

Spatial Datasets And Urban Applications

AAAI ICWSM TUTORIAL, MAY 15, 2017 MONTREAL, CANADA

<https://github.com/bmtgoncalves/ICWSM17>

Theory Session:

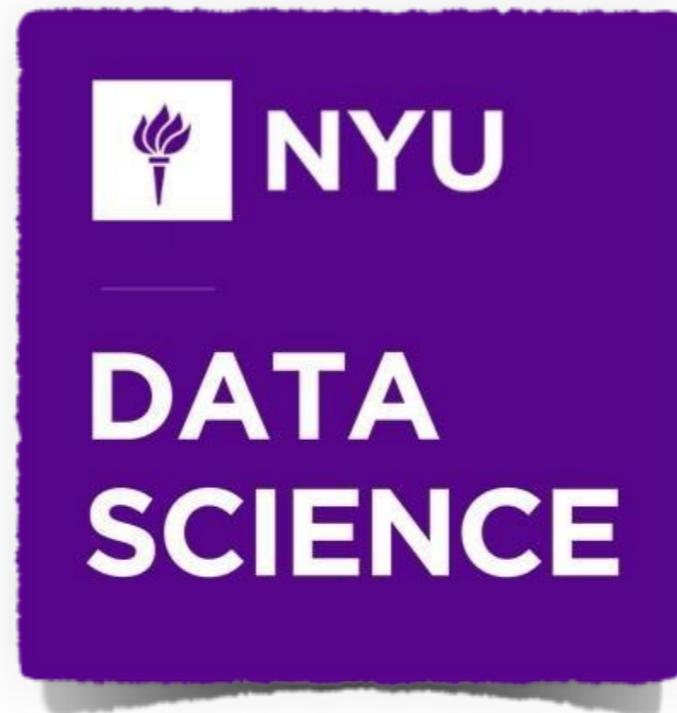
- **8:30 - 9:30** - *Bruno Gonçalves* - We will start with an overview of geolocated APIs such as those of Twitter and Foursquare, using Python focusing on their similarities and differences. Applications covered will range from visualizing Twitter usage and language geography , Inter-City Mobility through airline connections and the Global Language Network.
 - **9:30 - 10:30** - *Anastasios Noulas* - In this part of the tutorial attendees will be provided with an overview perspective on location-based technologies with a focus on urban applications. From human mobility models and theory, to neighborhood detection and finding an appropriate location for placing a shop or amenity attendees will be introduced to the mechanics and building blocks of a number of research works with a strong application orientation. Next, the spotlight will be put on newer advancements in the area of the so-call gig economy, introducing services like Uber and the mechanics and building blocks of a number of research works with a strong application orientation. Next, the spotlight will be put on newer advancements in the area of the so-call gig economy, introducing services like Uber and Airbnb, highlighting opportunities for research using corresponding transport and hospitality datasets.
 - **10:30 - 11:00** *Coffee break*
 - **11:00 - 12:00** - *Desislava Hristova* - Overview of the theory of social capital in cities and the social role of places. This will include an introduction to structural holes and open network structures in light of social capital and how this can be extended to places in cities. The role of social media will be discussed as a proxy for understanding urban development and gentrification in cities like London and New York. Concrete examples of measuring the digital footprints of wellbeing in cities and beyond will be provided alongside the theory of cultural capital and urban development.
 - **12:00 - 13:30** - *Lunch*
-

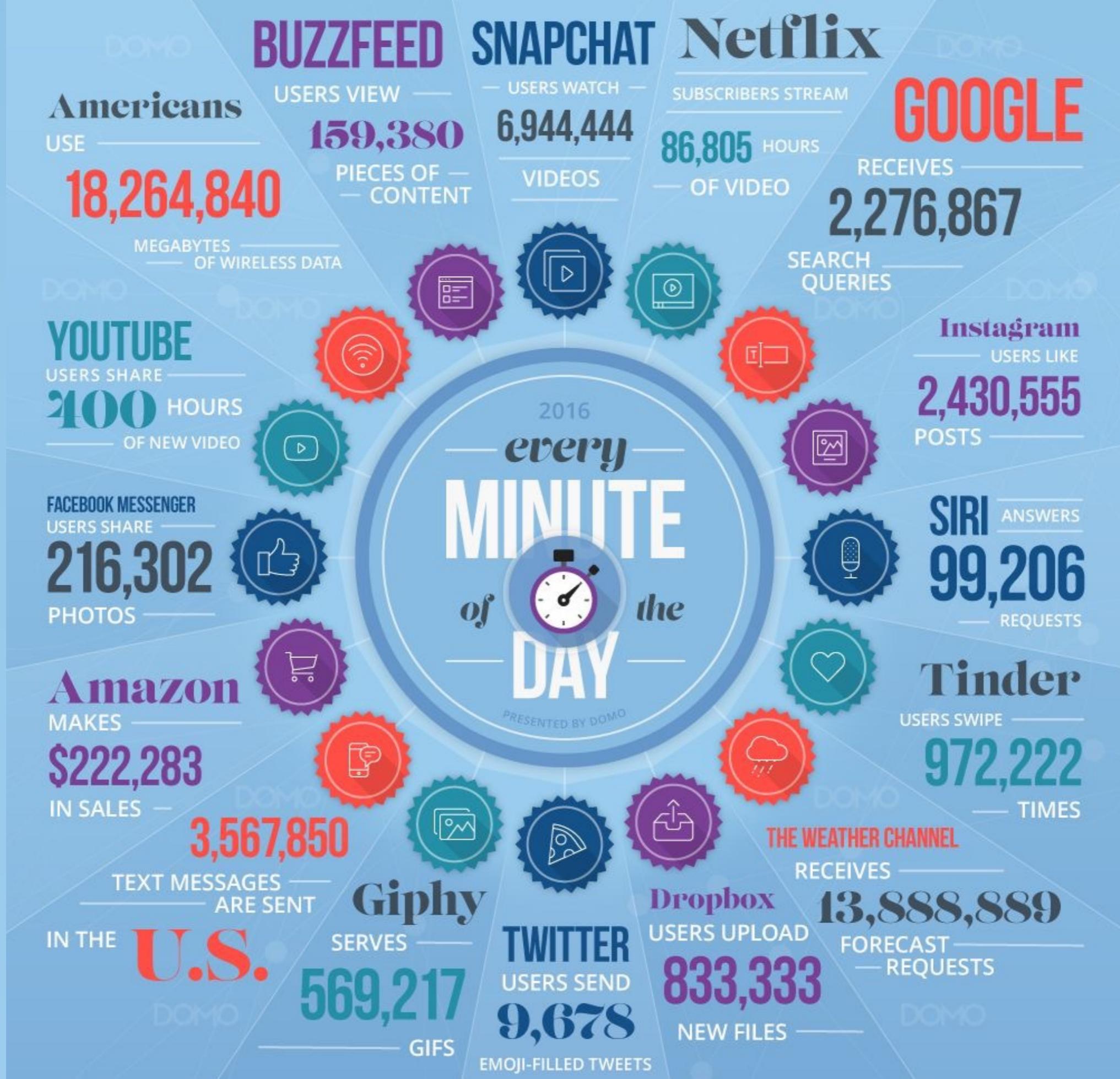
Practical Session:

- **13:30 - 14:30** - *Bruno Gonçalves* - After an introduction on how to collect data using the Twitter and Foursquare APIs we will proceed to perform several simple analyses to illustrate different important techniques and libraries. Using the airline transportation network (from the Bureau of Transportation Statistics) as a starting point, we will build a simple model of large scale mobility within the United States.
- **14:30 - 15:30** - *Anastasios Noulas* - In this part of the tutorial Anastasios will first focus on foundational methods for spatial retrieval using Python. Code snippets that allow for the organisation of spatial data, the calculation of geographic distances as well as the retrieval of points within a radius, given a pair of coordinates, will be examined in detail. Next, the Networkx library will be employed for the analysis of networks of places, whereas the session will close with a demo showcase of Uber's deck.gl framework for visualisation of spatial data.
- **15:30 - 16:00** *Coffee break*
- **16:00 - 17:00** - *Desislava Hristova* - Understanding Urban Development and Gentrification from Social Media. Social media is notoriously biased towards tech-savvy younger populations but can we exploit this bias to better understand processes of gentrification and predict urban development in neighbourhoods? We will look at geo-social data from Twitter and Foursquare alongside public data about deprivation in neighbourhoods to understand the contradictions which arise when complex urban processes such as gentrification take place. The role of symbolic capital such as social and cultural capital and the different ways of quantifying it will be explored as a follow-up to their theoretical definition in the first part of the tutorial. This will include hands-on exploration of predicting housing prices in developing neighbourhoods using standard statistical methods in Python.

Online Social Networks and Mobility - Theory

Bruno Gonçalves
www.bgoncalves.com





All this technology is making us antisocial



Social Media

SOCIAL MEOWDIA EXPLAINED



Mobile devices



Geolocated Activity



@bgoncalves

www.bgoncalves.com

@bgoncalves

www.bgoncalves.com

JAN
2016

GLOBAL DIGITAL SNAPSHOT

A SNAPSHOT OF THE WORLD'S KEY DIGITAL STATISTICAL INDICATORS



TOTAL
POPULATION



7.395
BILLION

URBANISATION: 54%

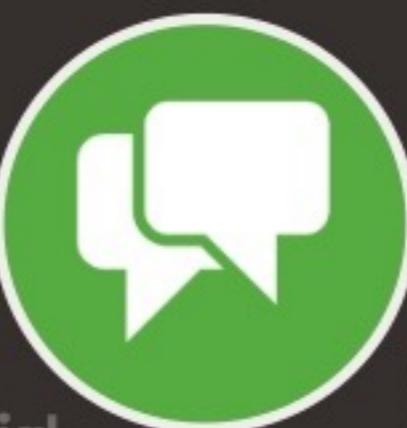
INTERNET
USERS



3.419
BILLION

PENETRATION: 46%

ACTIVE SOCIAL
MEDIA USERS



2.307
BILLION

PENETRATION: 31%

UNIQUE
MOBILE USERS



3.790
BILLION

PENETRATION: 51%

ACTIVE MOBILE
SOCIAL USERS



1.968
BILLION

PENETRATION: 27%

FIGURE REPRESENTS TOTAL GLOBAL POPULATION, INCLUDING CHILDREN

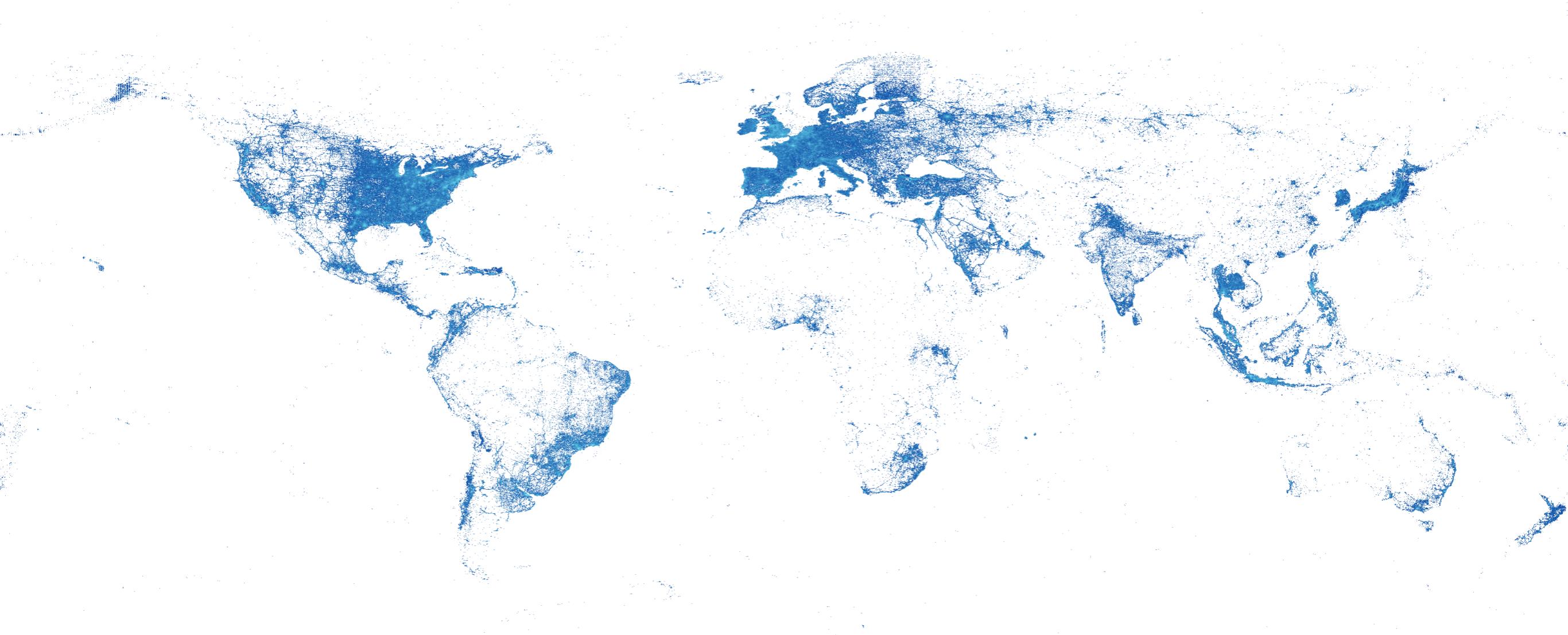
FIGURE INCLUDES ACCESS VIA FIXED AND MOBILE CONNECTIONS

FIGURE BASED ON ACTIVE USER ACCOUNTS, NOT UNIQUE INDIVIDUALS

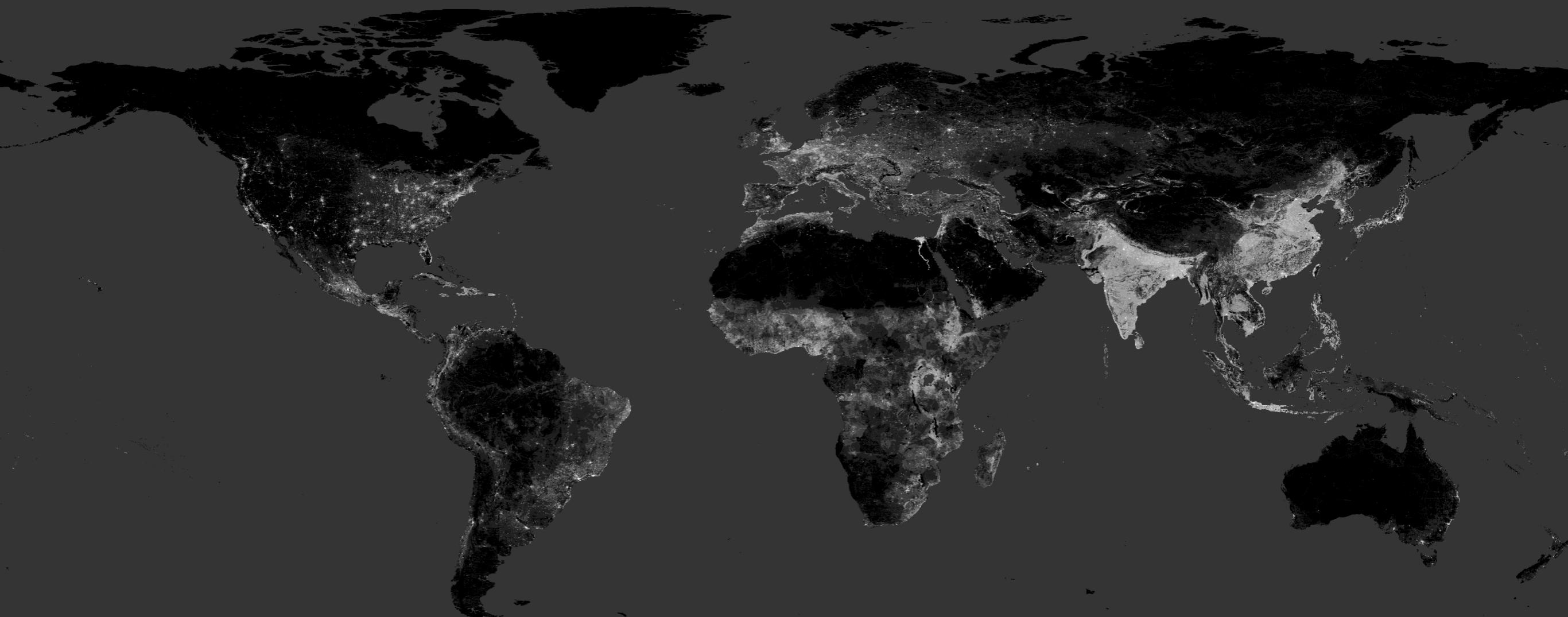
FIGURE REPRESENTS UNIQUE MOBILE PHONE USERS

FIGURE BASED ON ACTIVE USER ACCOUNTS, NOT UNIQUE INDIVIDUALS

World Coverage



World Population

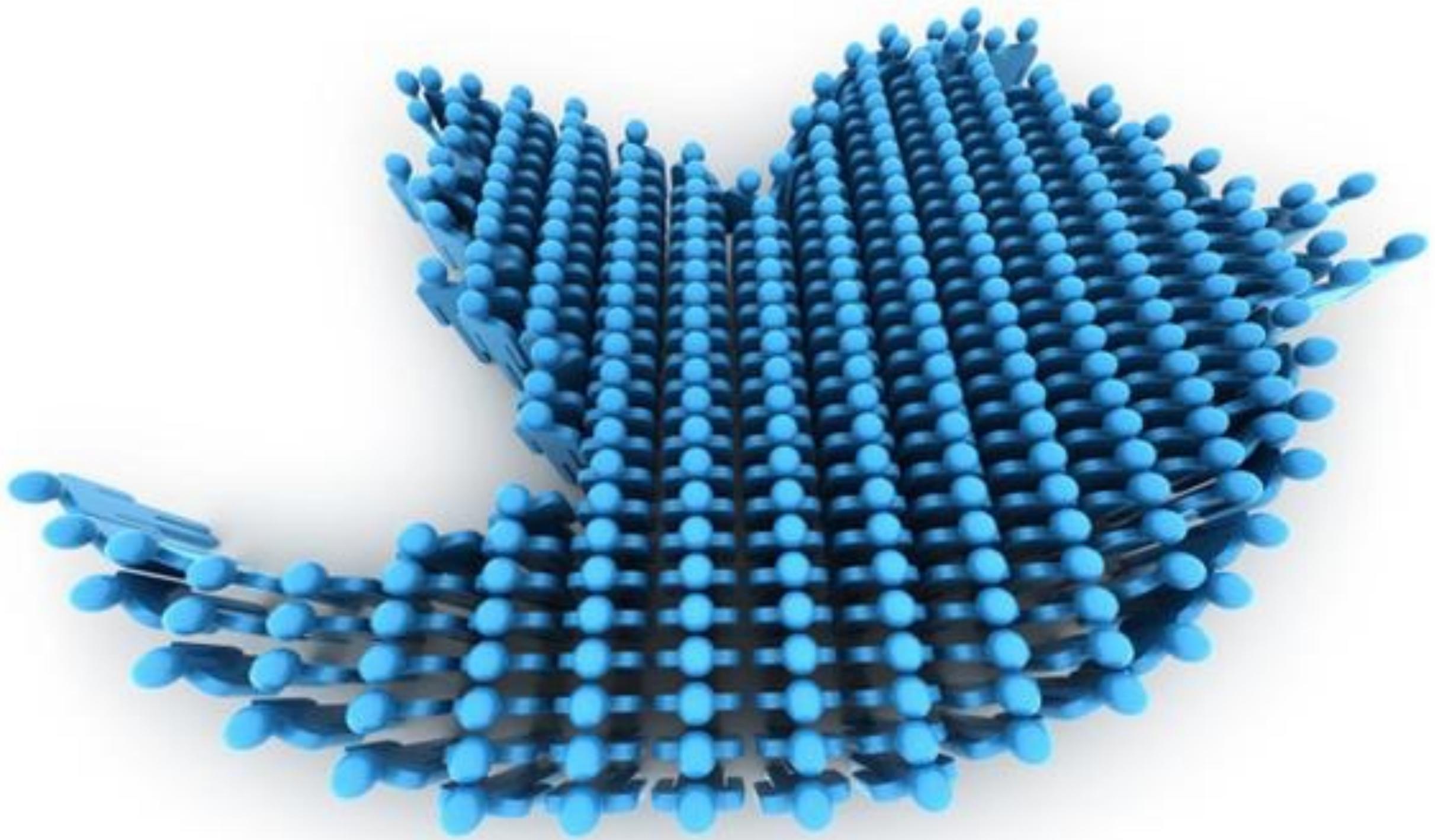


Geolocated Activity

PLoS One 8, E61981 (2013)

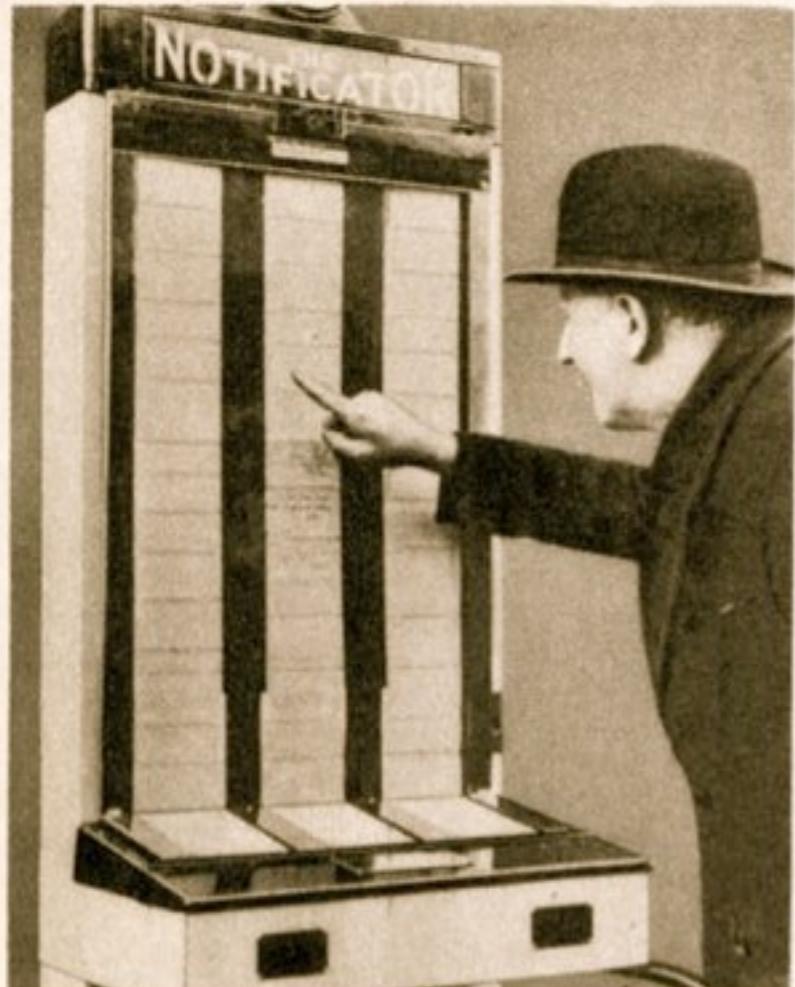


Twitter



Twitter

Robot Messenger Displays Person-to-Person Notes In Public



For a small sum Londoners may leave messages for friends in public places. When written on "notifier," message moves up behind window, remaining in view for two hours.

TO AID persons who wish to make or cancel appointments or inform friends of their whereabouts, a robot message carrier has been introduced in London, England.

Known as the "notifier," the new machine is installed in streets, stores, railroad stations or other public places where individuals may leave messages for friends.

The user walks up on a small platform in front of the machine, writes a brief message on a continuous strip of paper and drops a coin in the slot. The inscription moves up behind a glass panel where it remains in public view for at least two hours so that the person for whom it is intended may have sufficient time to observe the note at the appointed place. The machine is similar in appearance to a candy-vending device.

Source: Modern Mechanix (Aug, 1935)

twitter



Anatomy of a Tweet





Anatomy of a Tweet

```
[u'contributors',
 u'truncated',
 u'text',
 u'in_reply_to_status_id',
 u'id',
 u'favorite_count',
 u'source',
 u'retweeted',
 u'coordinates',
 u'entities',
 u'in_reply_to_screen_name',
 u'in_reply_to_user_id',
 u'retweet_count',
 u'id_str',
 u'favorited',
 u'user',
 u'geo',
 u'in_reply_to_user_id_str',
 u'possibly_sensitive',
 u'lang',
 u'created_at',
 u'in_reply_to_status_id_str',
 u'place',
 u'metadata']
```

Anatomy of a Tweet

```
[u'contributors',
 u'truncated',
 u'text',
 u'in_reply_to_status_id',
 u'id',
 u'favorite_count',
 u'source',
 u'retweeted',
 u'coordinates',
 u'entities',
 u'in_reply_to_screen_name',
 u'in_reply_to_user_id',
 u'retweet_count',
 u'id_str',
 u'favorited',
 u'user',  
    u'geo',
 u'in_reply_to_user_id_str',
 u'possibly_sensitive',
 u'lang',
 u'created_at',
 u'in_reply_to_status_id_str',
 u'place',
 u'metadata']
[u'follow_request_sent',
 u'profile_use_background_image',
 u'default_profile_image',
 u'id',
 u'profile_background_image_url_https',
 u'verified',
 u'profile_text_color',
 u'profile_image_url_https',
 u'profile_sidebar_fill_color',
 u'entities',
 u'followers_count',
 u'profile_sidebar_border_color',
 u'id_str',
 u'profile_background_color',
 u'listed_count',
 u'is_translator_enabled',
 u'utc_offset',
 u'statuses_count',
 u'description',
 u'friends_count',
 u'location',
 u'profile_link_color',
 u'profile_image_url',
 u'following',
 u'geo_enabled',
 u'profile_banner_url',
 u'profile_background_image_url',
 u'screen_name',
 u'lang',  
    u'profile_background_tile',
 u'favourites_count',
 u'name',
 u'notifications',
 u'url',
 u'created_at',
 u'contributors_enabled',
 u'time_zone',
 u'protected',
 u'default_profile',
 u'is_translator']
```

Anatomy of a Tweet

```
[u'contributors',
 u'truncated',
 u'text',
 u'in_reply_to_status_id',
 u'id',
 u'favorite_count',
 u'source',
 u'retweeted',
 u'coordinates',
 u'entities',
 u'in_reply_to_screen_name',
 u'in_reply_to_user_id',
 u'retweet_count',
 u'id_str',
 u'favorited',
 u'user',
 u'geo',
 u'in_reply_to_user_id_str',
 u'possibly_sensitive',
 u'lang',
 u'created_at',
 u'in_reply_to_status_id_str',
 u'place',
 u'metadata']

u"I'm at Terminal Rodovi\xelrio de Feira de Santana
(Feira de Santana, BA) http://t.co/WirvdHwYMq
```

Anatomy of a Tweet

```
[u'contributors',
 u'truncated',
 u'text',
 u'in_reply_to_status_id',
 u'id',
 u'favorite_count',
 u'source',
 u'retweeted',
 u'coordinates',
 u'entities',
 u'in_reply_to_screen_name',
 u'in_reply_to_user_id',
 u'retweet_count',
 u'id_str',
 u'favorited',
 u'user',
 u'geo',
 u'in_reply_to_user_id_str',
 u'possibly_sensitive',
 u'lang',
 u'created_at',
 u'in_reply_to_status_id_str',
 u'place',
 u'metadata']

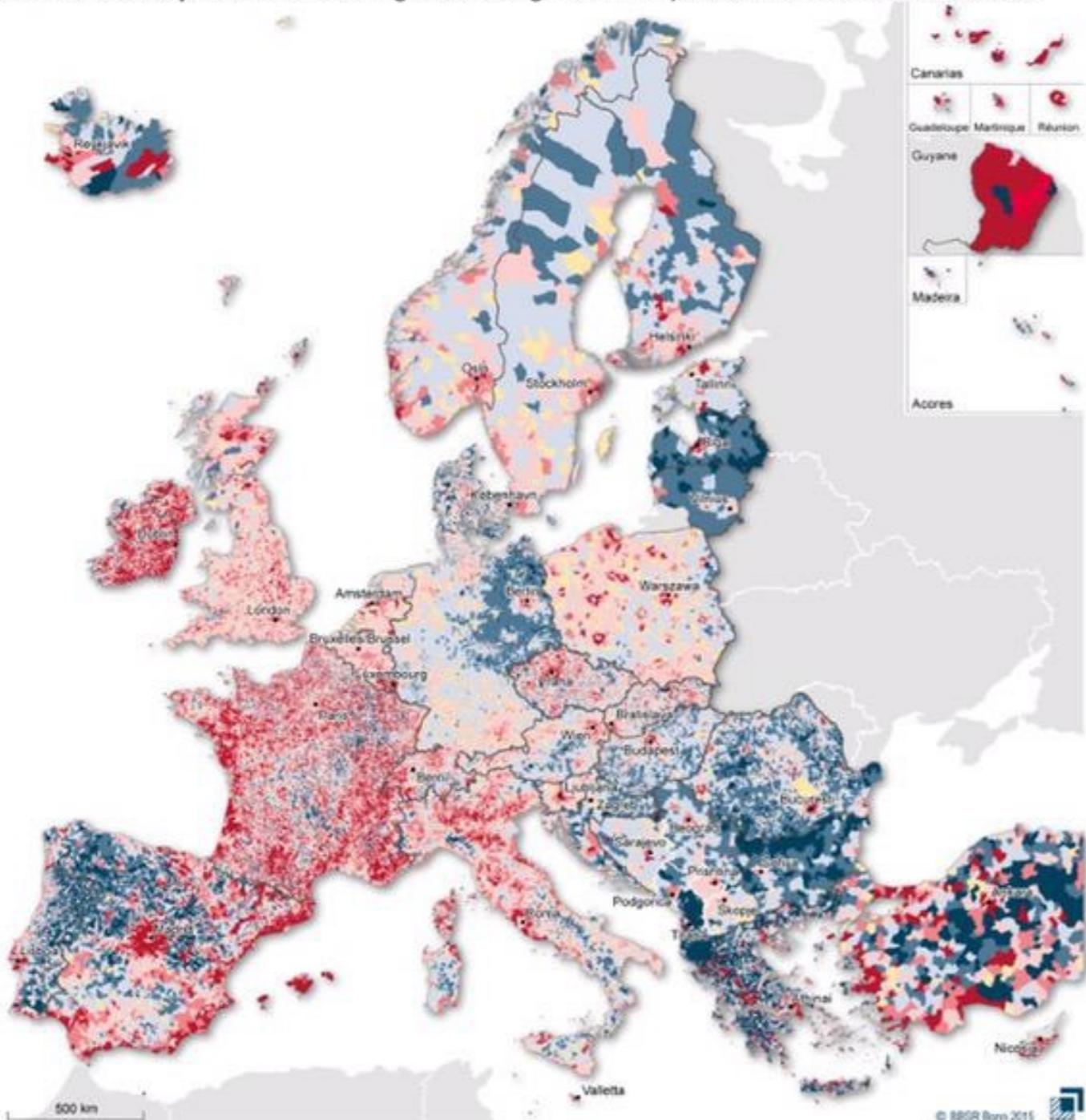
u"I'm at Terminal Rodovi\xcelrio de Feira de Santana
(Feira de Santana, BA) http://t.co/WirvdHwYMq

u'<a href="http://foursquare.com" rel="nofollow">
foursquare</a>'

[u'symbols',
 u'user_mentions',
 u'hashtags',
 u'urls' {u'display_url': u'4sq.com/1k5MeYF',
 u'expanded_url': u'http://4sq.com/1k5MeYF',
 u'indices': [70, 92],
 u'url': u'http://t.co/WirvdHwYMq'}
 u'coordinates']
```

Biases

Durchschnittliche jährliche Bevölkerungsentwicklung in den Europäischen Lokalen Gebietseinheiten



Durchschnittliche jährliche Bevölkerungsentwicklung von 2001-2011* in % in den Gemeinden (LAU2)**



*Bevölkerungsdaten: Dezember 2001, 2011.
**LAU2: 2001: FR, FRA, 2009: EL, NL, PL, SL, RO, 2011:
BA, 2007, 2010: ME, 2003, 2011: BE, 2007, 2011
Regierungsbezirke CH, 2007, 2010, KS, 2011, 2012
**Europa: Subkonsolidiert: LAU2: BG, LT, ME, AM, TR, LAU1:
Asien: Konsolidiert: LAU2: Azerbaidschan, AL, HZ, GL,
LAU1: Asien: BA, KR, KR

Stadtstaat: Laufende Bevölkerungszählung Europa
Bundesamt REU2
Geometrische Grundlage: GRI-Erfassung
Regionen: LAU2
Bevölkerung: R. Böller, L. Bödker, N. Körber-Wilgers,
T. Pannenberg, V. Schmid-Schwarz

Smartphone Ownership Highest Among Young Adults, Those With High Income/Education Levels

% of U.S. adults in each group who own a smartphone

All adults	64%
Male	66
Female	63
18-29	85
30-49	79
50-64	54
65+	27
White, non-Hispanic	61
Black, non-Hispanic	70
Hispanic	71
HS grad or less	52
Some college	69
College+	78
Less than \$30,000/yr	50
\$30,000-\$49,999	71
\$50,000-\$74,999	72
\$75,000 or more	84
Urban	68
Suburban	66
Rural	52

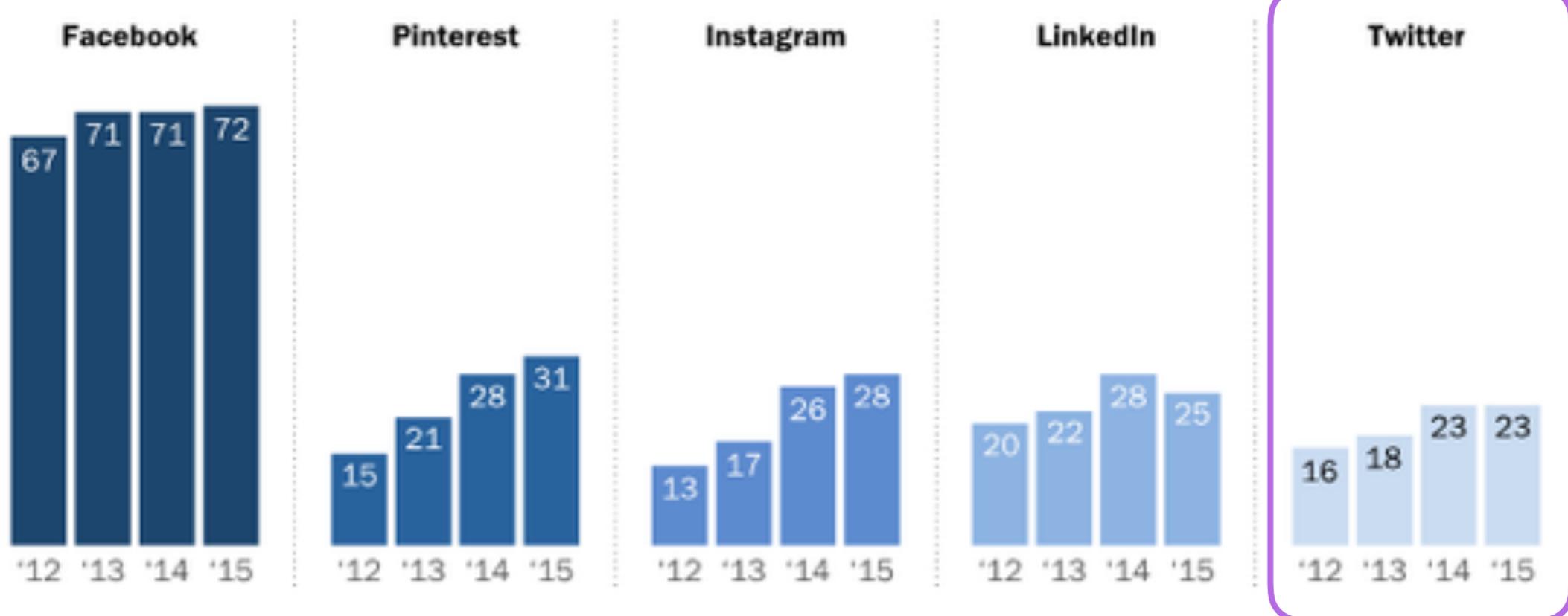
Combined analysis of Pew Research Center surveys conducted Dec. 4-7 and Dec. 18-21, 2014.

PEW RESEARCH CENTER

Biases

Pinterest and Instagram Usage Doubles Since 2012, Growth on Other Platforms is Slower

% of online adults who say they use the following social media platform, by year



Pew Research Center Survey, March 17-April 12, 2015.

PEW RESEARCH CENTER

Age Distribution

PLoS One 10, e0115545 (2015)

Twitter users

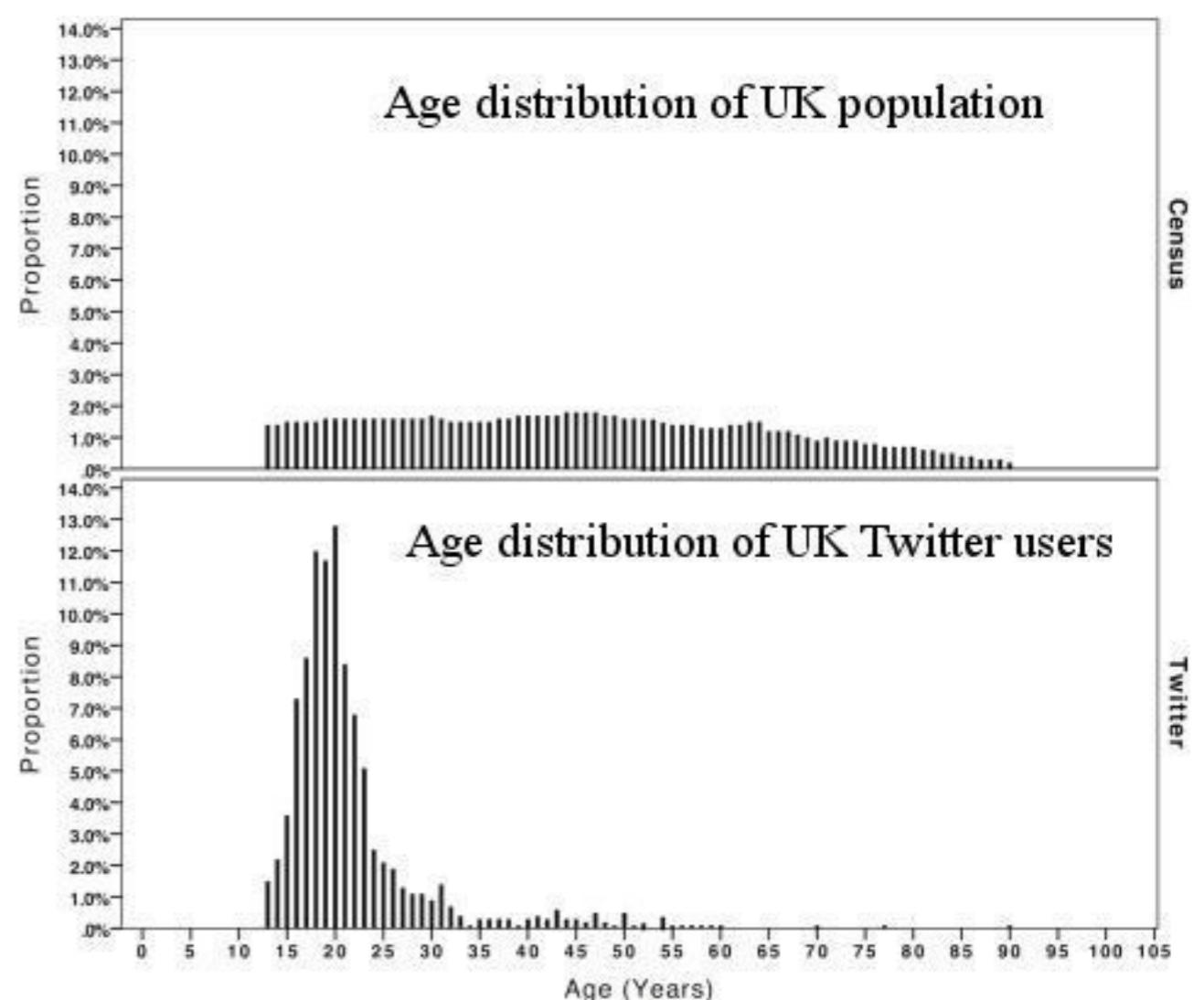
Among online adults, the % who use Twitter

	2013	2014
All internet users	18%	23%*
Men	17	24*
Women	18	21
White, Non-Hispanic	16	21 *
Black, Non-Hispanic	29	27
Hispanic	16	25
18-29	31	37
30-49	19	25
50-64	9	12
65+	5	10*
High school grad or less	17	16
Some college	18	24
College+ (n= 685)	18	30*
Less than \$30,000/yr	17	20
\$30,000-\$49,999	18	21
\$50,000-\$74,999	15	27*
\$75,000+	19	27*
Urban	18	25*
Suburban	19	23
Rural	11	17

Source: Pew Research Center's Internet Project September Combined Omnibus Survey, September 11-14 & September 18-21, 2014. N=1,597 internet users ages 18+. The margin of error for all internet users is +/- 2.9 percentage points. 2013 data from Pew Internet August Tracking Survey, August 07 - September 16, 2013, n= 1,445 internet users ages 18+.

Note: Percentages marked with an asterisk (*) represent a significant change from 2013.
Results are significant at the 95% confidence level using an independent z-test.

PEW RESEARCH CENTER



Age Distribution

Twitter users

Among online adults, the % who use Twitter

	2013	2014
All internet users	18%	23%*
Men	17	24*
Women	18	21
White, Non-Hispanic	16	21 *
Black, Non-Hispanic	29	27
Hispanic	16	25
18-29	31	37
30-49	19	25
50-64	9	12
65+	5	10*
High school grad or less	17	16
Some college	18	24
College+ (n= 685)	18	30*
Less than \$30,000/yr	17	20
\$30,000-\$49,999	18	21
\$50,000-\$74,999	15	27*
\$75,000+	19	27*
Urban	18	25*
Suburban	19	23
Rural	11	17

Source: Pew Research Center's Internet Project September Combined Omnibus Survey, September 11-14 & September 18-21, 2014. N=1,597 internet users ages 18+. The margin of error for all internet users is +/- 2.9 percentage points. 2013 data from Pew Internet August Tracking Survey, August 07 - September 16, 2013, n= 1,445 internet users ages 18+.

Note: Percentages marked with an asterisk (*) represent a significant change from 2013. Results are significant at the 95% confidence level using an independent z-test.

PEW RESEARCH CENTER

Facebook Demographics

Among internet users, the % who use Facebook

	Internet users
Total	72%
Men	66
Women	77
White, Non-Hispanic	70
Black, Non-Hispanic (n=85)	67
Hispanic	75
18-29	82
30-49	79
50-64	64
65+	48
High school grad or less	71
Some college	72
College+	72
Less than \$30,000/yr	73
\$30,000-\$49,999	72
\$50,000-\$74,999	66
\$75,000+	78
Urban	74
Suburban	72
Rural	67

Source: Pew Research Center, March 17-April 12, 2015.

PEW RESEARCH CENTER

Instagram Demographics

Among internet users, the % who use Instagram

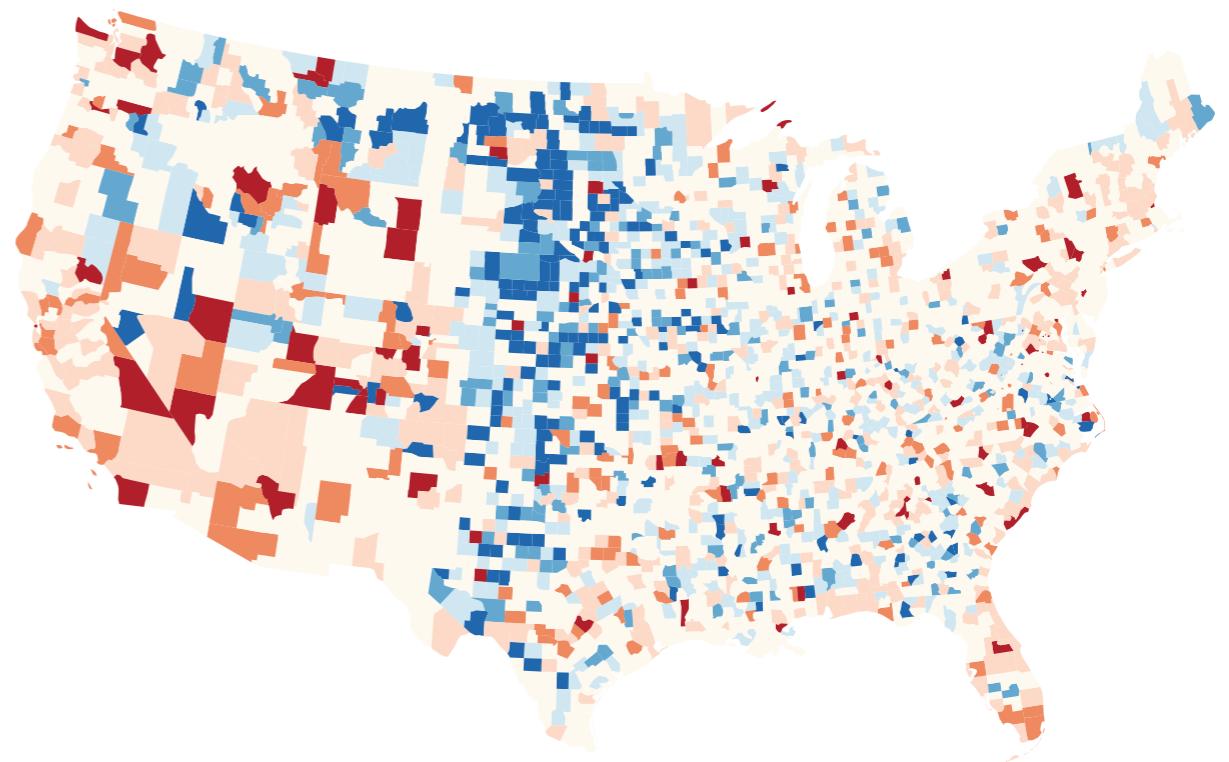
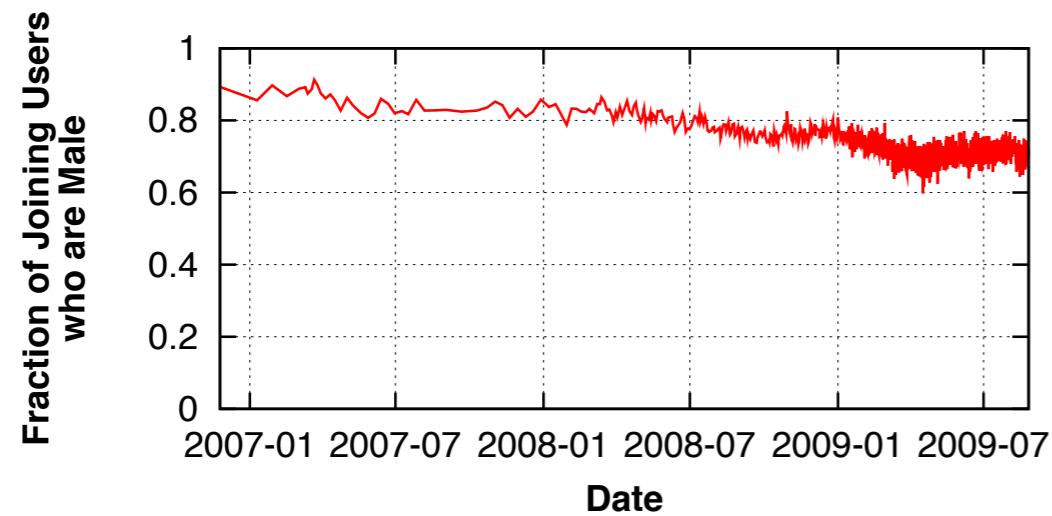
	Internet user
Total	28%
Men	24
Women	31
White, Non-Hispanic	21
Black, Non-Hispanic (n=85)	47
Hispanic	38
18-29	55
30-49	28
50-64	11
65+	4
High school grad or less	25
Some college	32
College+	26
Less than \$30,000/yr	26
\$30,000-\$49,999	27
\$50,000-\$74,999	30
\$75,000+	26
Urban	32
Suburban	28
Rural	18

Source: Pew Research Center, March 17-April 12, 2015.

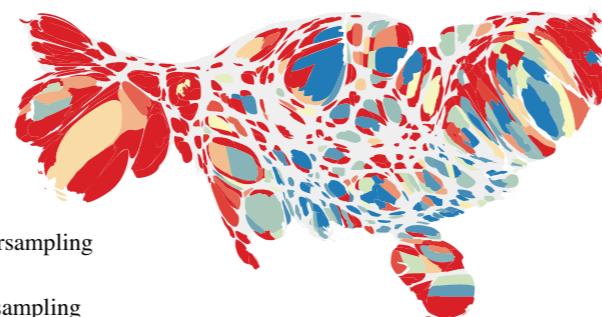
PEW RESEARCH CENTER

Demographics

ICWSM'11, 375 (2011)



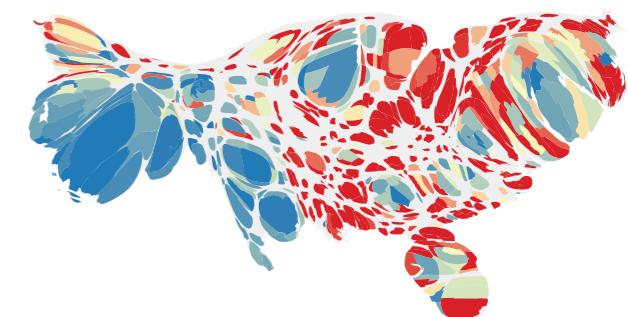
(a) Caucasian (non-hispanic)



(b) African-American

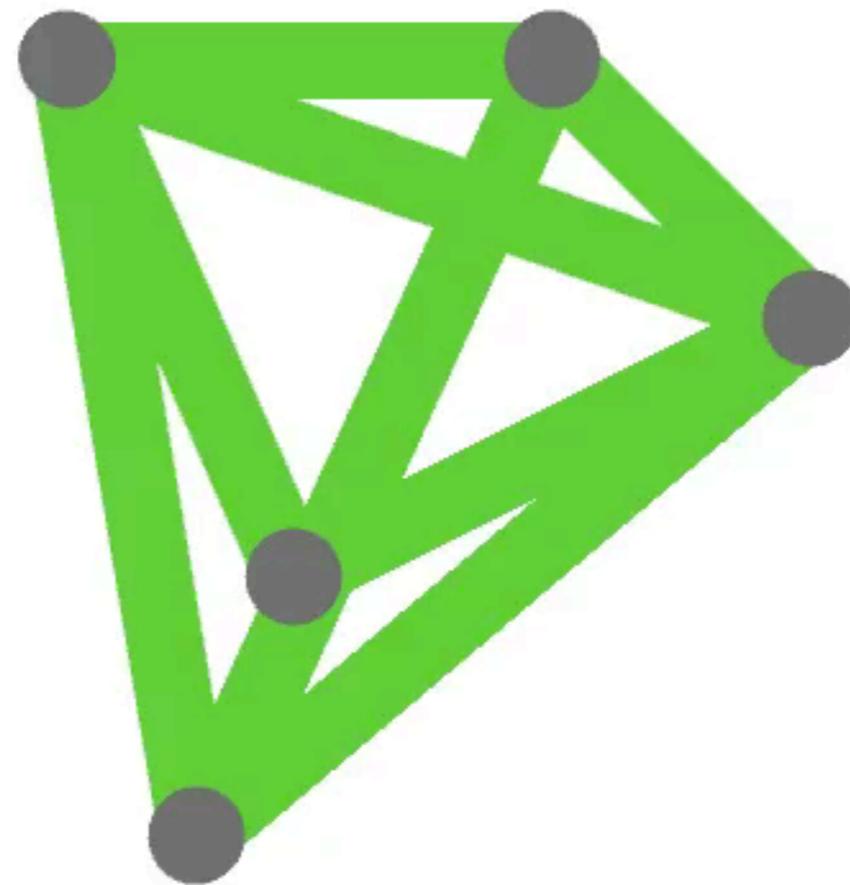


(c) Asian or Pacific Islander



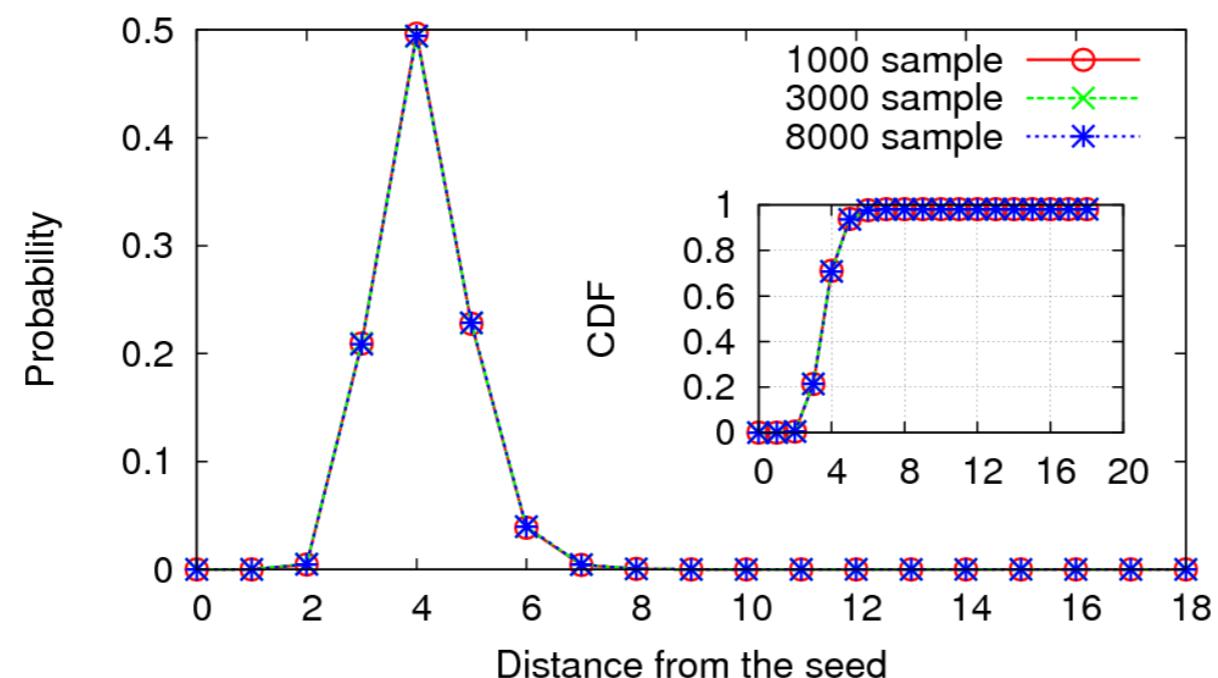
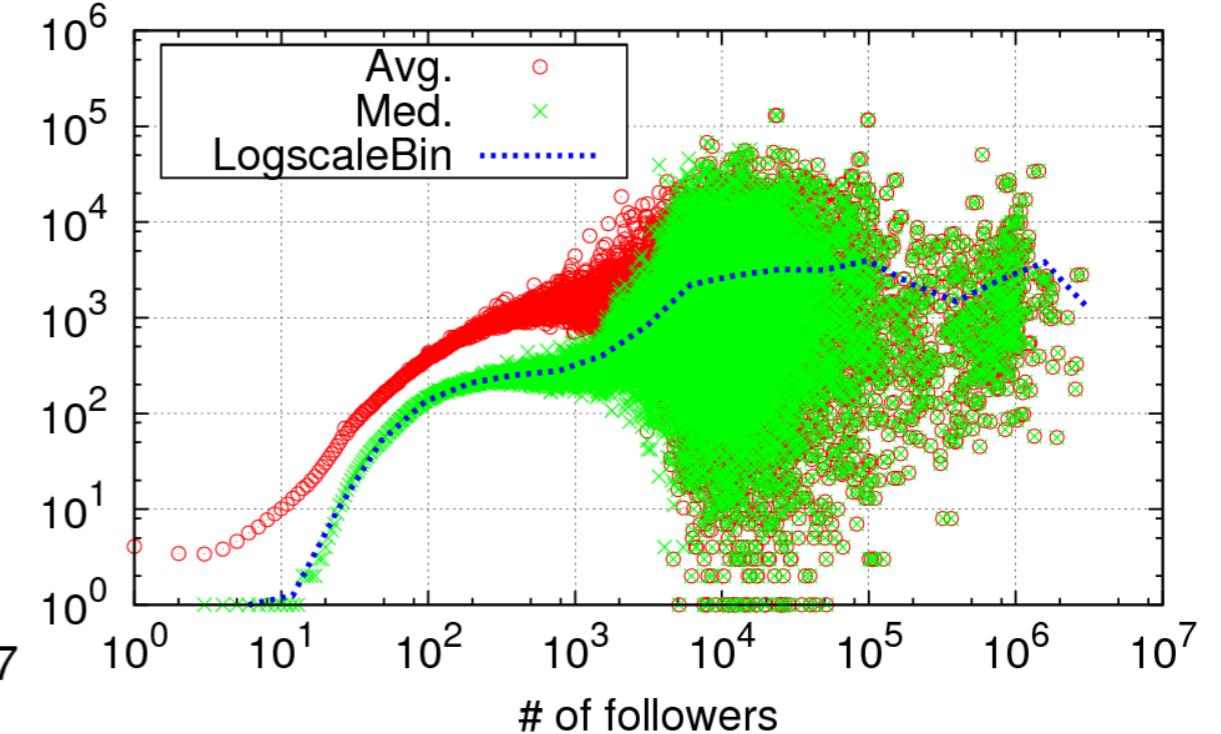
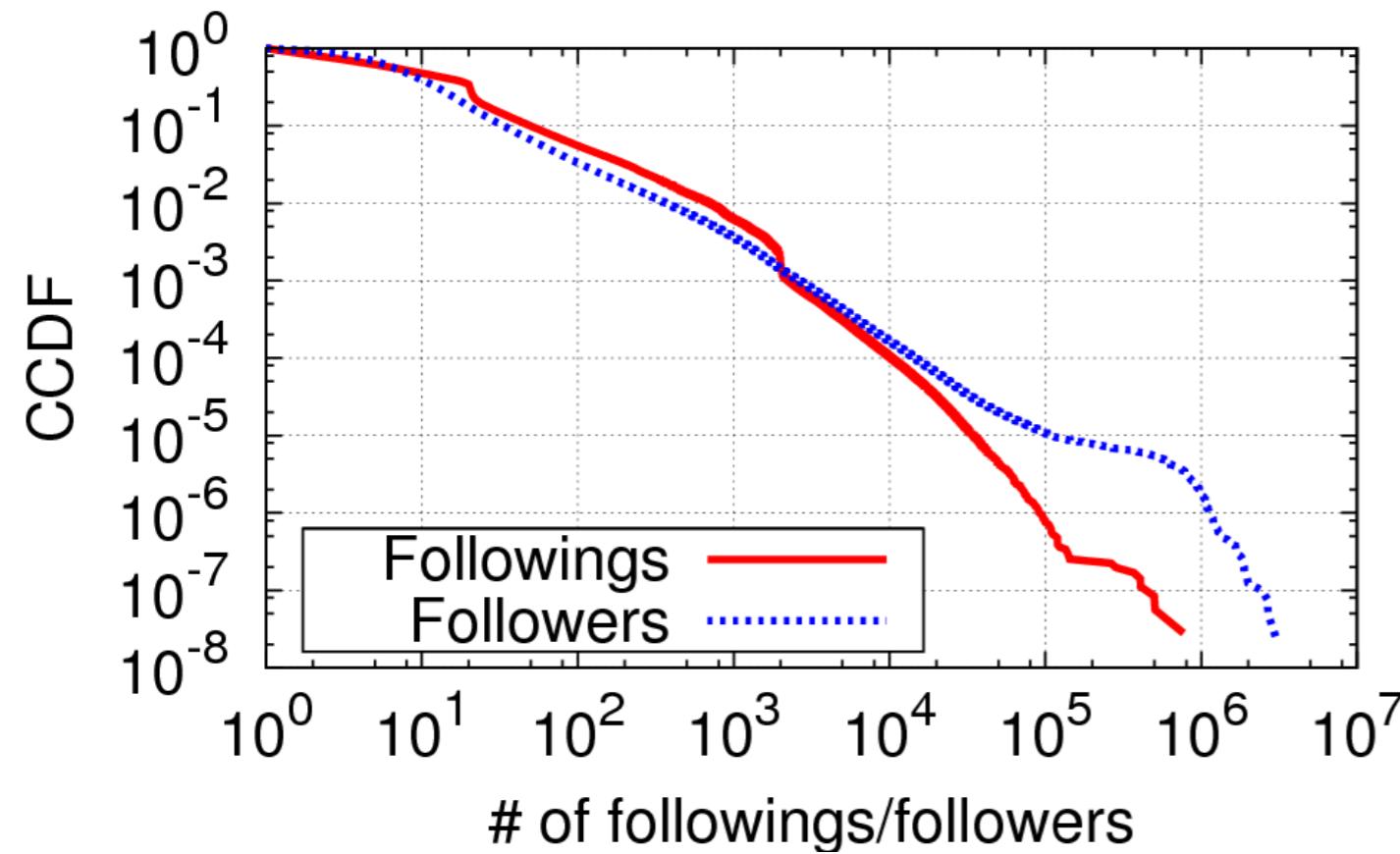
(d) Hispanic

Multilayer Network



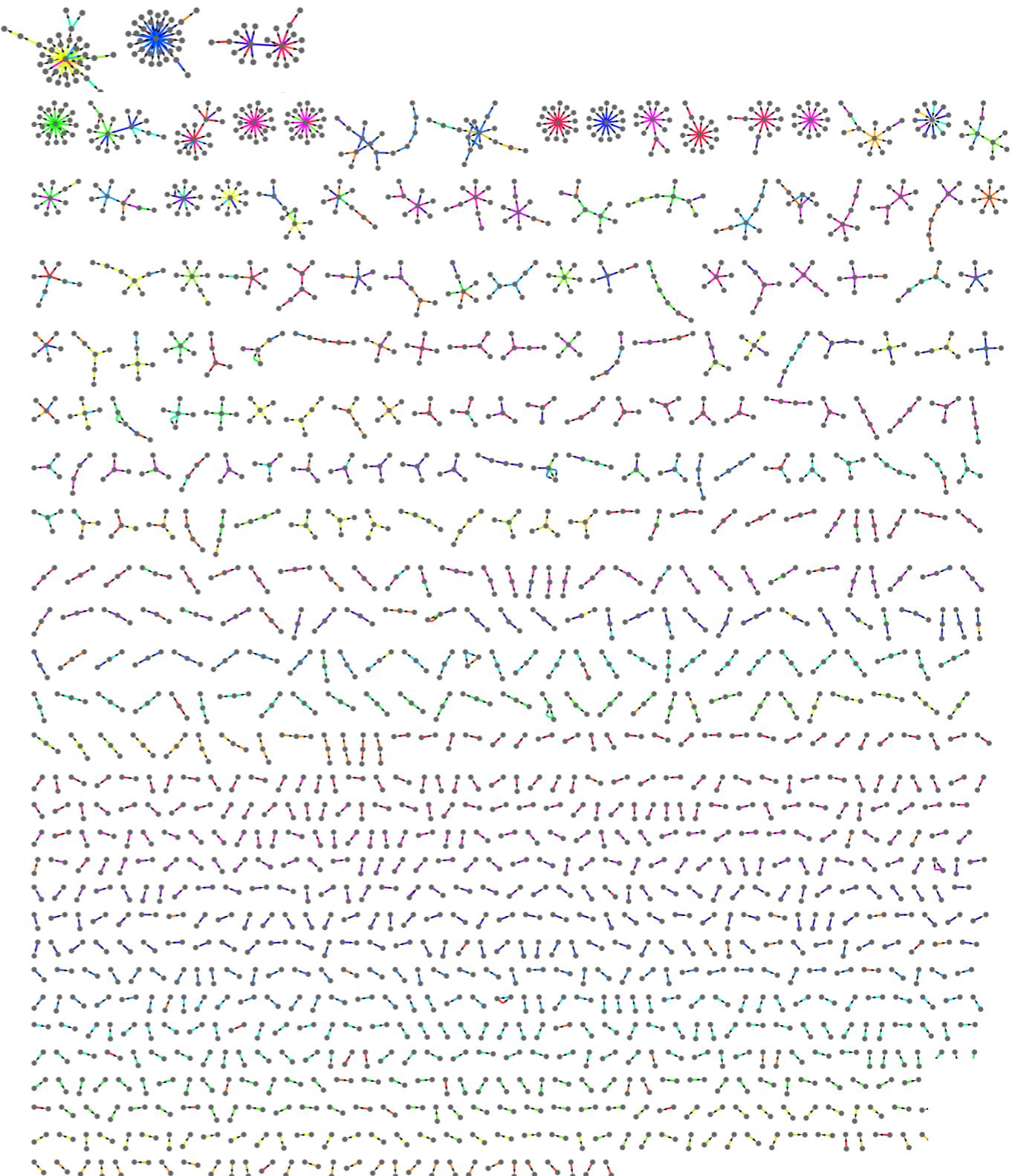
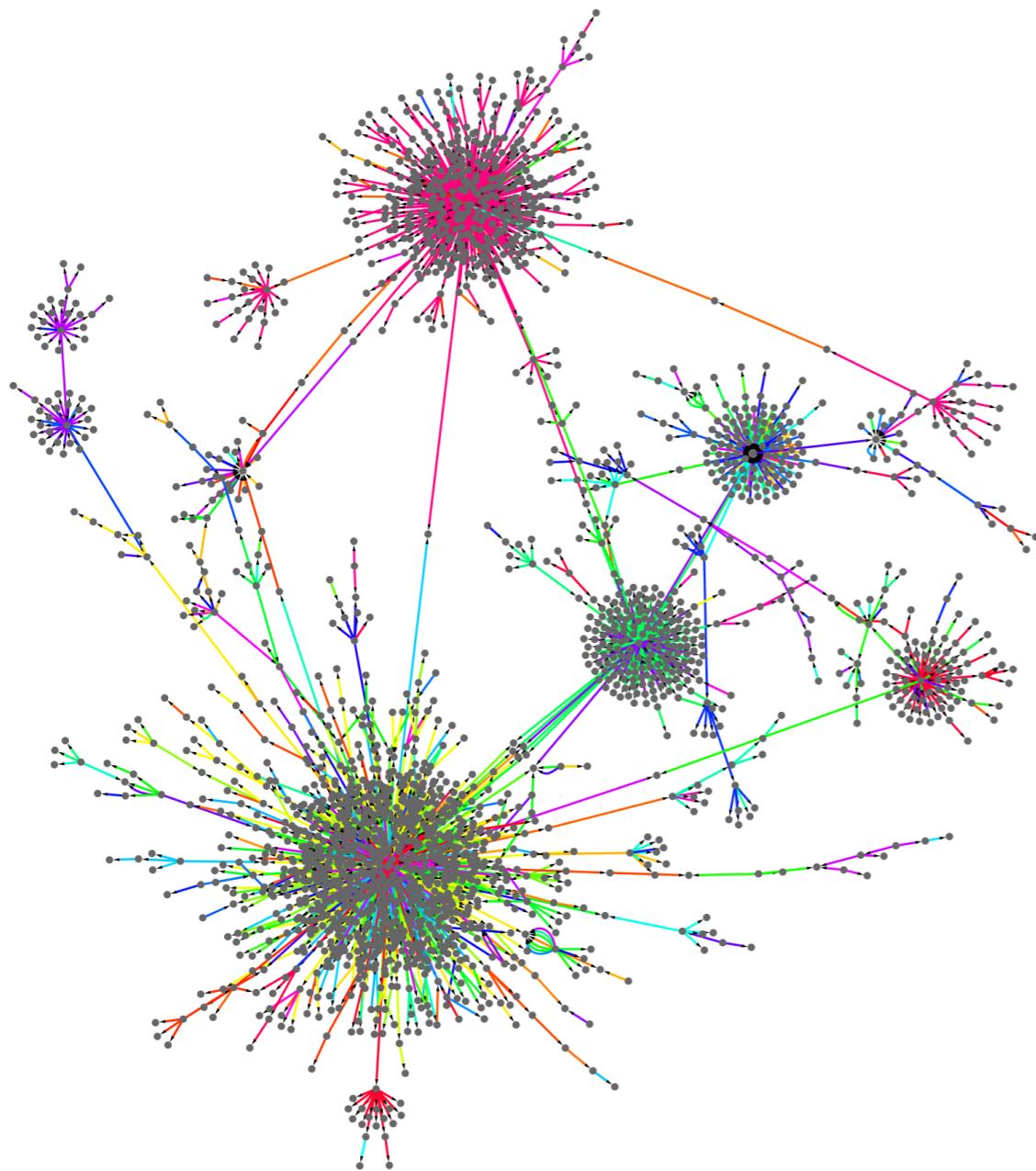
Twitter Network

WWW'10, 591 (2010)



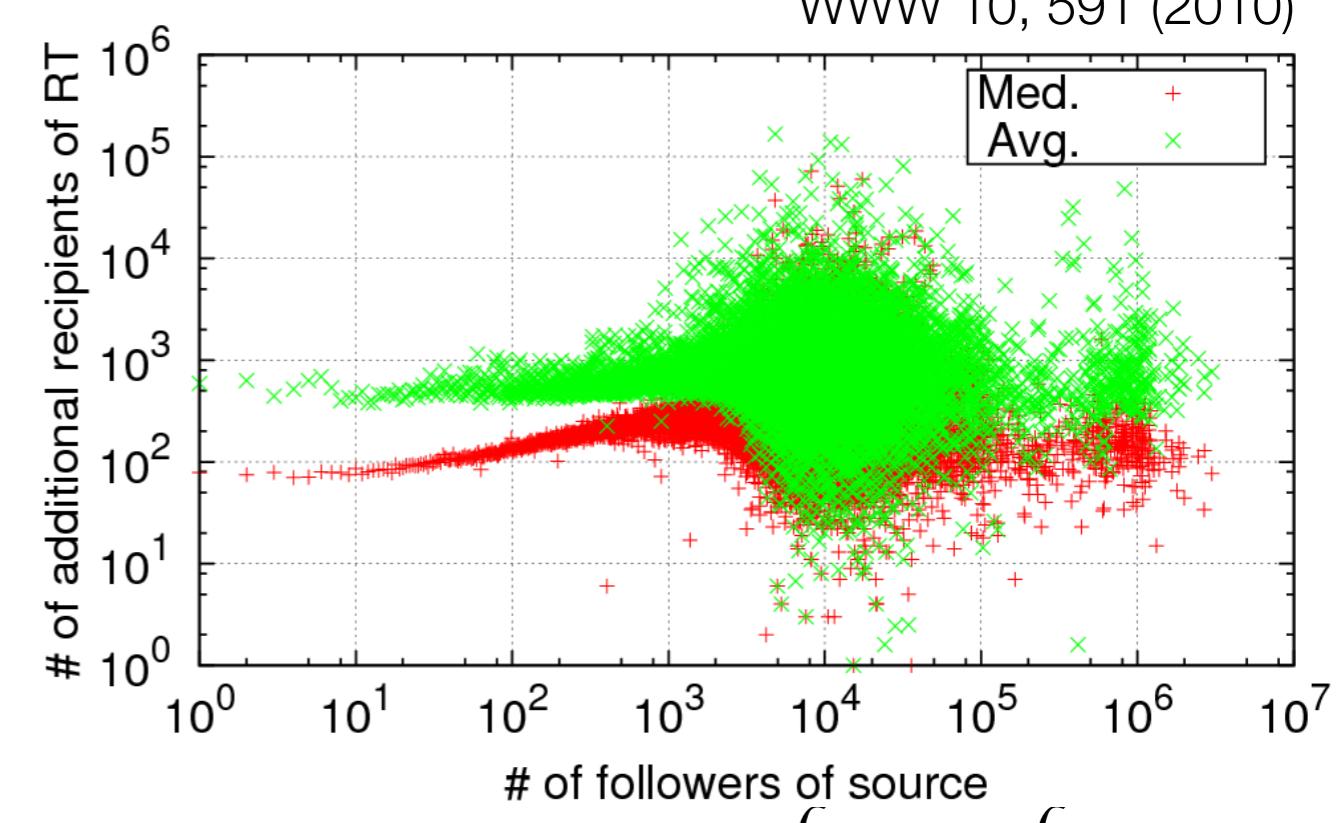
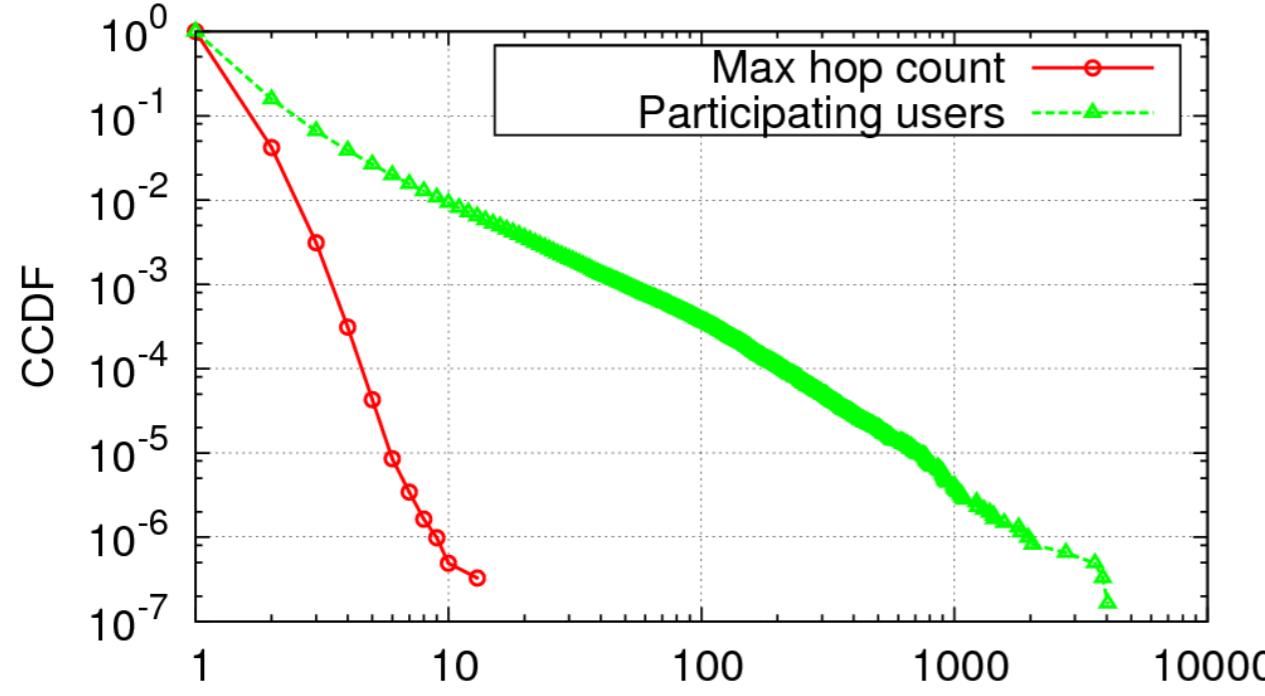
Retweet Trees

WWW'10, 591 (2010)



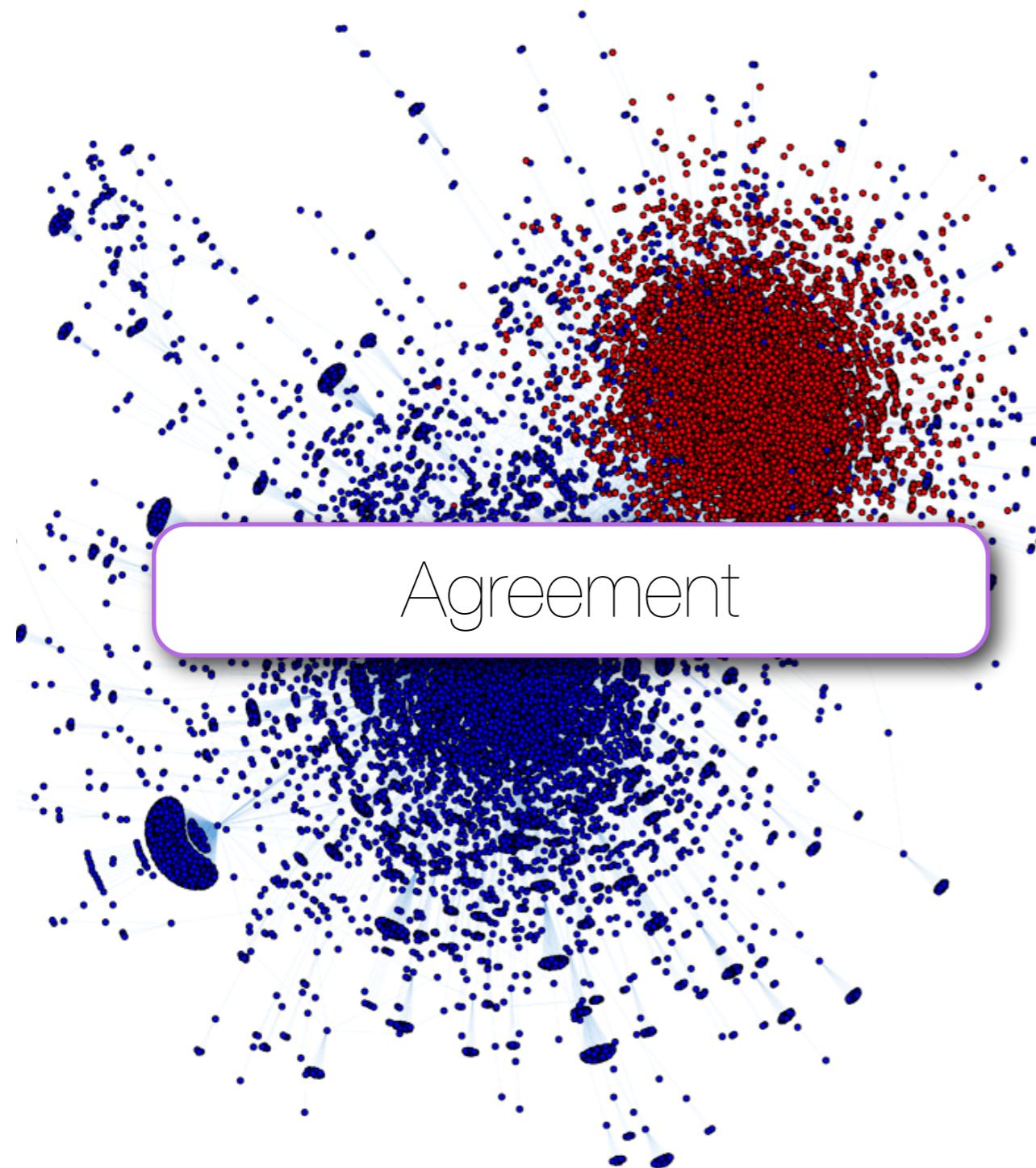
Retweets Trees

WWW'10, 591 (2010)

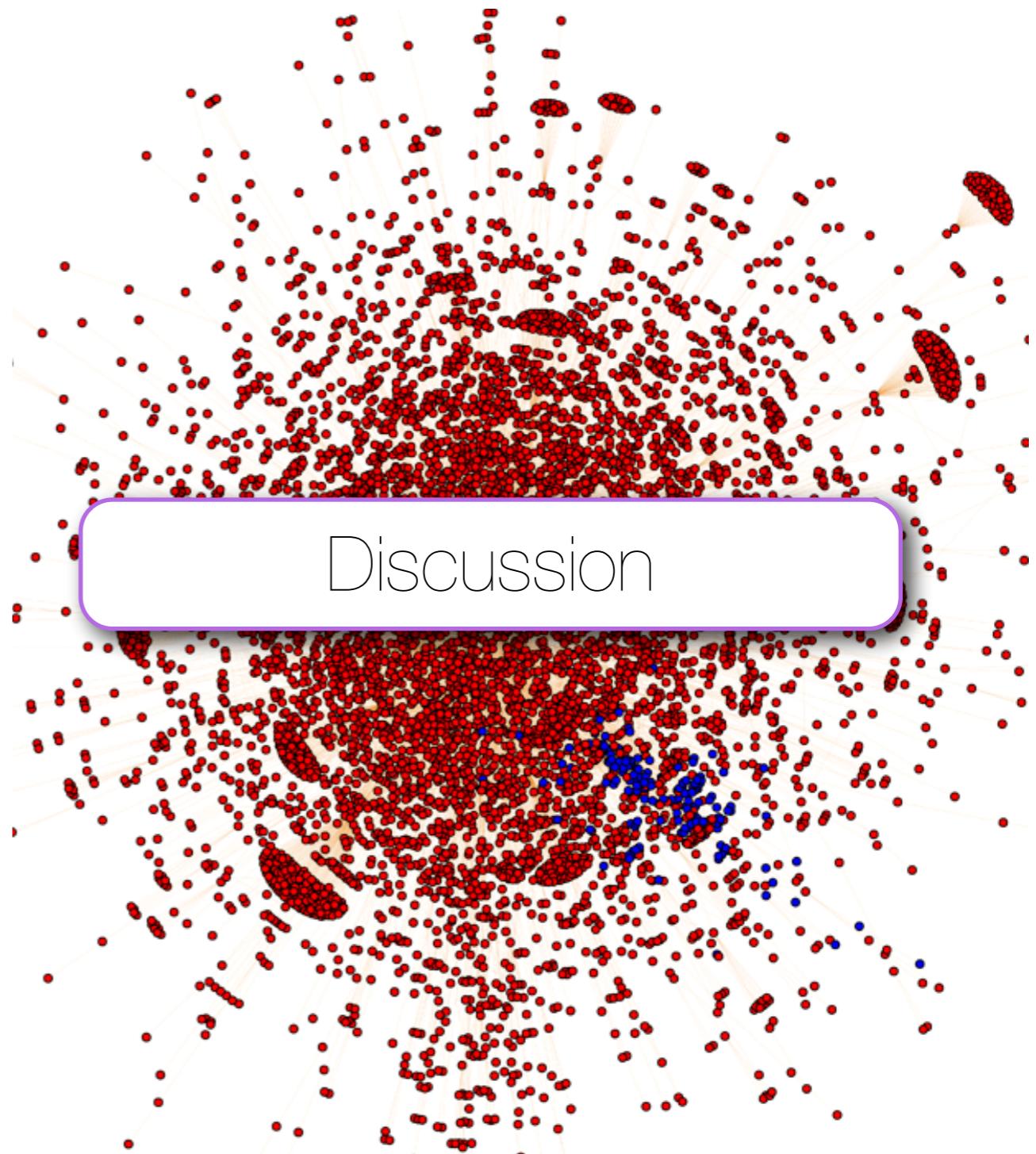


Link Function

ICWSM'11, 89 (2011)



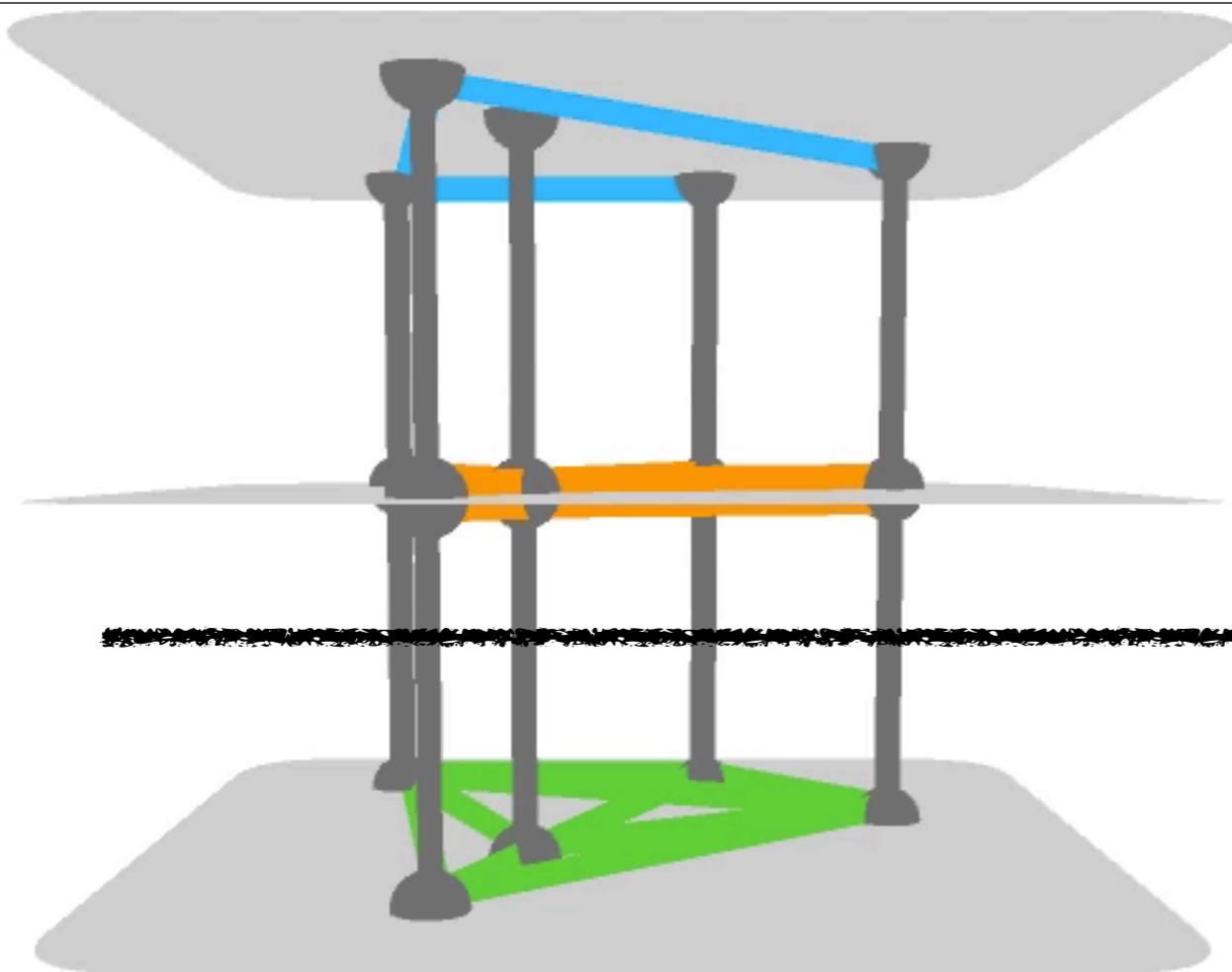
Agreement



Discussion

Multilayer Network

Retweet



Information
Layer(s)

Mention

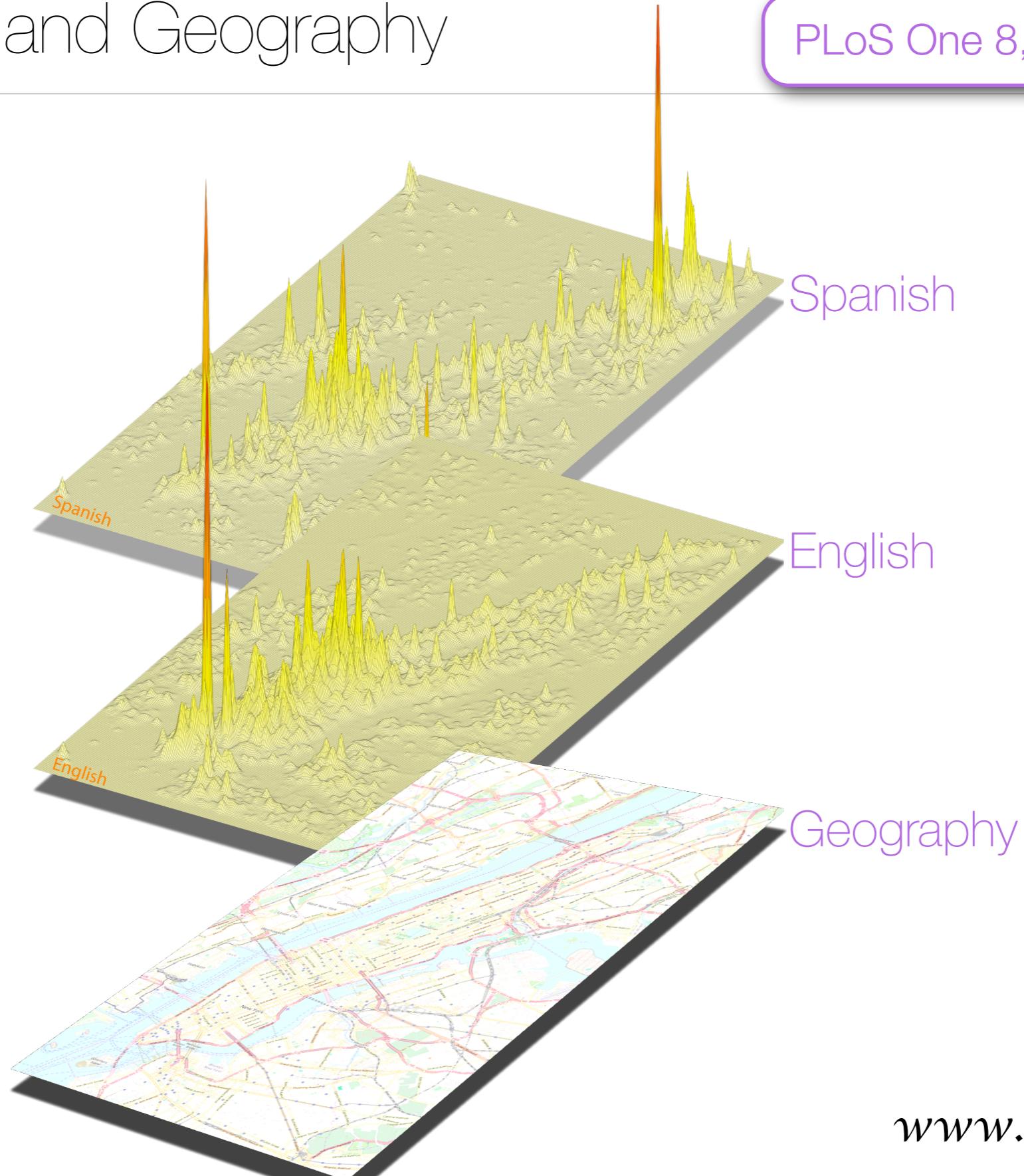
Social
Layer(s)

Follower

Geographical
Layer(s)

Language and Geography

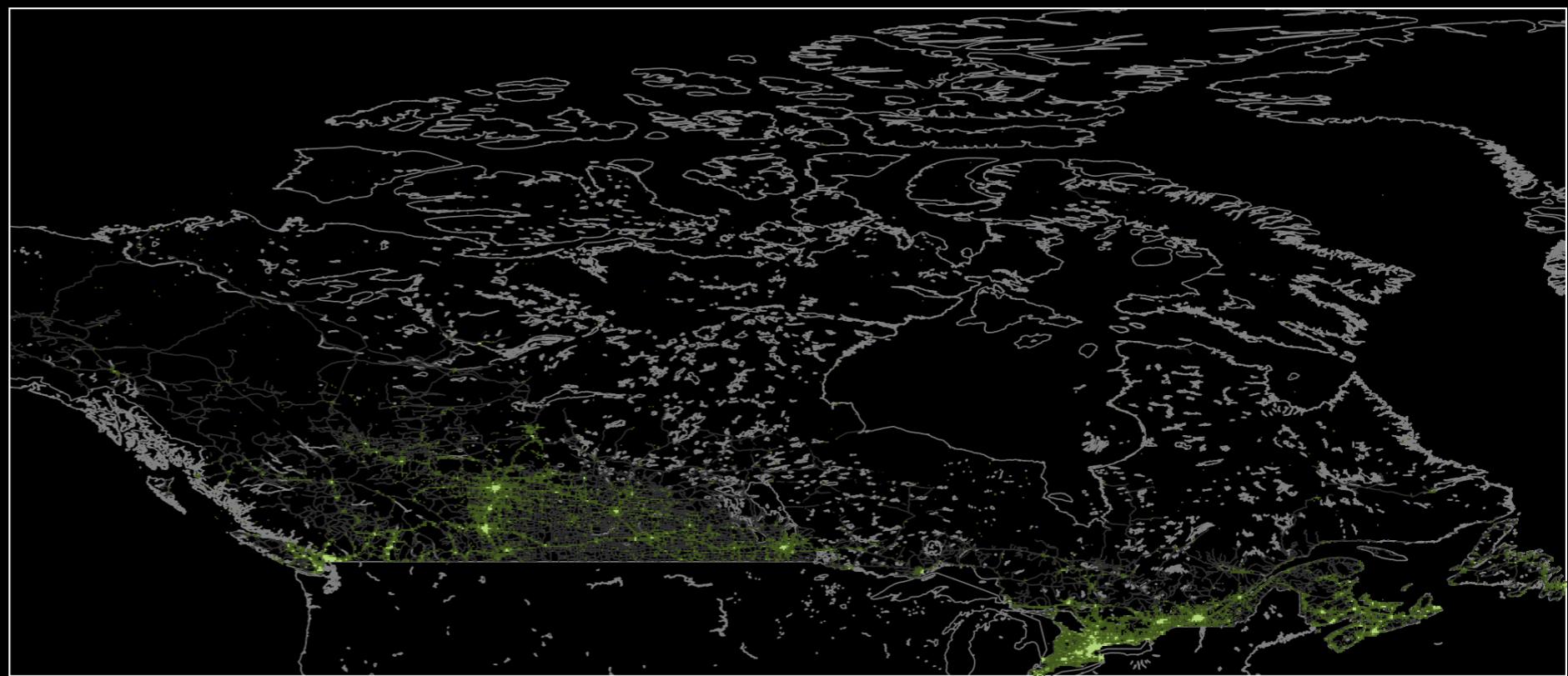
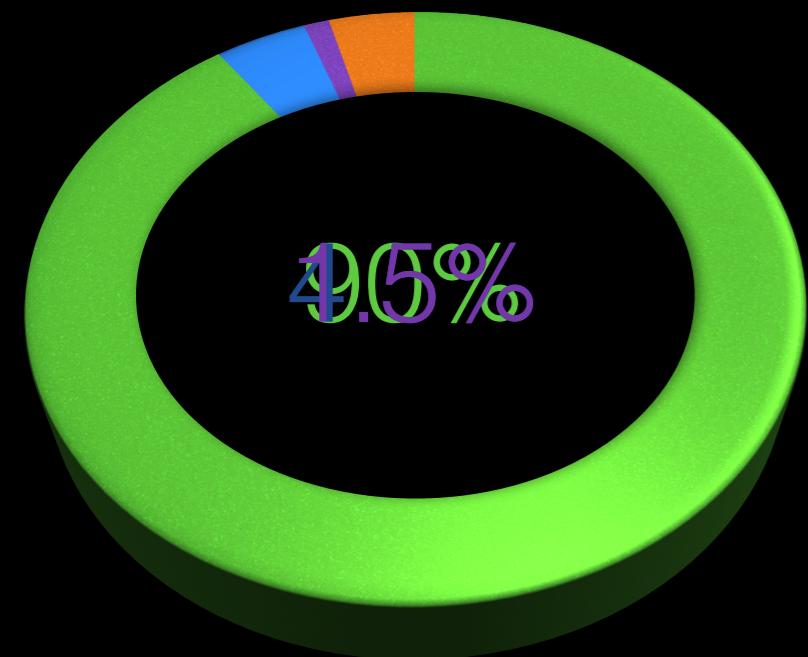
PLoS One 8, E61981 (2013)



Signal By Language

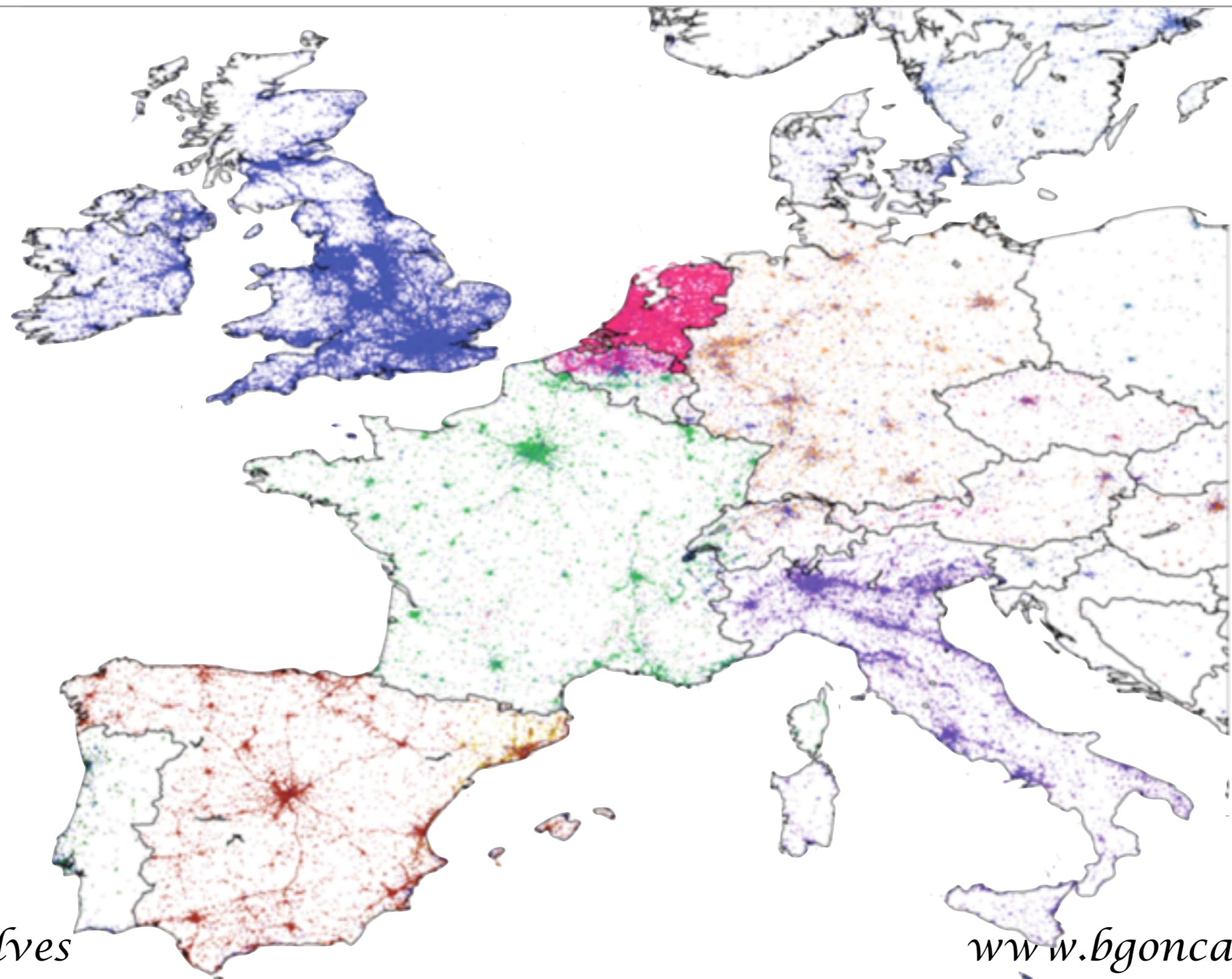
PLoS One 8, E61981 (2013)

- English
- French
- Spanish
- Other



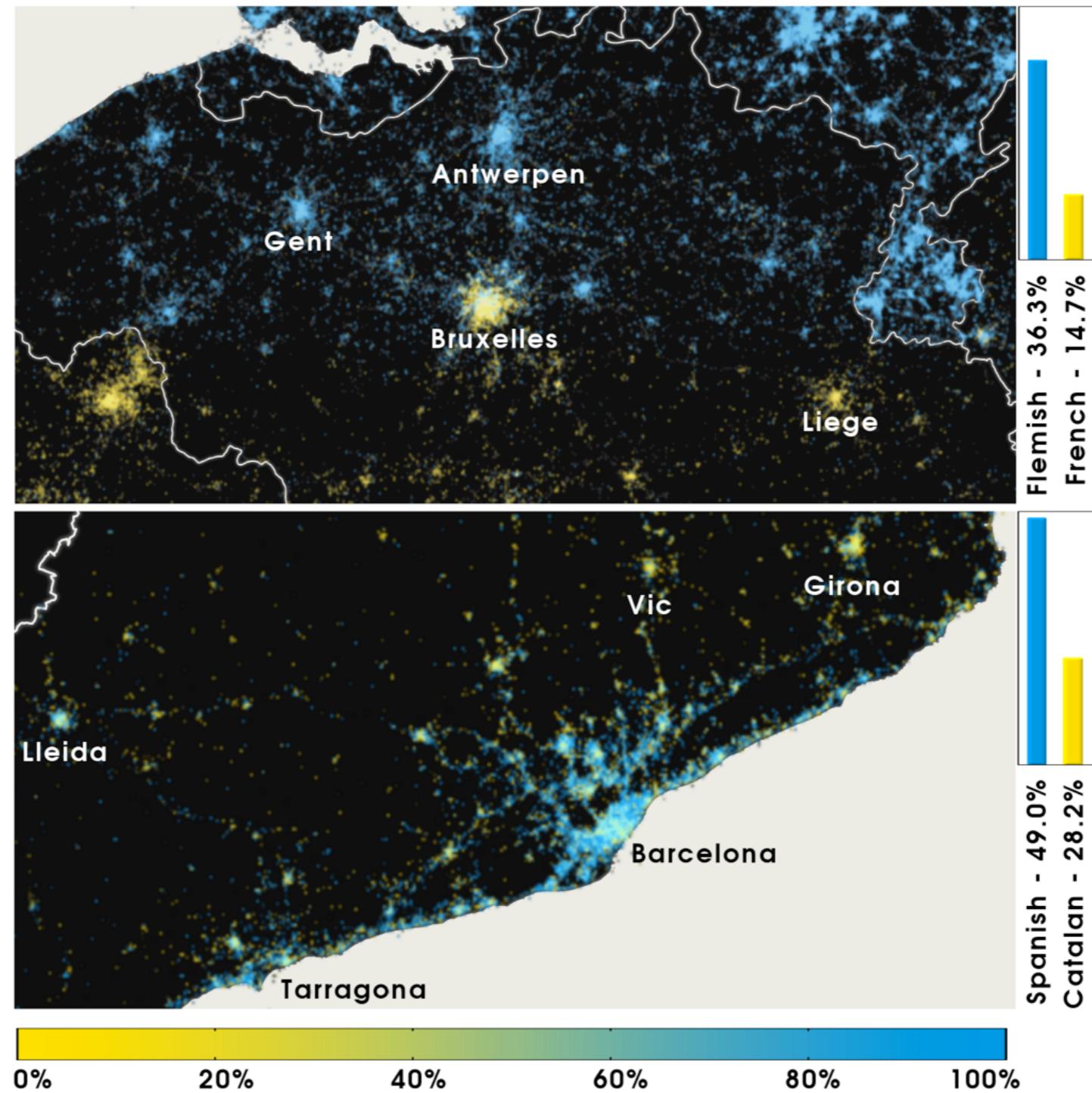
Language Geography

PLoS One 8, E61981 (2013)

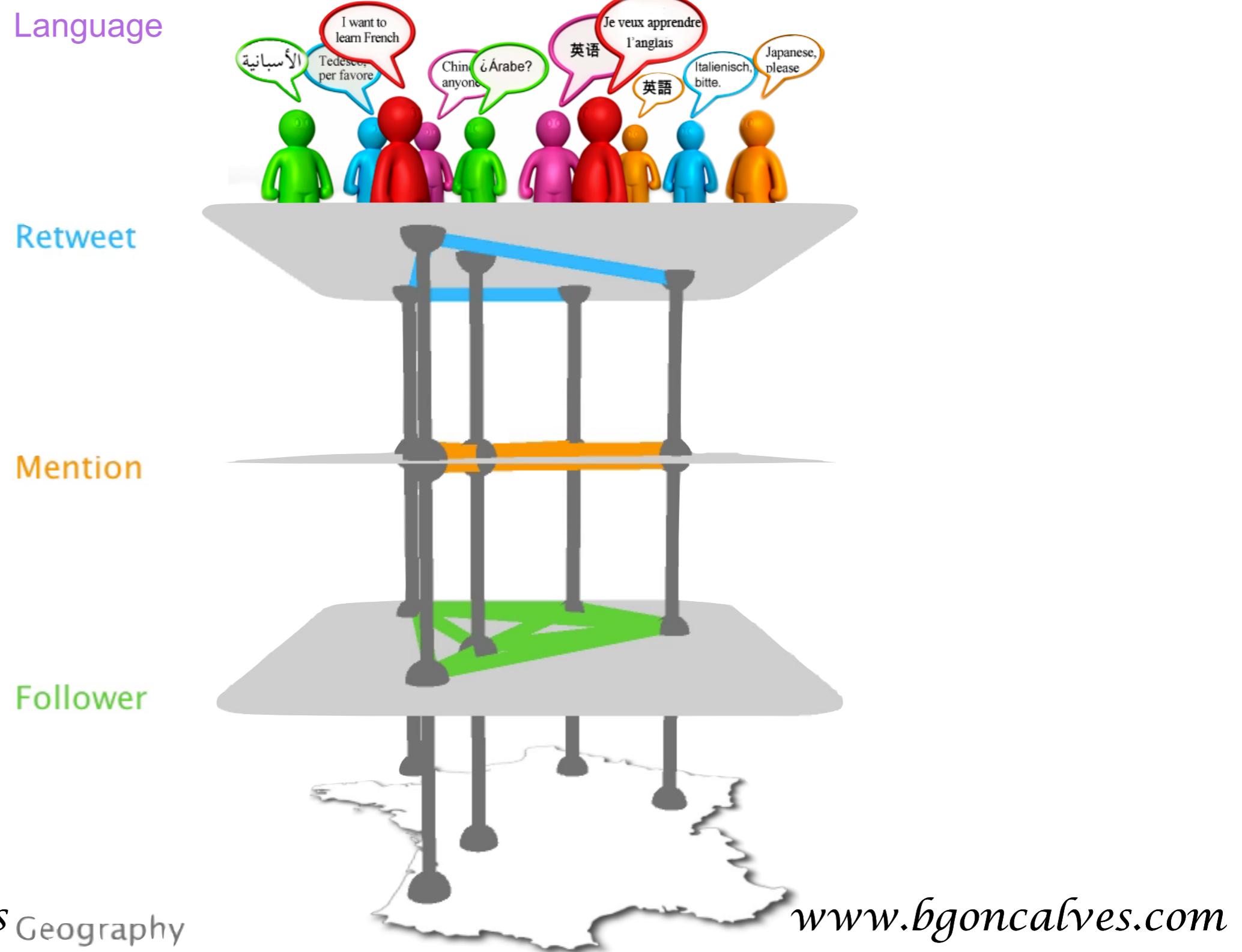


Language Distribution

PLoS One 8, E61981 (2013)



Multilayer Network



Foursquare



Foursquare



Anatomy of a Checkin

NASA checked in at **NASA HQ**
Washington, D.C. | December 27, 2012 via foursquare Web

LIKE - 1837 LIKES

We've got a new 'Curiosity Explorer' badge. Explore your curiosity at science museum s & planetariums to earn it <http://go.nasa.gov/9kpN5g>

Only NASA's friends can see comments and add their own.

1837 people like this

NASA HQ
Washington, D.C.
Government Building

SAVE

Switch Language: English

Anatomy of a Checkin

```
[u'venue',
 u'like',
 u'photos',
 u'source',
 u'vesibility',
 u'entities',
 u'shoot',
 u'timeZoneOffset',
 u'type',
 u'id',
 u'createdAt',
 u'likes']
```

Anatomy of a Checkin

```
[u'verified',
 u'name',
 u'url',
 u'like',
 u'contact',
 u'location',
 u'stats',
 u'id',
 u'categories',
 u'likes']

[ {u'indices': [112, 137],
  u'object': {u'url': u'http://go.nasa.gov/9kpN5g'},
  u'type': u'url'}]

u"We've got a new 'Curiosity Explorer' badge. Explore your
curiosity at science museums & planetariums to earn it http://go.nasa.gov/9kpN5g"
```



```
{u'count': 1837,
 u'groups': [{u'count': 1837, u'items': [], u'type': u'others'}],
 u'summary': u'1837 likes'}
```

Objects

- Two main types of objects:
 - Users
 - Venues
- Multiple possible actions (by Users on Venues)
 - Checkin
 - Like
 - Tip

API Limitations

- Users have privacy concerns with respect to publicly sharing their location.
 - "Stalker apps"
 - "Please Rob Me"
- Privacy is a big concern for Foursquare
- API structure reflects this
- "Easy" to get information on users and on venues. Connecting users to venues much harder to do.

Objects

<https://developer.foursquare.com/overview/venues.html>

Venues

- Venues correspond to physical locations
- Are perhaps the most important object in the Foursquare universe
- API is particularly generous, allowing for **5000** requests per hour.

Objects

<https://developer.foursquare.com/docs/users/users>

Venues

- Venues correspond to physical locations
- Are perhaps the most important object in the Foursquare universe
- API is particularly generous, allowing for **5000** requests per hour.

Users

- Users interact with venues in multiple ways:
 - checkin
 - leaving a “tip”
 - “liking”
- Users connect to each other through friendship/colocation

Objects

<https://developer.foursquare.com/docs/checkins/checkins.html>

Checkins

- Checkins are the *Raison d'être* of Foursquare.
- They connect **Users** with **Venues** providing valuable temporal and demographic information.
- Users have the option of sharing their checkins through social media
- The status text is shared along with the URL of the web version of the checkin.
- Currently, particularly hard to obtain

Mobility and Social Structure



Airline Flights

NATS

US Flight dataset

https://www.transtats.bts.gov/databases.asp?Mode_ID=1&Mode_Desc=Aviation&Subject_ID2=0

The screenshot shows the Bureau of Transportation Statistics (BTS) website. At the top, there's a navigation bar with links for 'Explore Topics and Geography', 'Browse Statistical Products and Data', 'Learn About BTS and Our Work', and 'Newsroom'. Below the navigation is a search bar with a magnifying glass icon. The main content area features a title 'Bureau of Transportation Statistics' and a sub-section titled 'Data Library: Aviation'. This section lists various databases with their descriptions and profile links. On the left side, there's a sidebar with links for 'Resources' (Database Directory, Glossary, Upcoming Releases), 'Data Finder' (By Mode: Aviation, Maritime, Highway, Transit, Rail, Pipeline, Bike/Pedestrian, Other), and 'By Subject' (Safety, Freight Transport, Passenger Travel, Infrastructure). The URL in the browser bar is https://www.transtats.bts.gov/databases.asp?Mode_ID=1&Mode_Desc=Aviation&Subject_ID2=0.

OST-R > BTS



Search this site:

 Go

Advanced Search

Resources

Database Directory

Glossary

Upcoming Releases

Data Release History

Data Finder

By Mode

Aviation

Maritime

Highway

Transit

Rail

Pipeline

Bike/Pedestrian

Other

By Subject

Safety

Freight Transport

Passenger Travel

Infrastructure

Data Library: Aviation

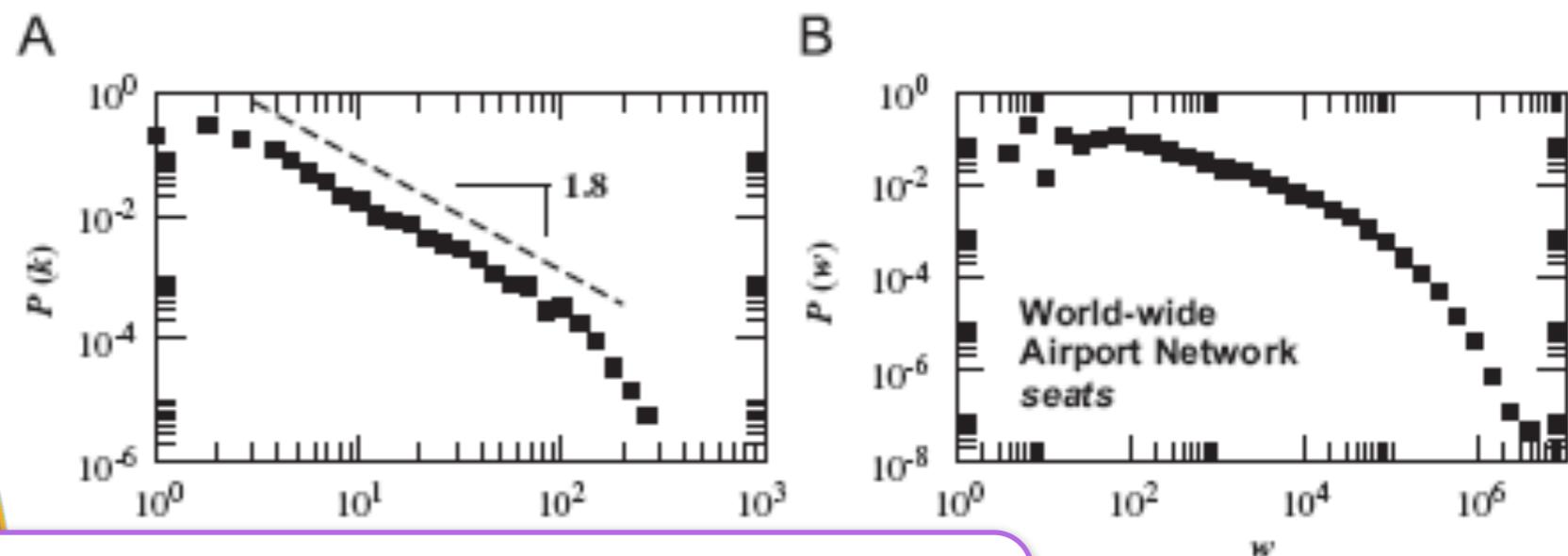
Databases Summary Tables Glossary Filter Subject All Subjects Go

<<Prev Rows 1 to 15 of 26 Next>>

Database Name	Description	Profile
Air Carrier Financial Reports (Form 41 Financial Data)	Form 41 Financial Schedule consists of financial information on large U.S. certified air carriers--includes balance sheet, cash flow, employment, income statement, fuel cost and consumption, aircraft operating expenses, and operating expenses. Note: Numbers presented on B1, B11 Balance Sheet and P11, P12 Statement of Operations now follow the format of common public financial documents. This format reverses signs from the accounting format in which numbers appeared prior to 10/18/2006 (Examples).	Profile
Air Carrier Statistics (Form 41 Traffic)- U.S. Carriers	Monthly data reported by certificated U.S. air carriers on passengers, freight and mail transported. Also includes aircraft type, service class, available capacity and seats, and aircraft hours ramp-to-ramp and airborne.	Profile
Air Carrier Statistics (Form 41 Traffic)- All Carriers	Monthly data reported by certificated U.S. and foreign air carriers on passengers, freight and mail transported. Also includes aircraft type, service class, available capacity and seats, and aircraft hours ramp-to-ramp and airborne.	Profile
Air Carrier Summary Data (Form 41 and 298C Summary Data)	Summary data of the non-stop segment and on-flight market data reported by air carriers on Form 41 and Form 298C	Profile
Airline On-Time Performance Data	Monthly data reported by US certified air carriers that account for at least one percent of domestic scheduled passenger revenues--includes scheduled and actual arrival and departure times for flights.	Profile
Airline Origin and Destination Survey (DB1B)	Origin and Destination Survey (DB1B) is a 10% sample of airline tickets from reporting carriers. Data includes origin, destination and other itinerary details of passengers transported.	Profile
American Travel Survey (ATS) 1995	National data on the nature and characteristics of long-distance personal travel, from a household survey conducted by BTS about every five years.	Profile

Long Range Mobility

PNAS 106, 21484 (2009)



What about Short Range Mobility?

Complete IATA and OAG databases:

3362 airports worldwide

220 countries

> 20,000 connections

w_{ij} #passengers on connection i-j

>99% total traffic

Commuting



@bgoncalves

<http://youtube.com/watch?v=FqkinE8khZs>

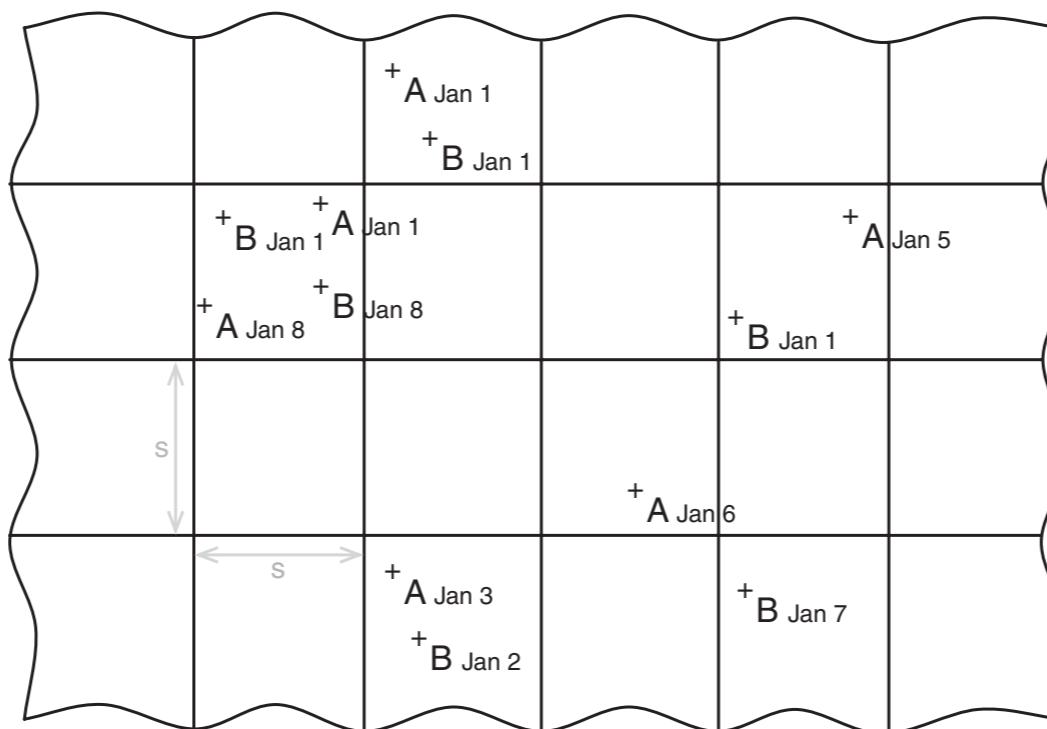
Subscribe:
[youtube.com/tanvideo11](https://www.youtube.com/tanvideo11)

www.bgoncalves.com

Co-occurrences and Social Ties

PNAS 107, 22436 (2010)

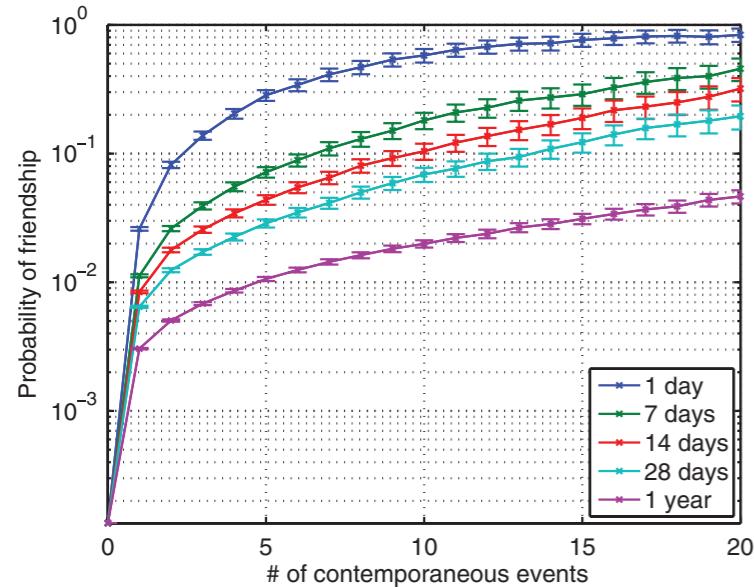
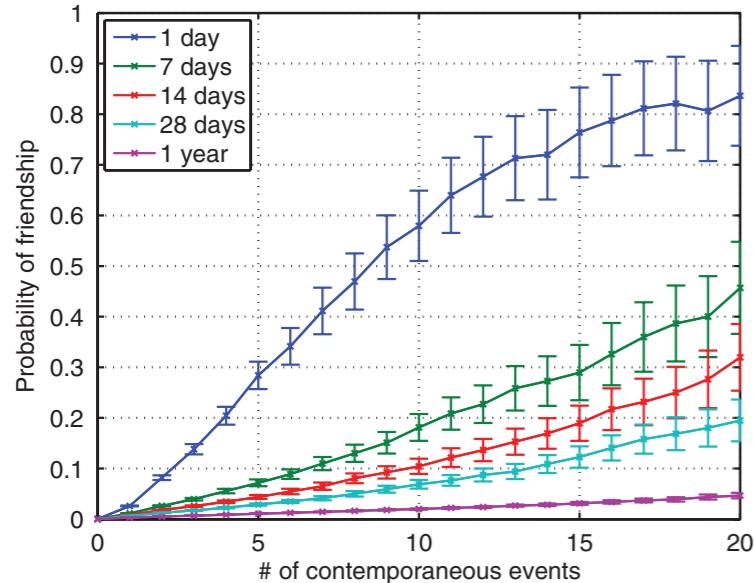
- Geotagged Flickr Photos
- Divide the world into a grid



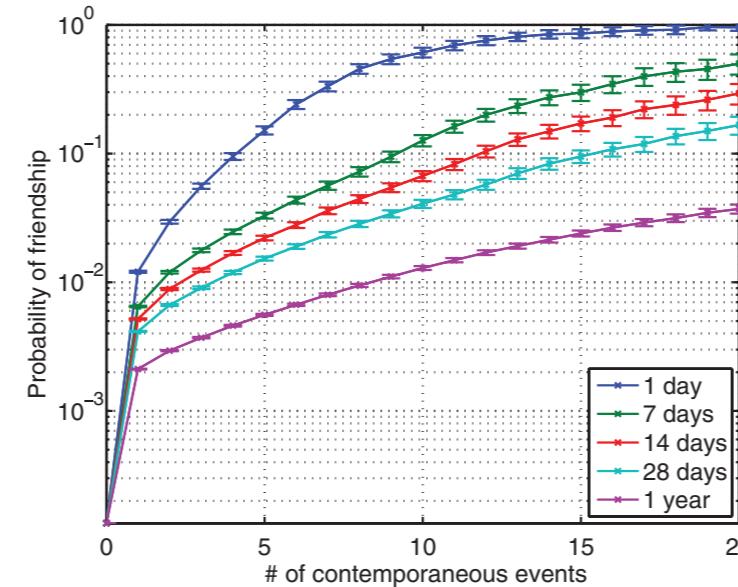
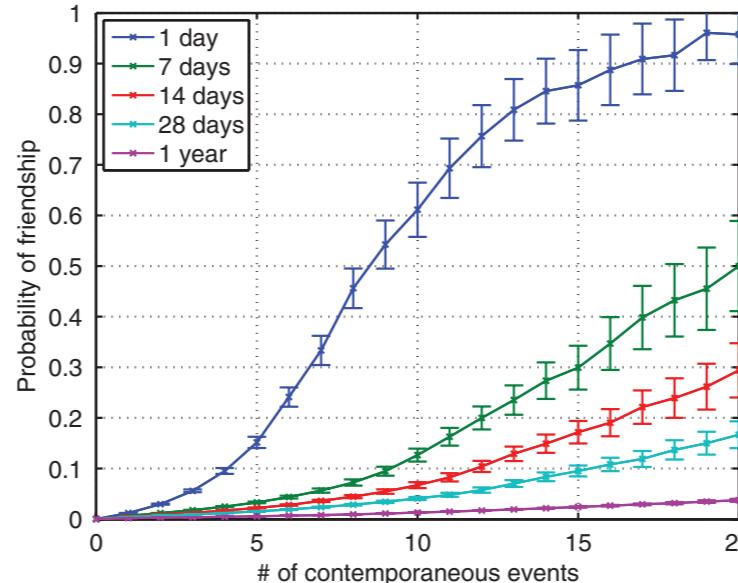
Count number of cells on which two individuals were within a given interval

Co-occurrences and Social Ties

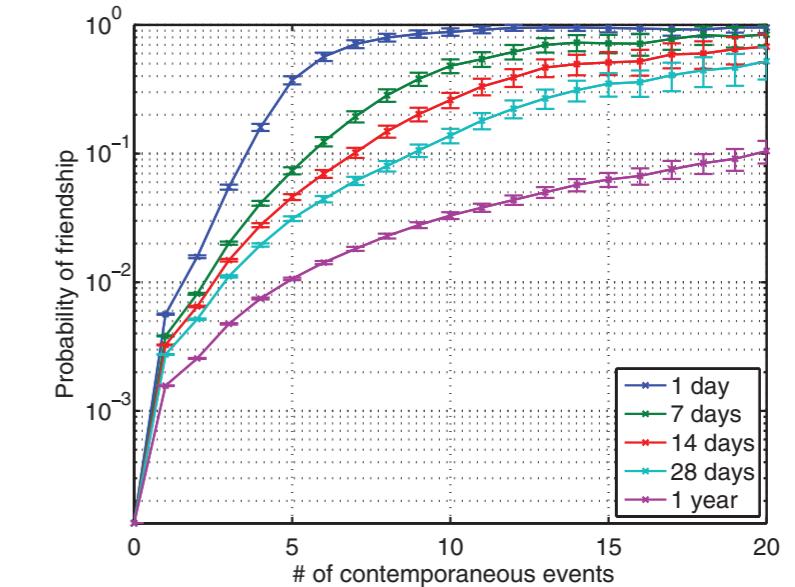
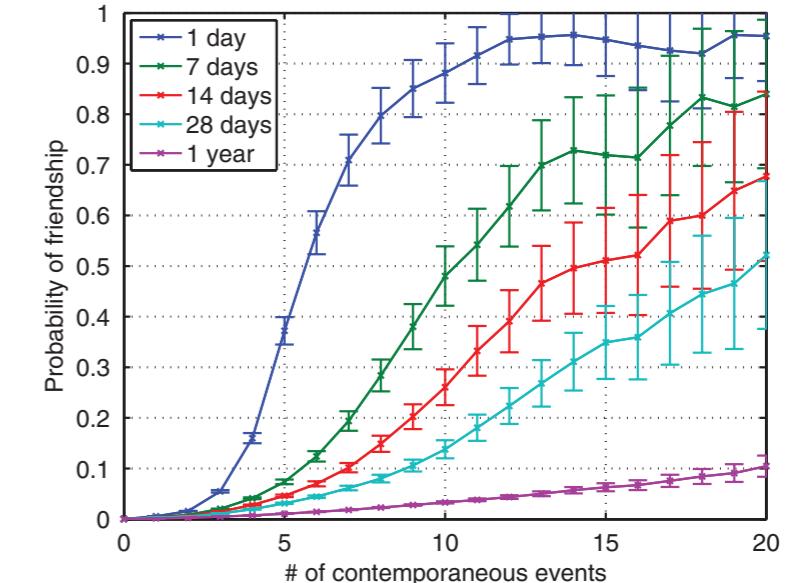
PNAS 107, 22436 (2010)



A $s = 0.001^\circ$



B $s = 0.01^\circ$



C $s = 0.1^\circ$