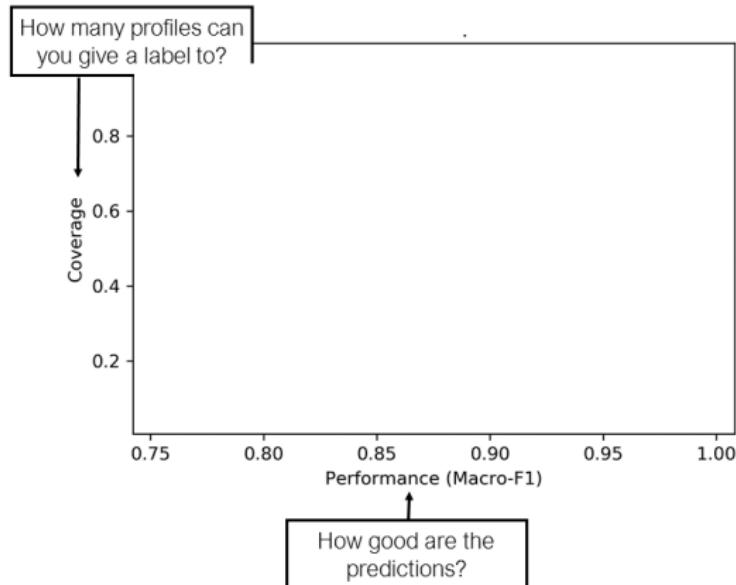
State of the art for gender inference

Gender from Images



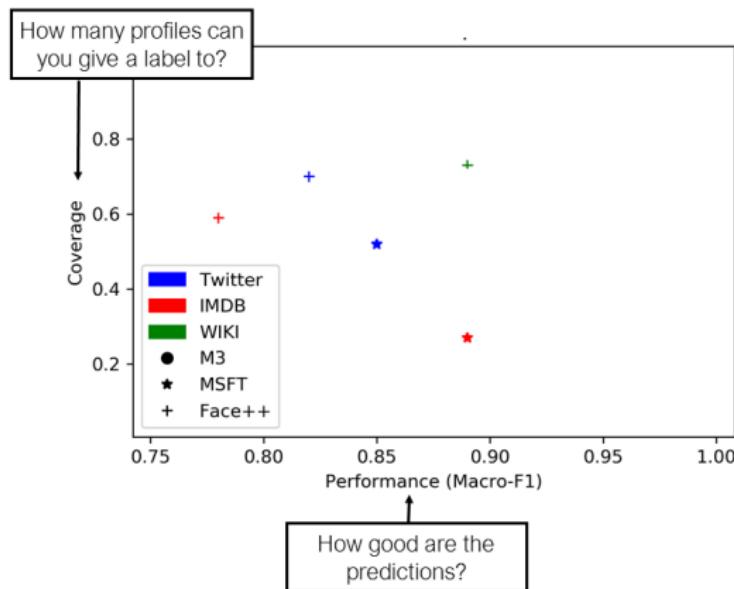
State of the art for gender inference

Gender from Images



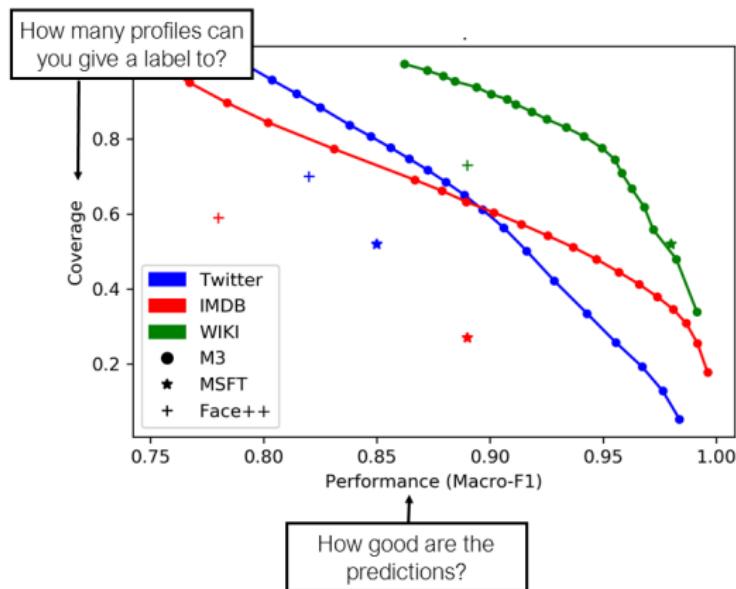
State of the art for gender inference

Gender from Images



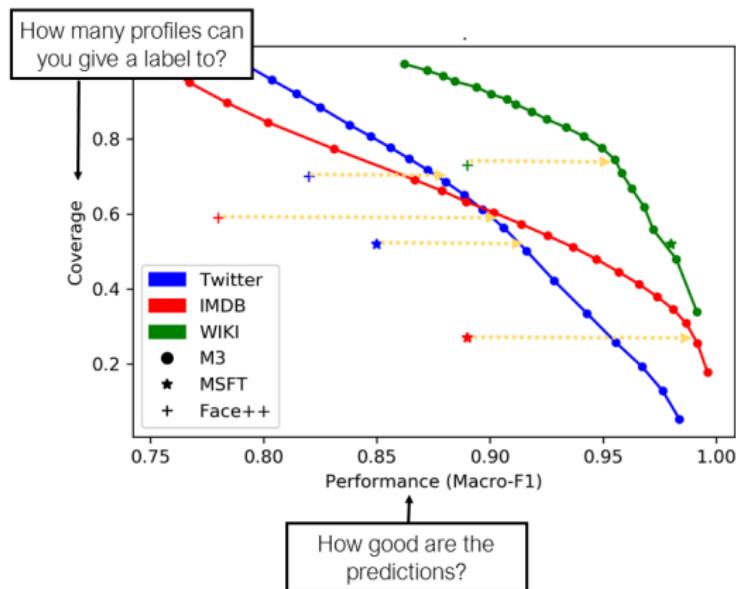
State of the art for gender inference

Gender from Images



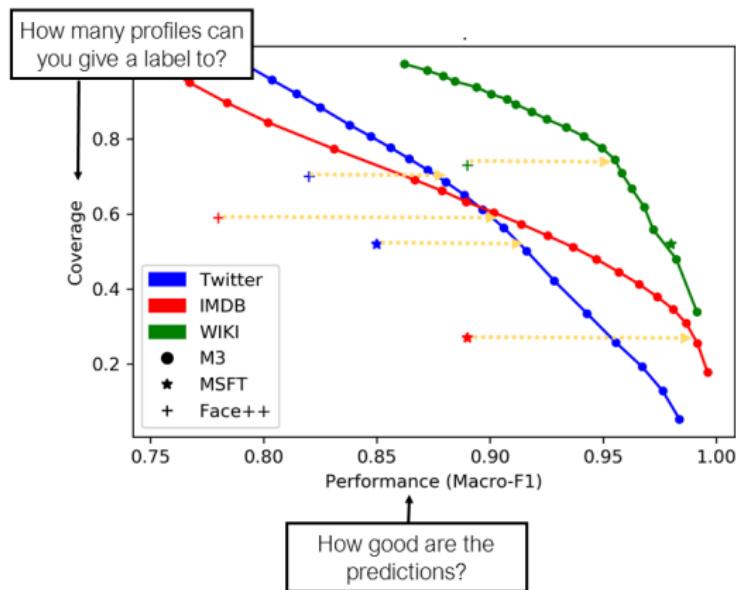
State of the art for gender inference

Gender from Images

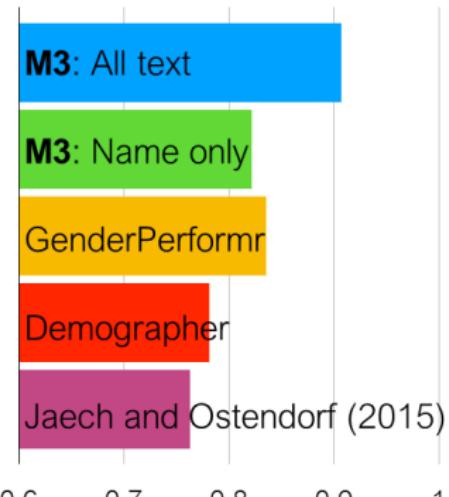


State of the art for gender inference

Gender from Images



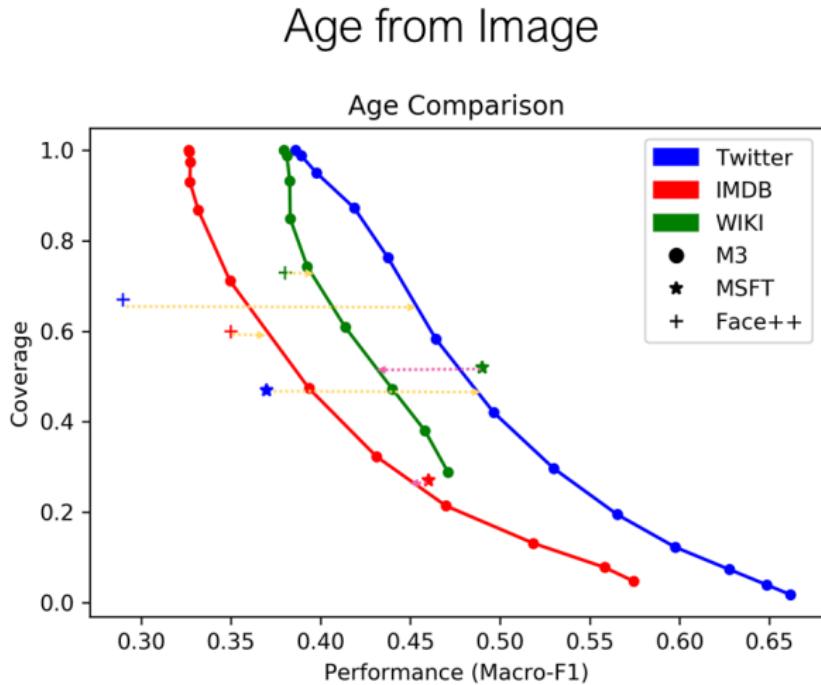
Gender from Text



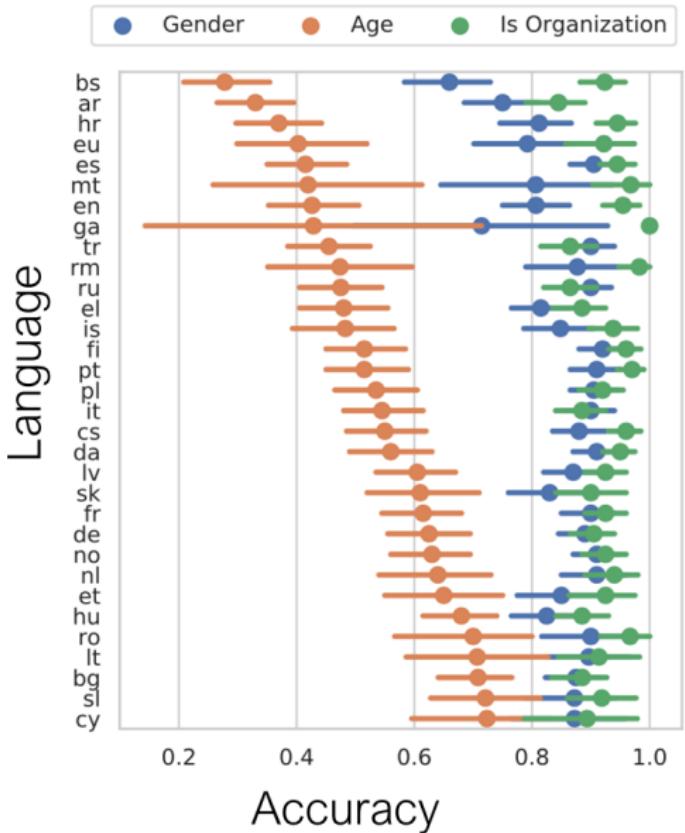
Quick aside: why the increase
in coverage with M3?



State of the art for age inference

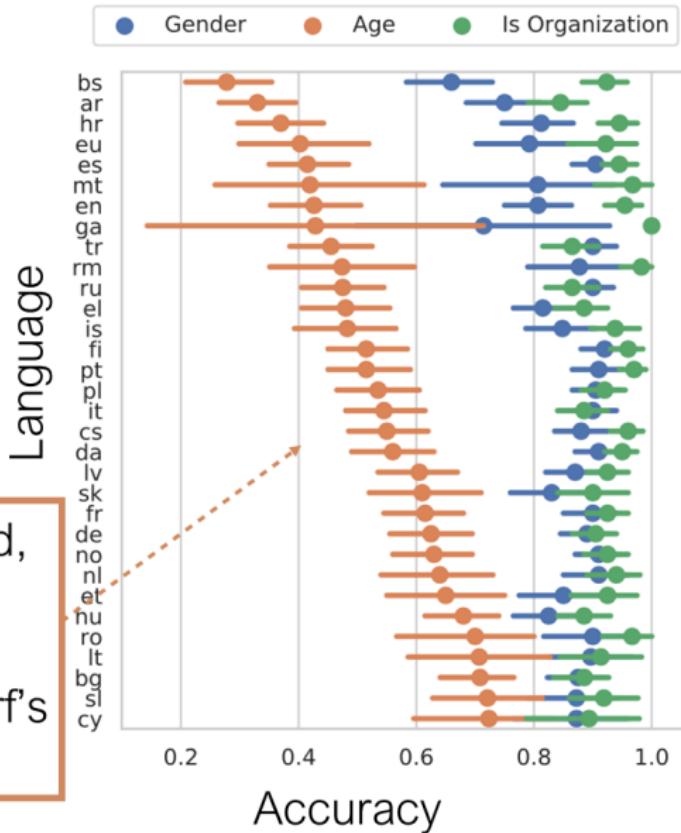


Works pretty well for multiple languages



200 profiles from
32 languages,
3 annotators,
majority voting

Works pretty well for multiple languages



200 profiles from
32 languages,
3 annotators,
majority voting

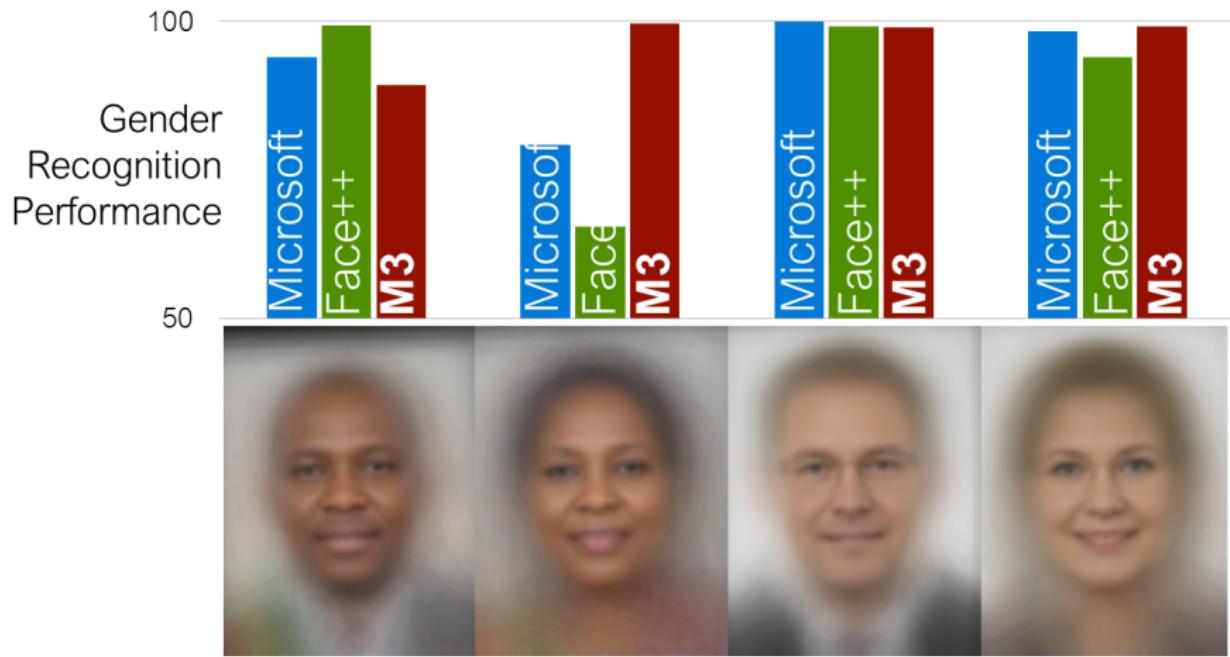
Far less algorithmic bias than current commercial methods

Gender
Recognition
Performance



data from Buolamwini and Gebru (2018)

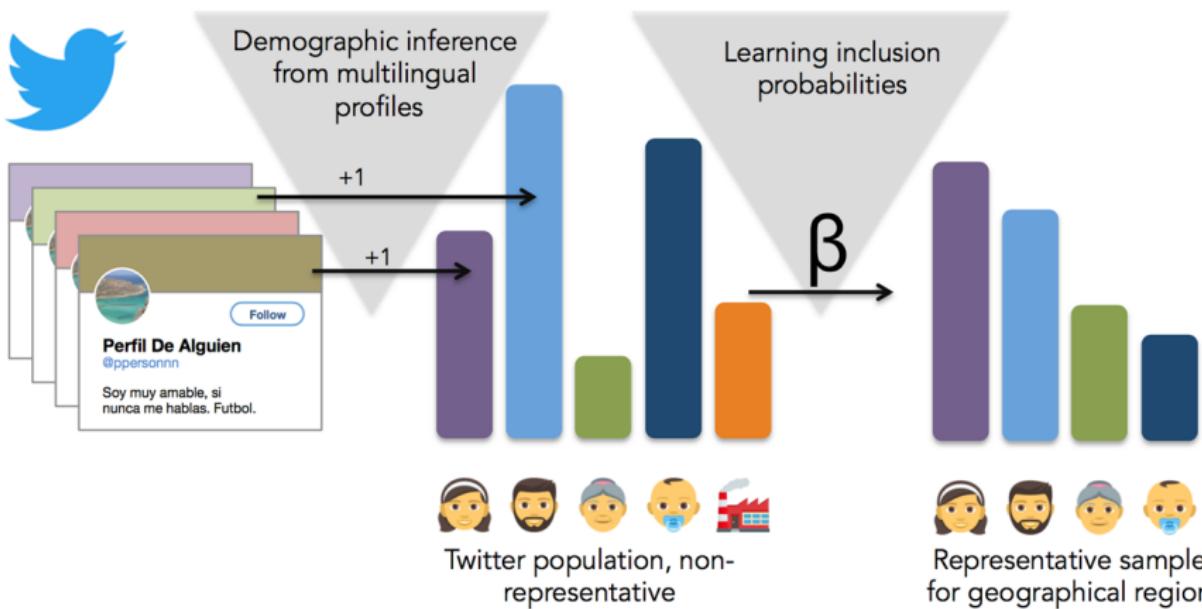
Far less algorithmic bias than current commercial methods



data from Buolamwini and Gebru (2018)

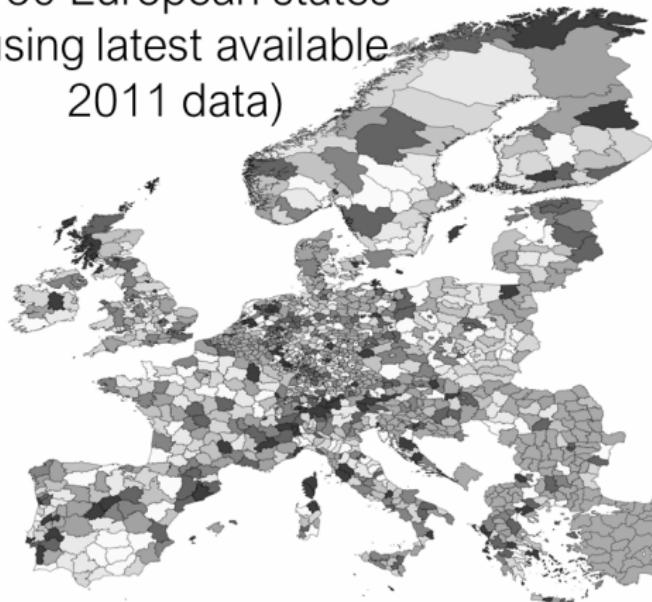
- State-of-the-art results, generally outperforms Face++ and Microsoft Face
- Less biased by skin tones (GenderShades)
- Text-only mode also exceeds comparison systems
- Age is hard, even for humans

Learn the inclusion probabilities for poststratification through population estimation



Population prediction

NUTS3 census regions
in 30 European states
(using latest available
2011 data)



3,202,964 EU-based
Twitter users from
September-December
2016



Population prediction

We run simple regression-based approaches using census data (N) and Twitter (M) for NUTS3 level regions across the EU.

NUTS3 regions range in population from 150k to 800k.

We run leave-one-NUTS3-region-out prediction and measure error with *mean absolute percentage error* (MAPE).

An aside: Location

Online social ties can reveal substantial information about a person

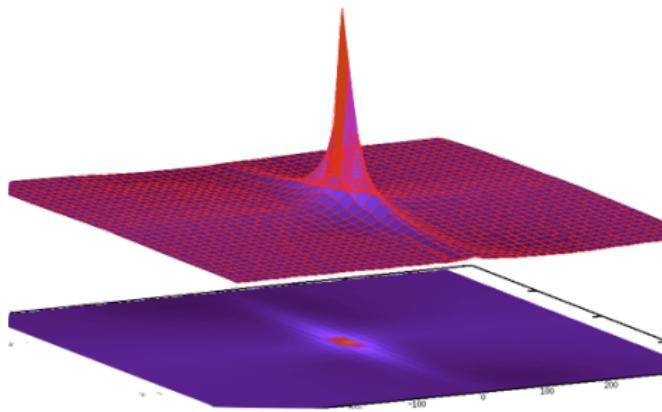
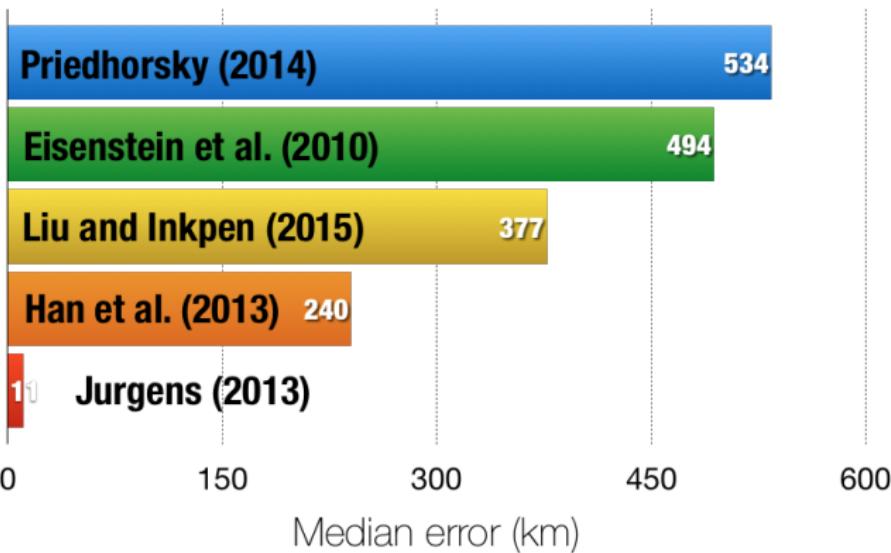


Figure Frequency distribution of where your friends are relative to you

An aside: Location

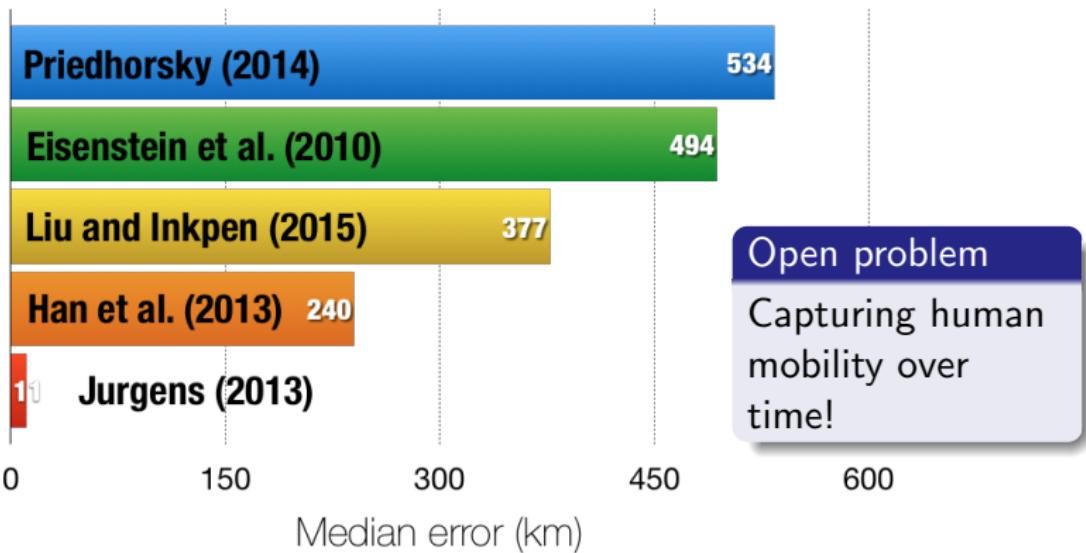
Inferring location on Twitter



That's what friends are for: Inferring location in online communities based on social relationships. David Jurgens. Proceedings of the 7th International AAAI Conference on Weblogs and Social Media (ICWSM) 2013.

An aside: Location

Inferring location on Twitter



That's what friends are for: Inferring location in online communities based on social relationships. David Jurgens. Proceedings of the 7th International AAAI Conference on Weblogs and Social Media (ICWSM) 2013.

Models

$N \sim M$ is the base model using only total population from the census (N) and Twitter (M)

$N \sim \sum_g M(g)$ uses gender marginal counts only (i.e., the total counts of males and females not broken down by ages)

$N \sim \sum_a M(a)$ uses age marginal counts only

$N \sim \sum_{a,g} M(a,g)$ uses the joint histograms inferred from Twitter but only the total population values from the census

$\log N(a,g) \sim \log M(a,g) + a + g$ uses the joint histograms inferred from Twitter and the joint histograms from the census

Results

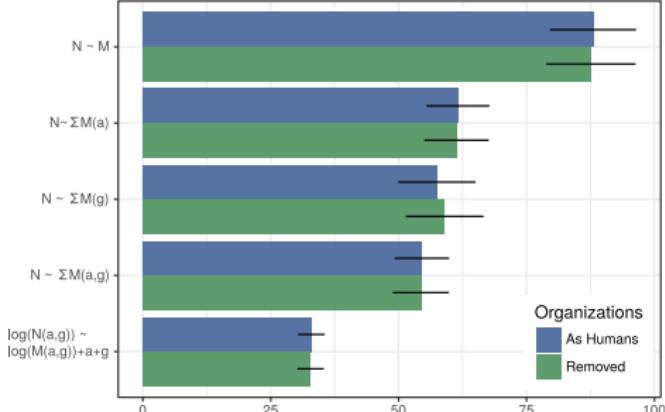
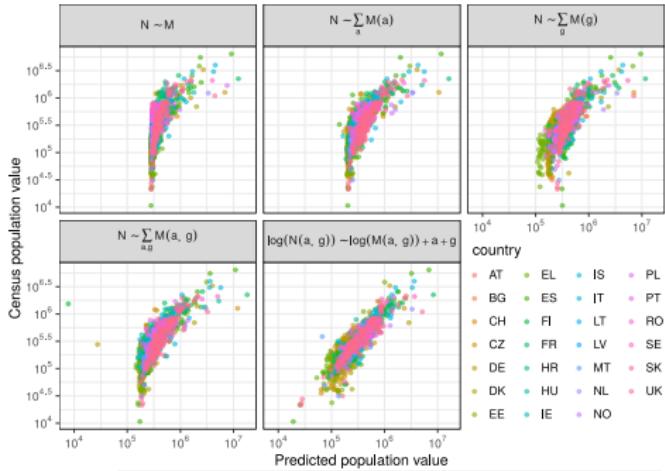
$N \sim M$ is the base model using only total population from the census (N) and Twitter (M)

$N \sim \sum_g M(g)$ uses gender marginal counts only (i.e., the total counts of males and females not broken down by ages)

$N \sim \sum_a M(a)$ uses age marginal counts only

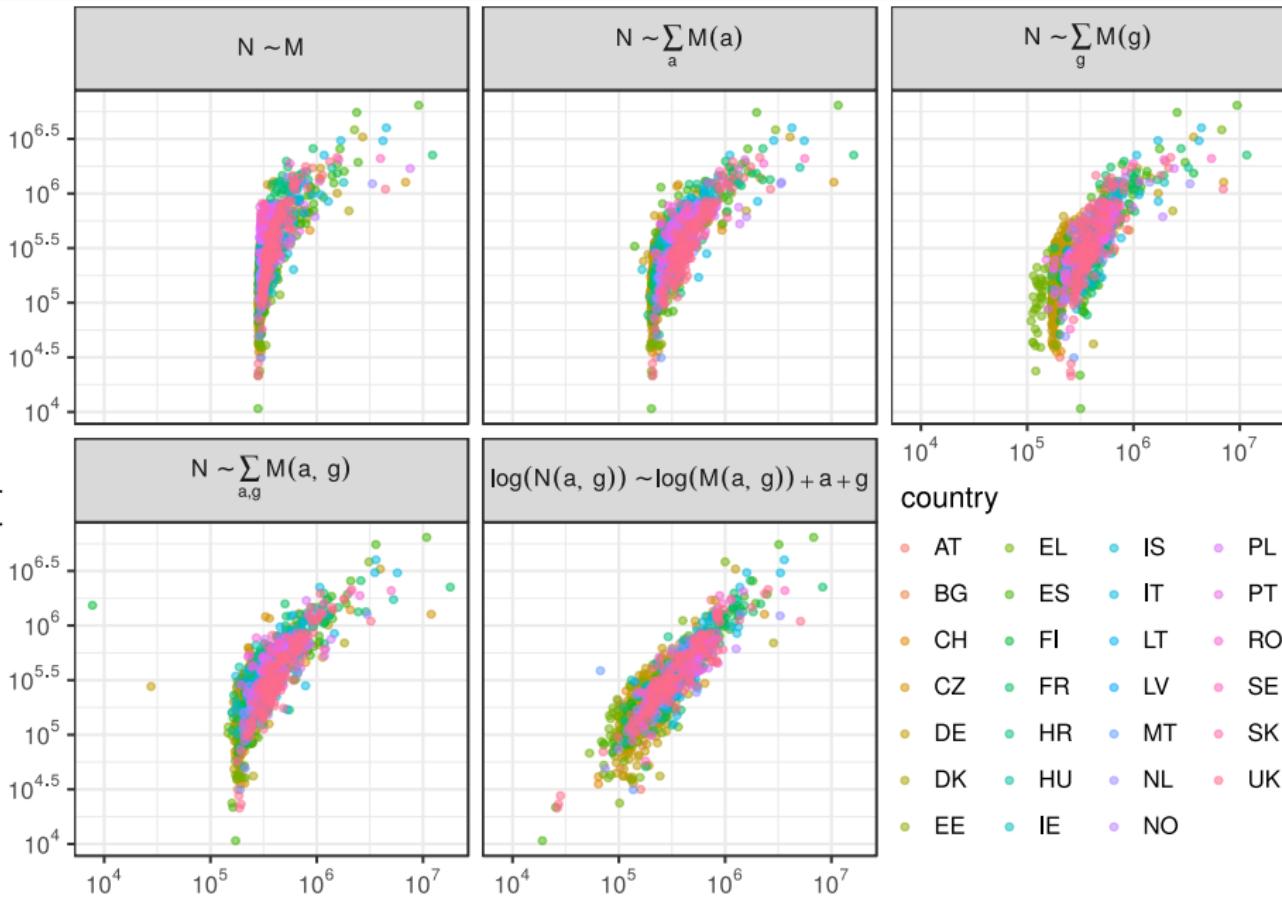
$N \sim \sum_{a,g} M(a,g)$ uses the joint histograms inferred from Twitter but only the total population values from the census

$\log N(a,g) \sim \log M(a,g) + a + g$ uses the joint histograms inferred from Twitter and the joint histograms from the census

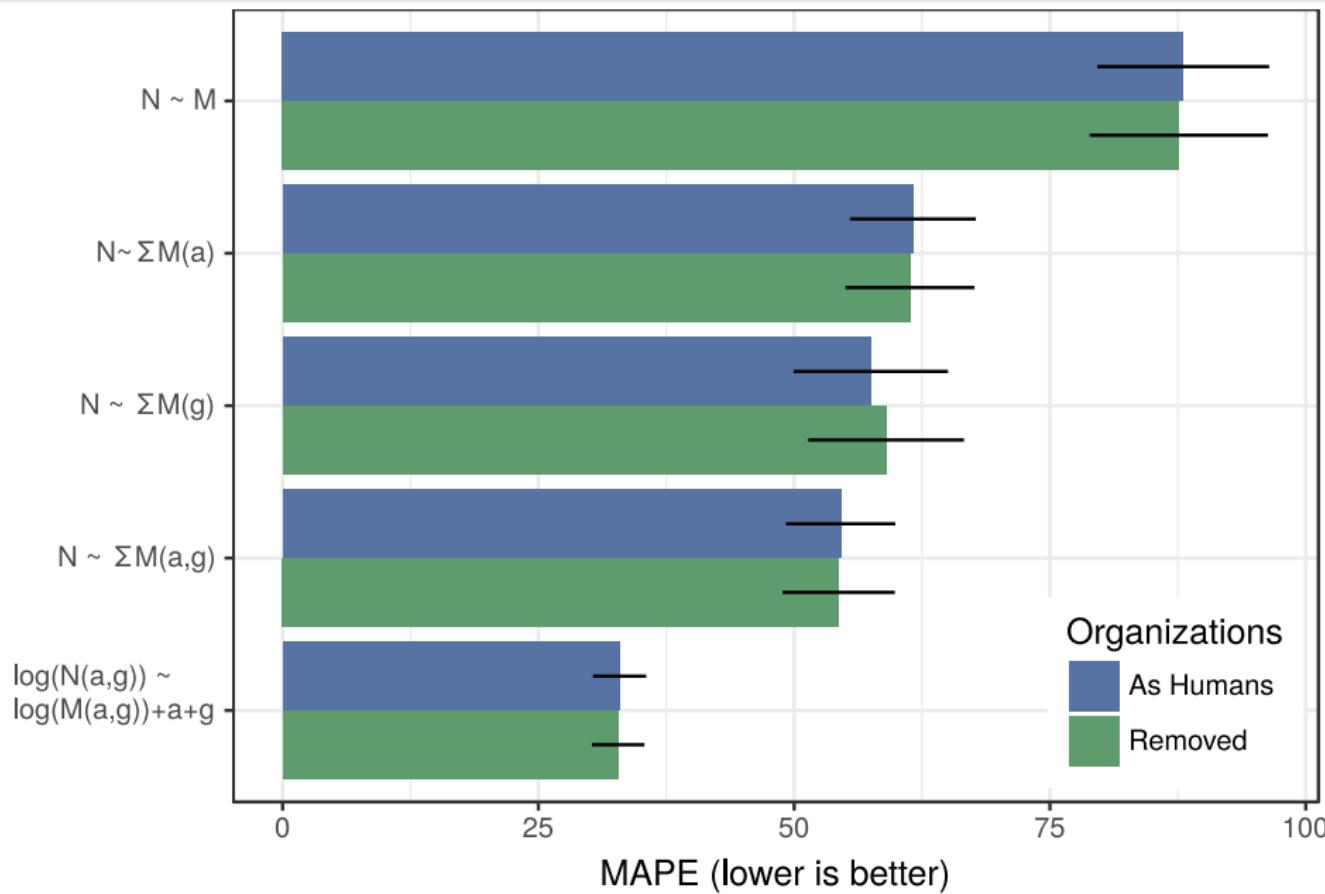


Results

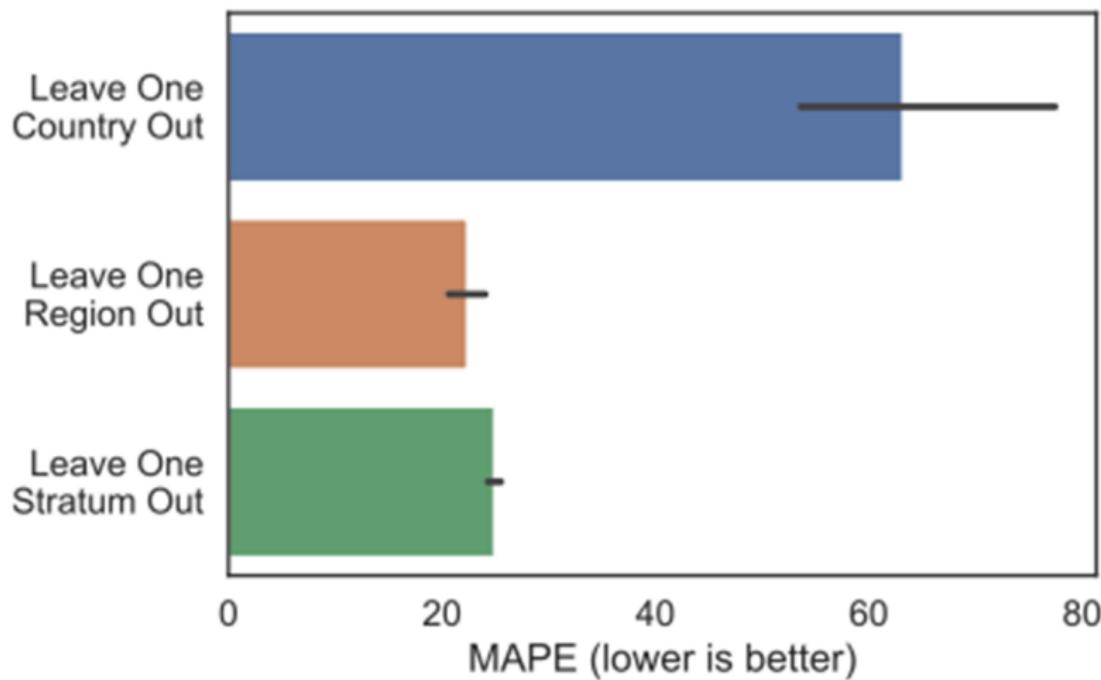
Census population value



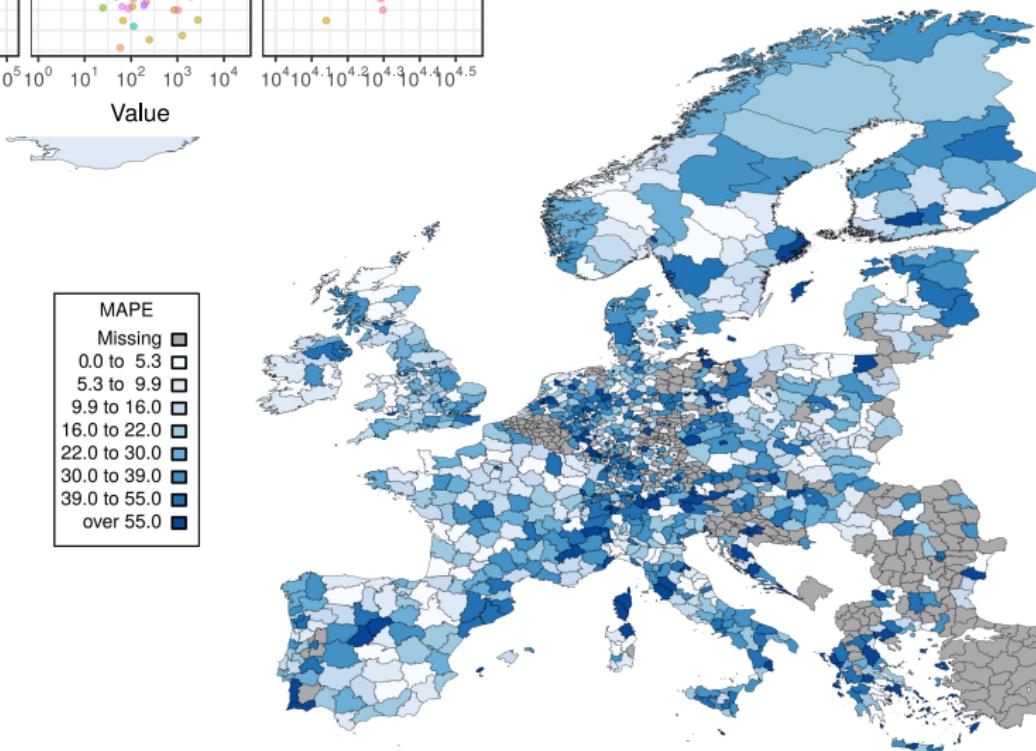
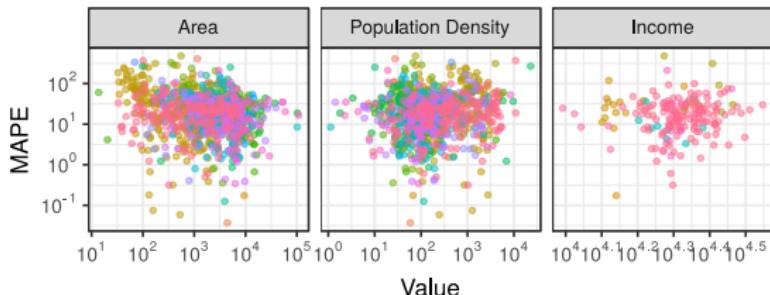
Results



Inclusion weights generalize to population estimation on held-out data



Sources of error



What have we learnt?

- Possible to infer demographics/org-status fairly reliably in multilingual settings

What have we learnt?

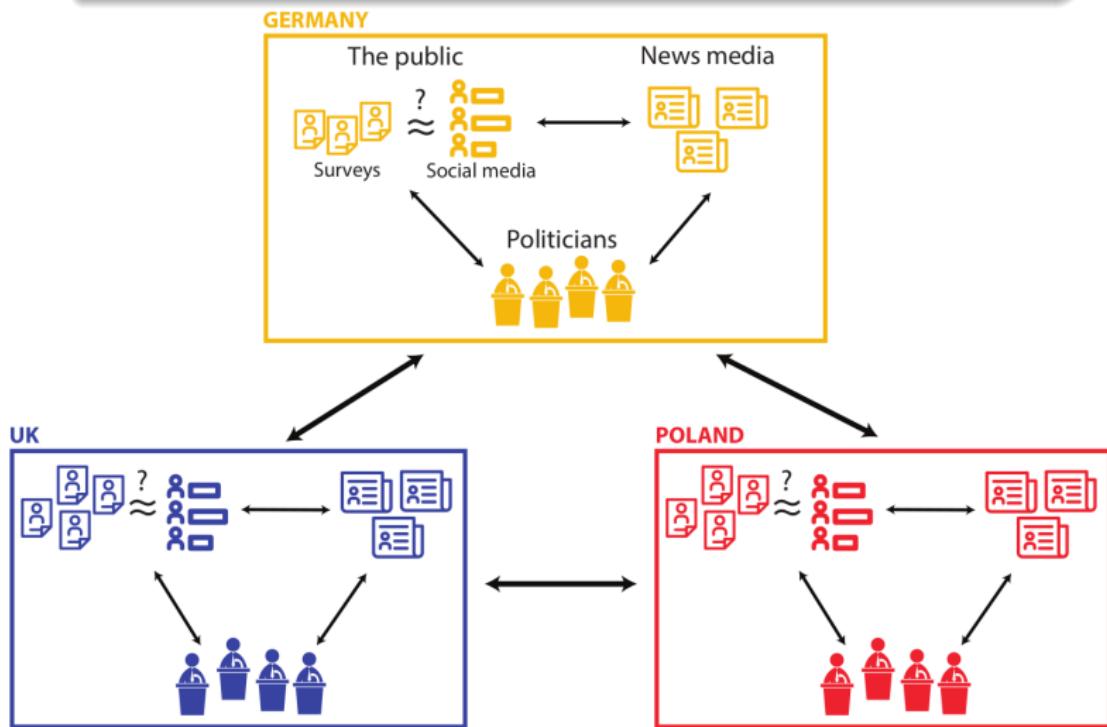
- Possible to infer demographics/org-status fairly reliably in multilingual settings
- Reweighting helps significantly, but...
 - Heterogeneity in social media inclusion probabilities that needs to be accounted for, especially in multicountry setting
 - More analysis needed to understand how reweighting helps with specific dependent variables

What have we learnt?

- Possible to infer demographics/org-status fairly reliably in multilingual settings
- Reweighting helps significantly, but...
 - Heterogeneity in social media inclusion probabilities that needs to be accounted for, especially in multicountry setting
 - More analysis needed to understand how reweighting helps with specific dependent variables
- Open problems
 - Capturing human mobility over time
 - Non-independence of social media observations (e.g., homophily)
 - How to best include probability output of M3 inference into downstream tasks

Back to measuring the attention to societal issues

Comparing political, media, and public 'agendas' across languages and countries within the EU, over time.





Demographic Inference and Representative Population Estimates from Multilingual Social Media Data

Zijian
Wang



Scott
Hale



David
Adelani



Przemyslaw
Grabowicz



Timo
Hartmann



Fabian
Flöck



David
Jurgens



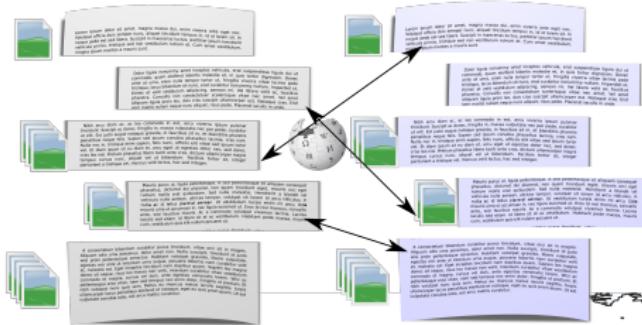
New resources for the community

- M3 model as off-the-shelf python library
 pip install m3inference
<https://github.com/euagendas/m3inference>
- Inclusion probabilities for Europe and eventually other regions
<http://euagendas.org/inclusionprobs>
- Live demo of M3 here: <http://m3.euagendas.org/>
- Everything is at <https://github.com/euagendas/>

Measuring the media agenda



Weekly/nightly news across EU+Russia



Segment
and match
news items

Compare similarity
of news within and
between countries



Measuring the media agenda

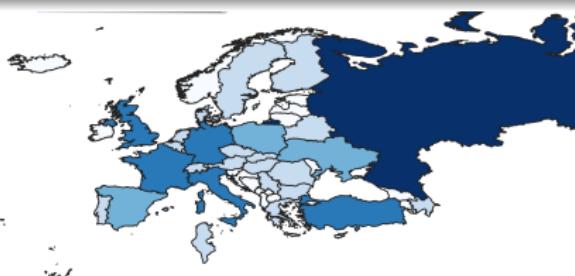


Weekly/nightly news across EU+Russia

Cross-national media agenda similarity

- Shared language
- Geographic proximity
- Past migration
- Economic similarity
- Temporal factors (e.g., GDP)

Compare similarity
of news within and
between countries



Thank you!

THE ALAN
TURING
INSTITUTE



VolkswagenStiftung



Scott Hale



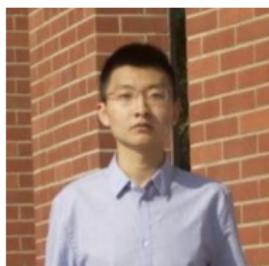
David Jurgens



Fabian Flöck



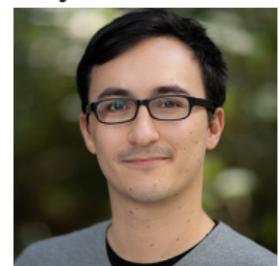
Przemyslaw Grabowicz



Zijian Wang



David Adelani



Timo Hartmann

Chico Camargo

Research questions

- RQ1 What differences exist in the media, political, and public agendas across nations and languages and how do they evolve over time in comparison to each other?
- RQ2 To what extent do social media reflect a broader public agenda?
- RQ3 How do the agendas of actors in a specific national or linguistic context influence how agendas are adopted in other contexts, and what explains these influences?