

BIG DATA ANALYTICS

Fosca Giannotti and Roberto Trasarti
Pisa KDD Lab, ISTI-CNR & Univ. Pisa

<http://www-kdd.isti.cnr.it/>



DIPARTIMENTO DI INFORMATICA Università di Pisa
anno accademico 2015/2016

BRIGHT 2015 - La notte dei ricercatori in Toscana

••• 25 settembre 2015 •••



PROGRAMMA
DEGLI EVENTI
AL CNR DI PISA

<http://nottedeiricercatori.pisa.it/> • via Giuseppe Moruzzi 1 - 56124 PISA

- 17-19 “O Bog data è meglio Pelè”
- 18:30-19:30 La COOP sei tu..dimmi cosa compri ti dirà chi sei

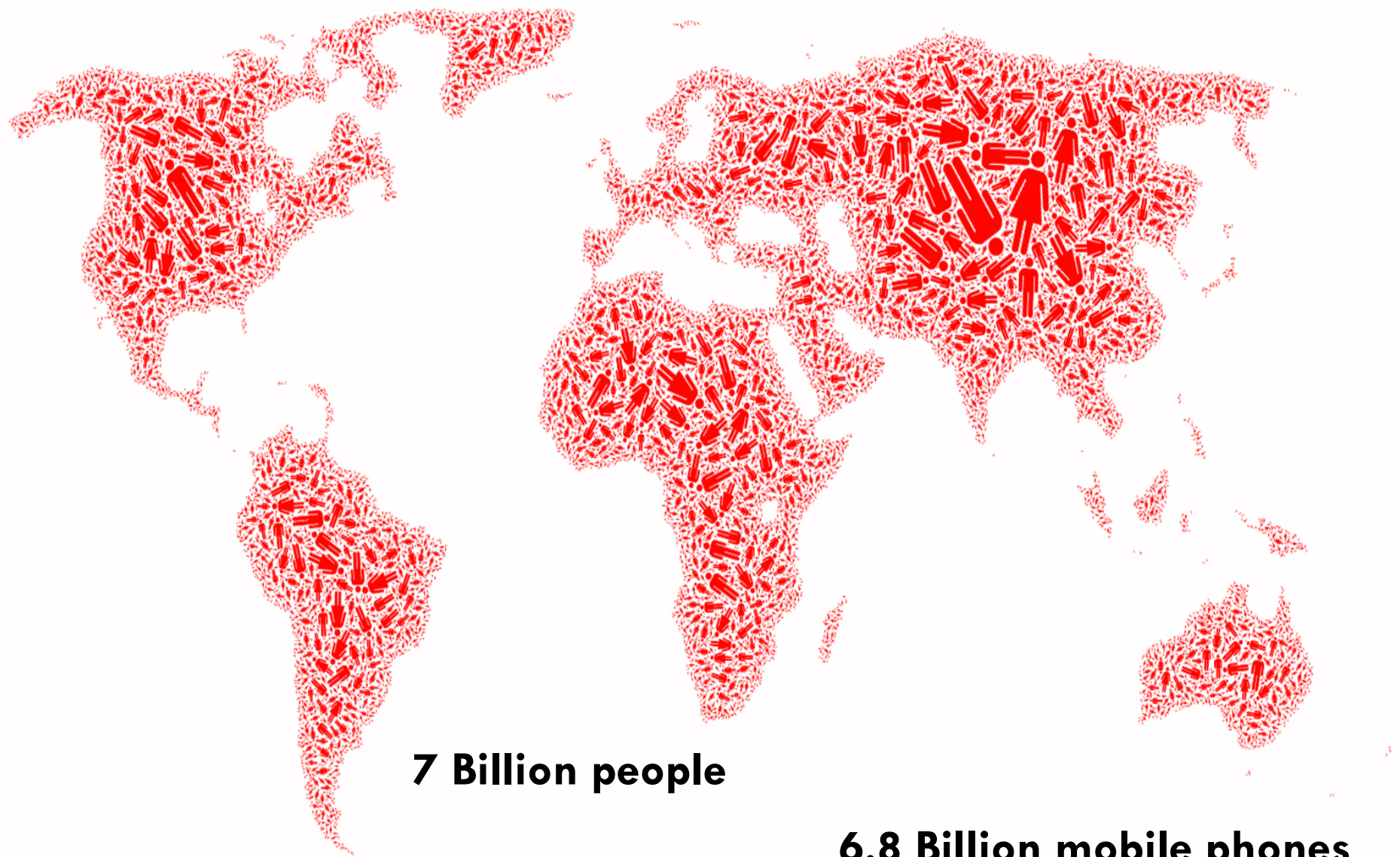
LECTURE 1: SOCIAL MINING & BIG DATA

OPPORTUNITY & RISKS

FOSCA GIANNOTTI

KNOWLEDGE DISCOVERY & DATA
MINING LAB. – ISTI CNR, PISA





7 Billion people

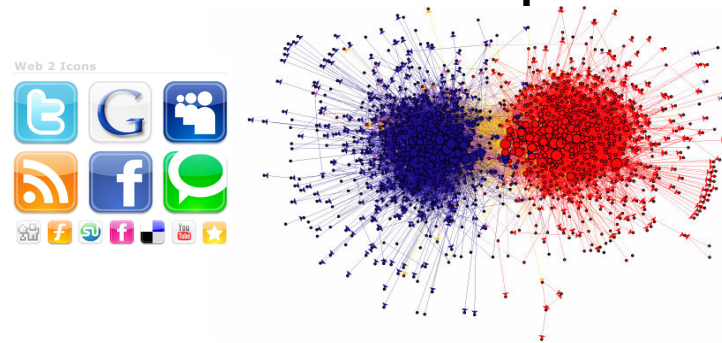
6.8 Billion mobile phones

Digital Footprints of Human Activities

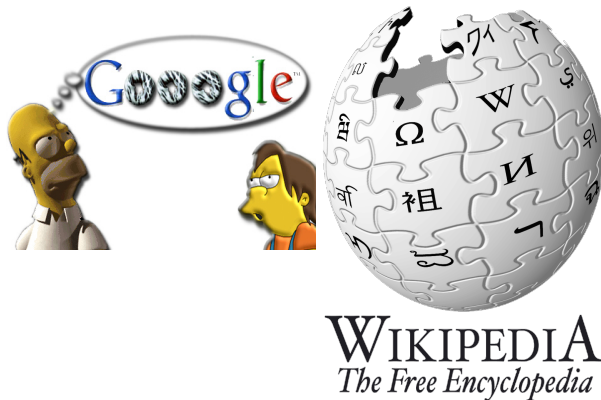
Shopping patterns & lifestyle



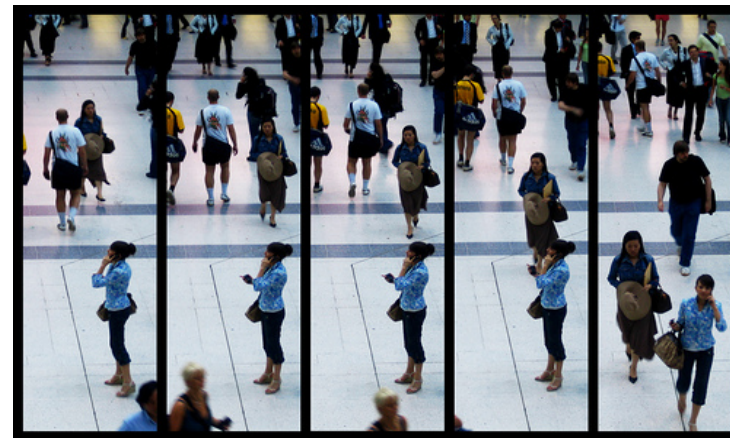
Relationships & social ties

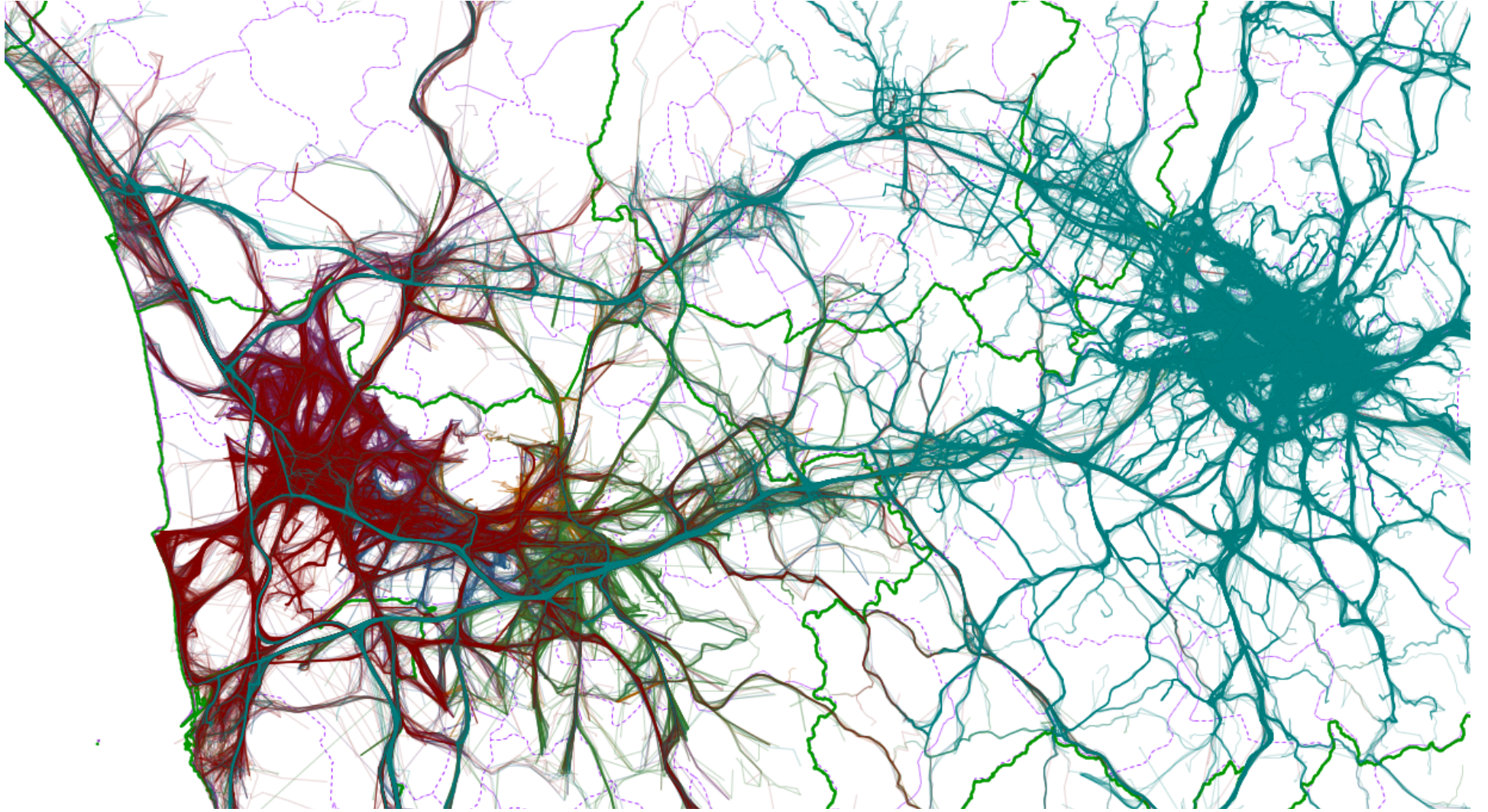


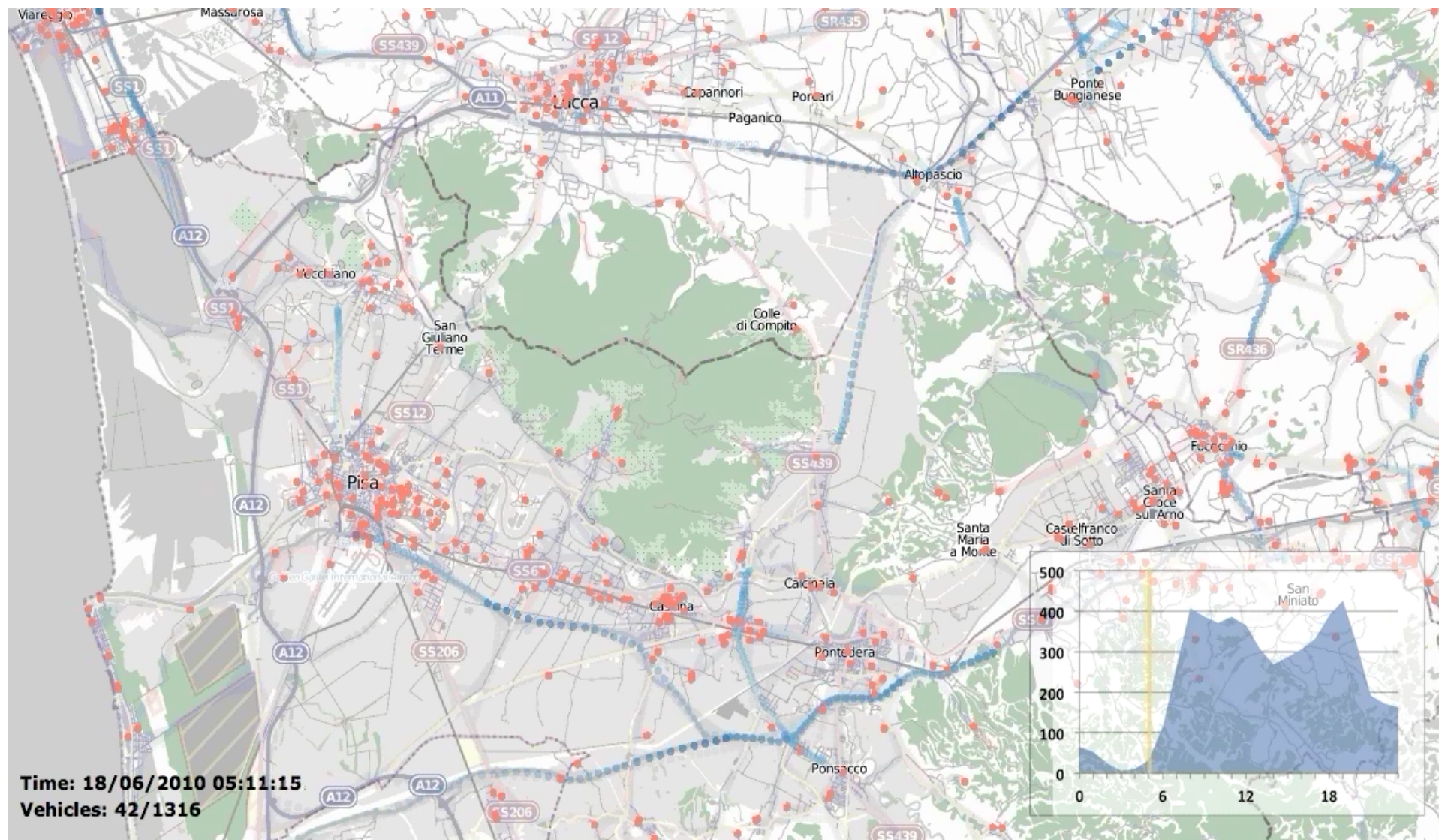
Desires, opinions, sentiments

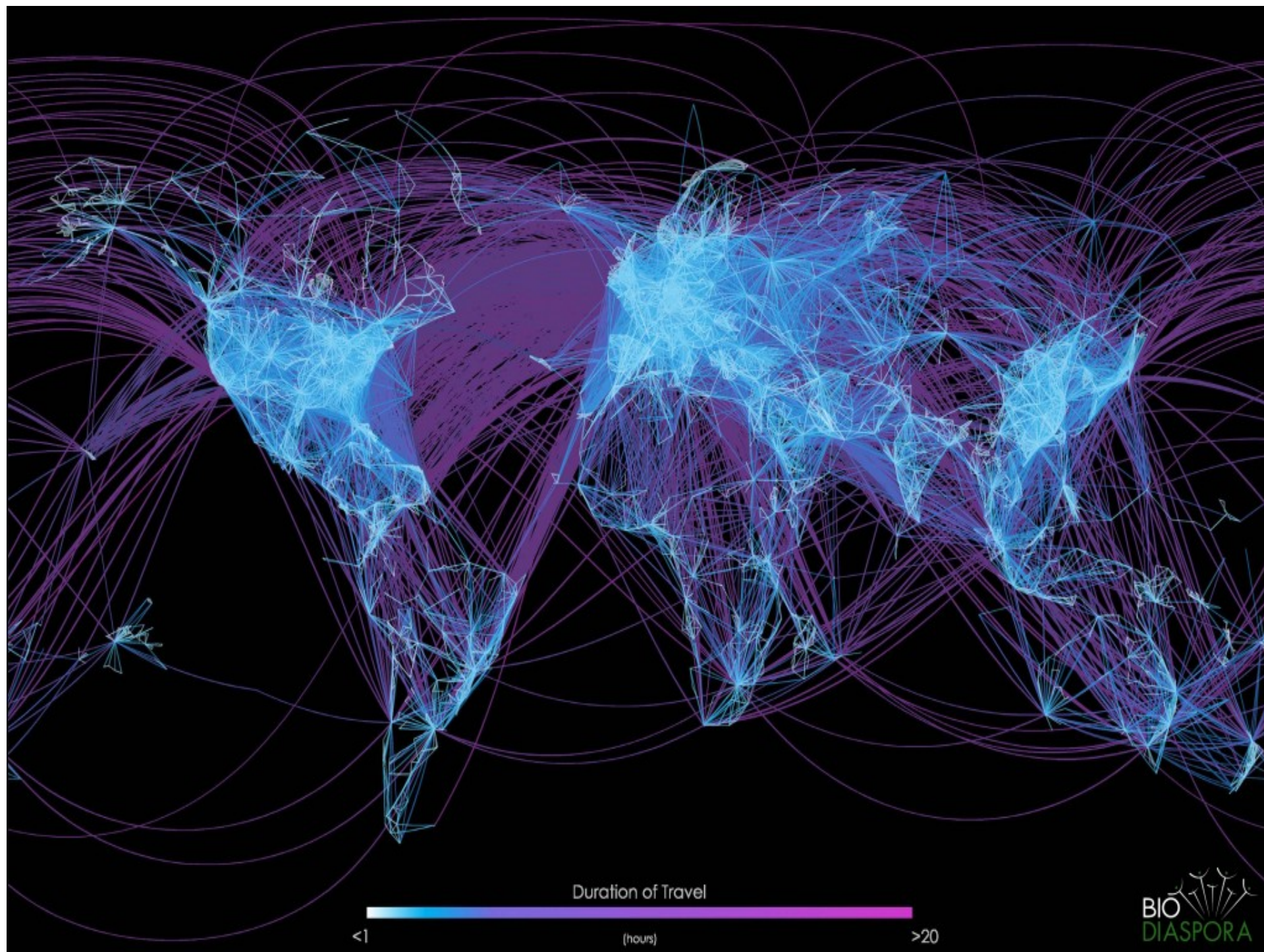


Movements









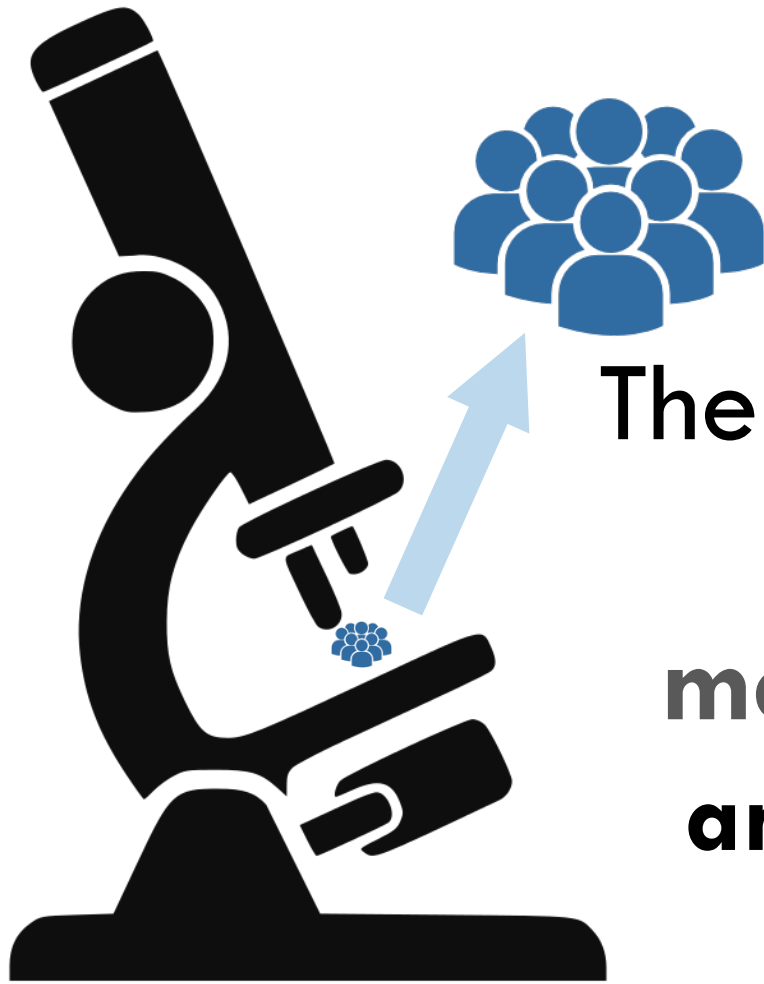




POLLICINI DIGITALI

La Vita Nova, e-magazine de Il Sole 24 Ore
Fosca Giannotti, Dino Pedreschi

Big Data & Social Mining



The Social Microscope:
a tool to
measure, understand,
and possibly predict
human behavior

SOCIAL MINING: MAKING SENSE OF BIG DATA



BIG DATA & NEW QUESTIONS TO ASK



Google Flu Trends



Detecting influenza epidemics using search engine query data

Jeremy Ginsberg¹, Matthew H. Mohebbi¹, Rajan S. Patel¹, Lynnette Brammer², Mark S. Smolinski¹ & Larry Brilliant¹

¹Google Inc. ²Centers for Disease Control and Prevention

Nature 457, 1012-1014 (2009)

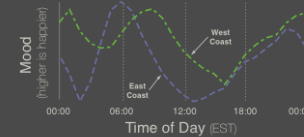
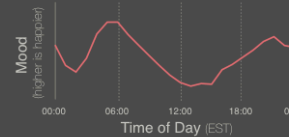
Pulse of the Nation: U.S. Mood Throughout the Day inferred from Twitter

All times are Eastern Standard Time (EST)



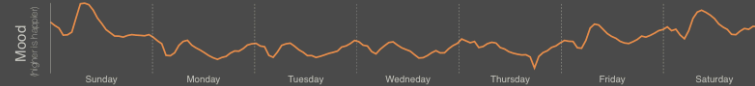
Mood Variations

A number of interesting trends can be observed in the data. First, overall daily variations can be seen (first graph), with the early morning and late evening having the highest level of happiness. Second, geographic variations can be observed (second graph), with a significantly happier west coast that is consistently three hours behind the east coast.



Weekly Variations

Weekly trends can be observed as well, with weekends much happier than weekdays.

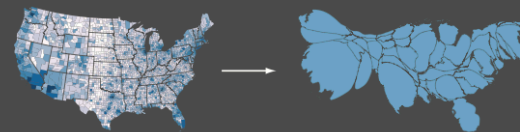


About the Data and Visualization

The plots were calculated using over 300 million tweets (Sep 2006 – Aug 2009) collected by MPI-SWS researchers, represented as density-preserving cartograms. The mood of each tweet was inferred using ANEW word list (Bradley, M.M., & Lang, P.J. *Affective norms for English words (ANEW): Stimuli, instruction manual and affective ratings*. Technical report C-1, The Center for Research in Psychophysiology, University of Florida). County area data was taken from the U.S. Census Bureau at <http://factfinder.census.gov>, and the base U.S. map was taken from Wikimedia Commons. User locations were inferred using the Google Maps API, and mapped into counties using PostGIS and U.S. county maps from the U.S. National Atlas. Mood colors were selected using Color Brewer 2.

About Cartograms

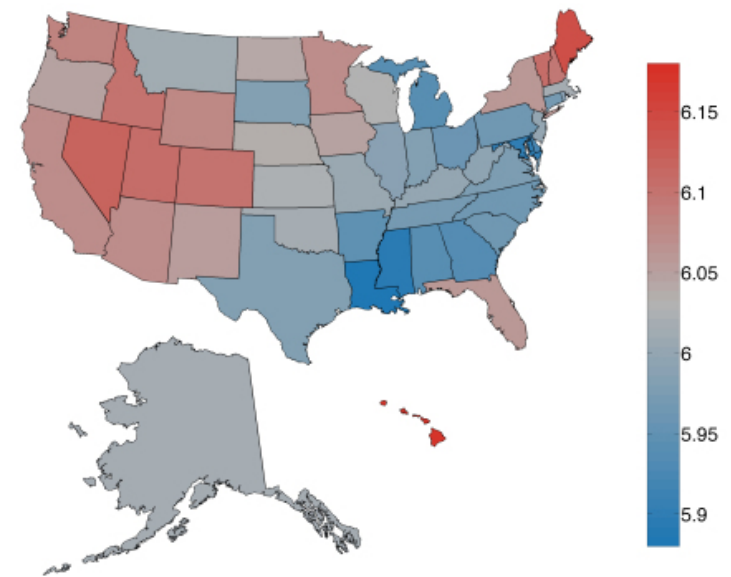
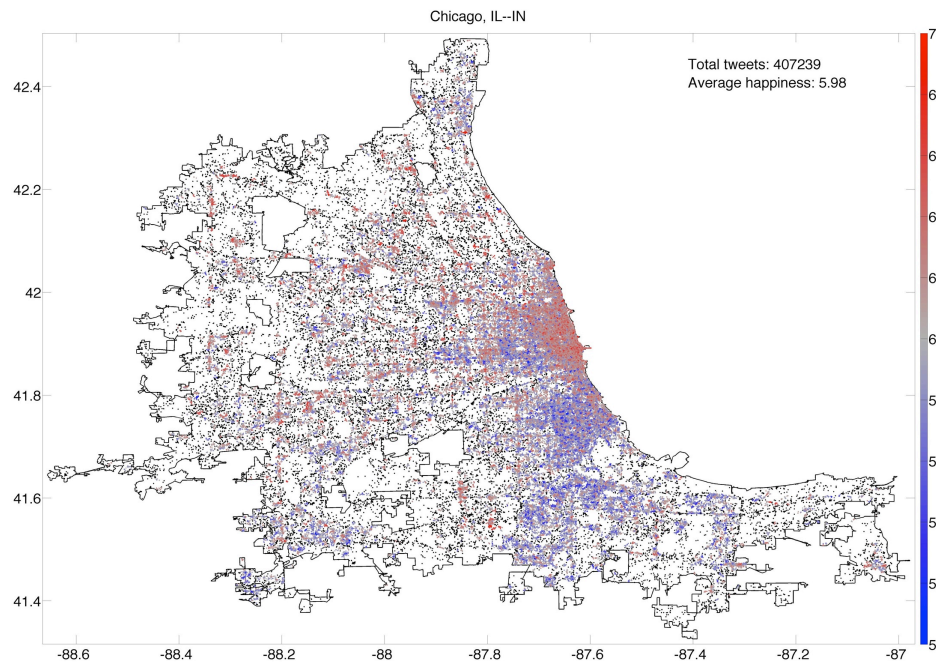
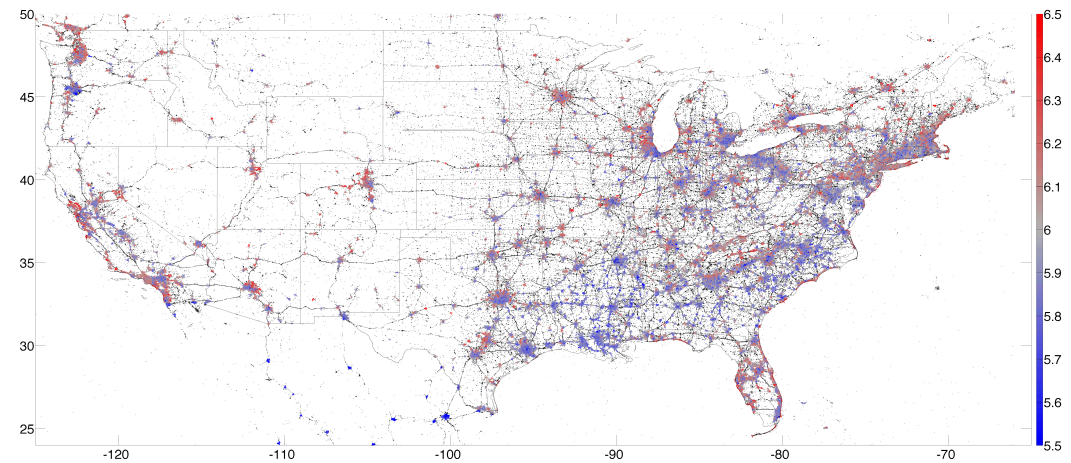
A cartogram is a map in which the mapping variable (in this case, the number of tweets) is substituted for the true land area. Thus, the geometry of the actual map is altered so that the shape of each region is maintained as much as possible, but the area is scaled in order to be proportional to the number of tweets that originate in that region. The result is a density-equalizing map. The cartograms in this work were generated using the cart software by Mark E. J. Newman.



Northeastern University
College of Computer and Information Science[†]
Center for Complex Network Research[‡]

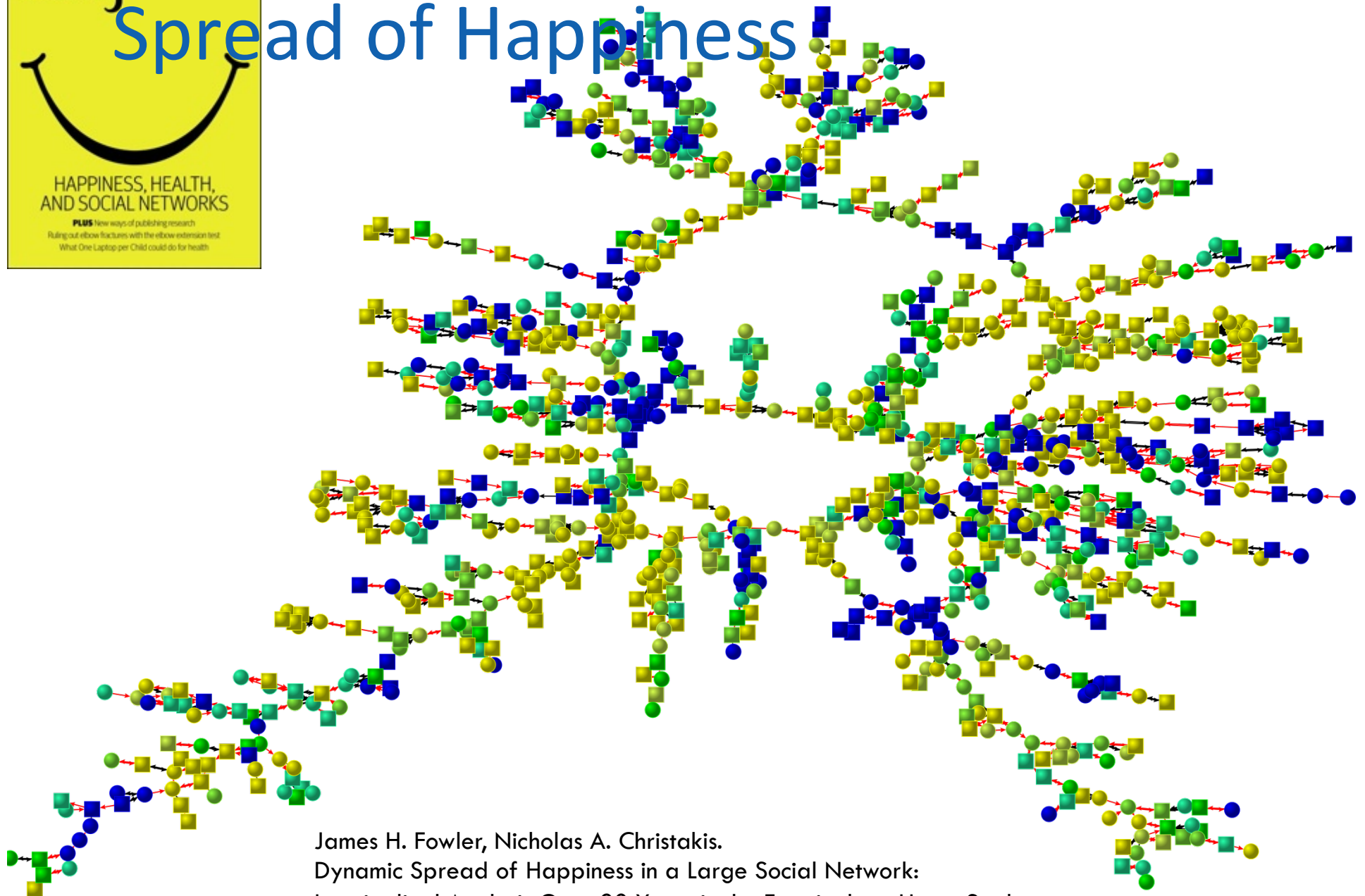
HARVARD UNIVERSITY[§]

where-is-the-happiest-city-in-the-usa?

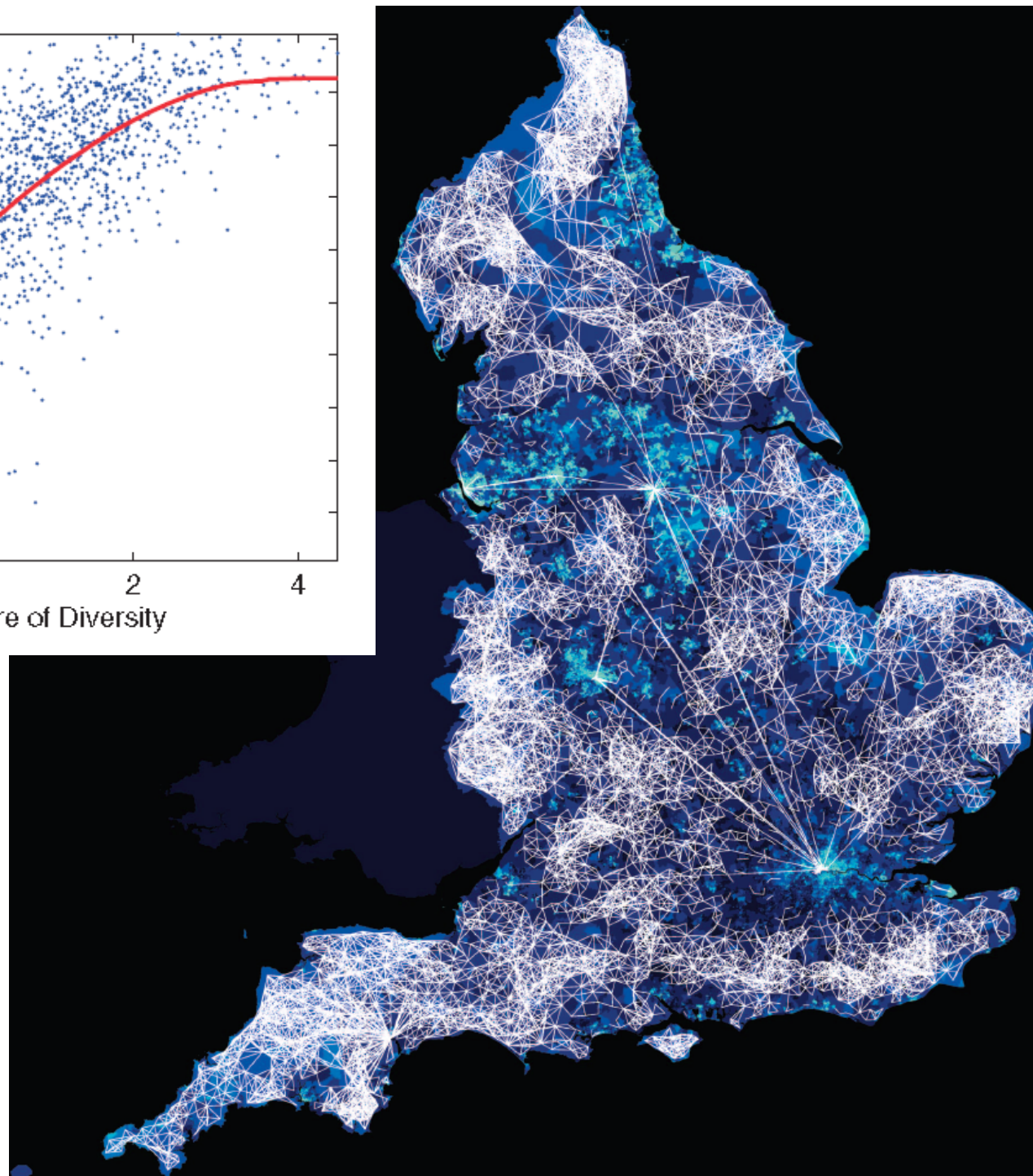
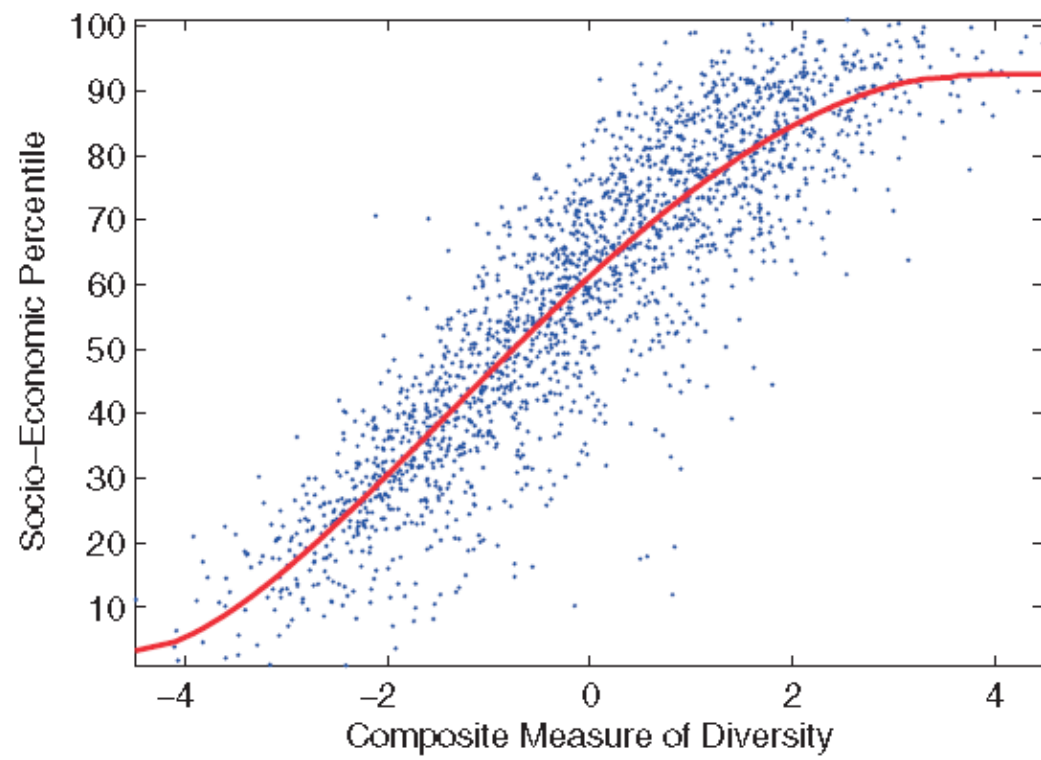


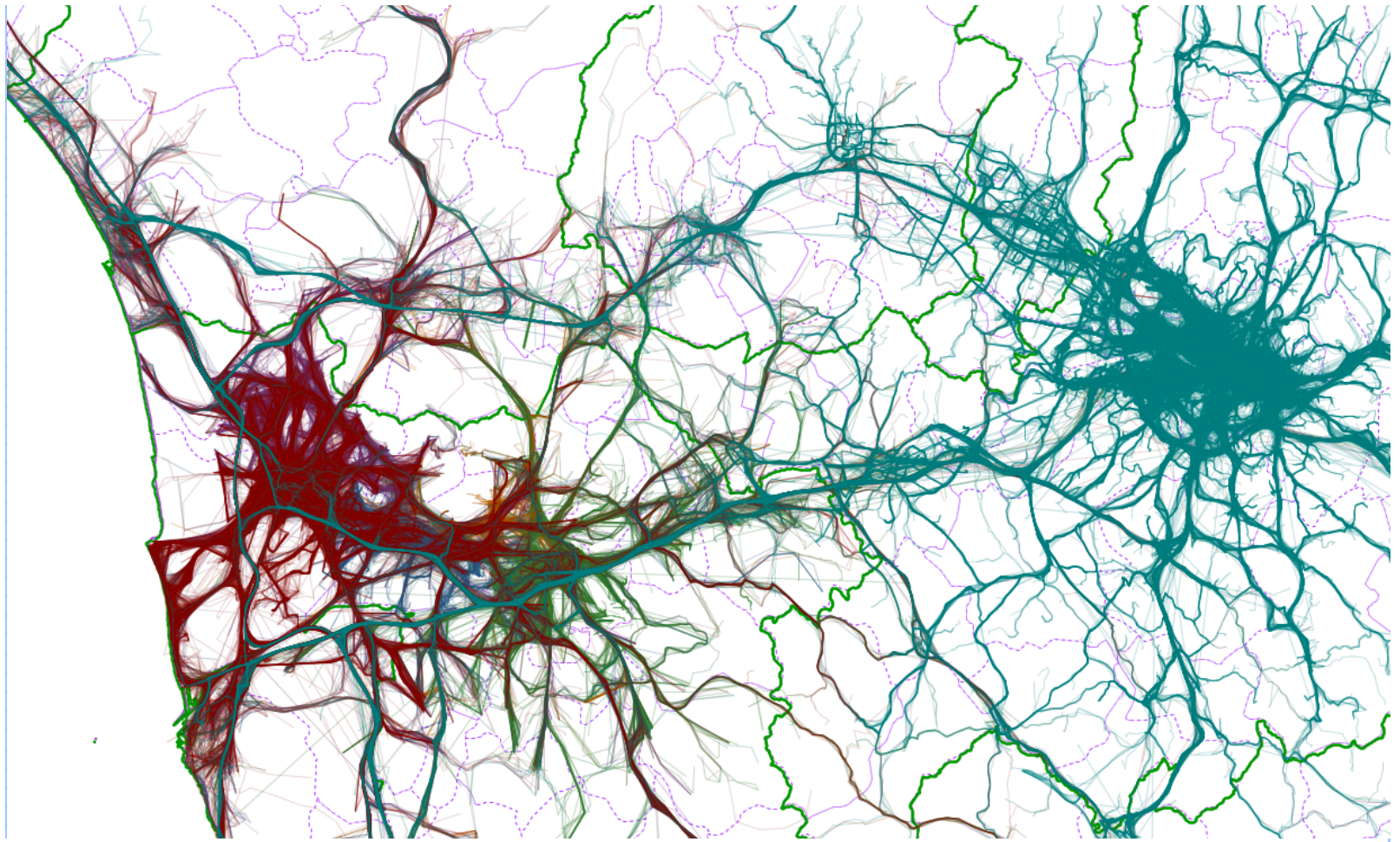


Spread of Happiness



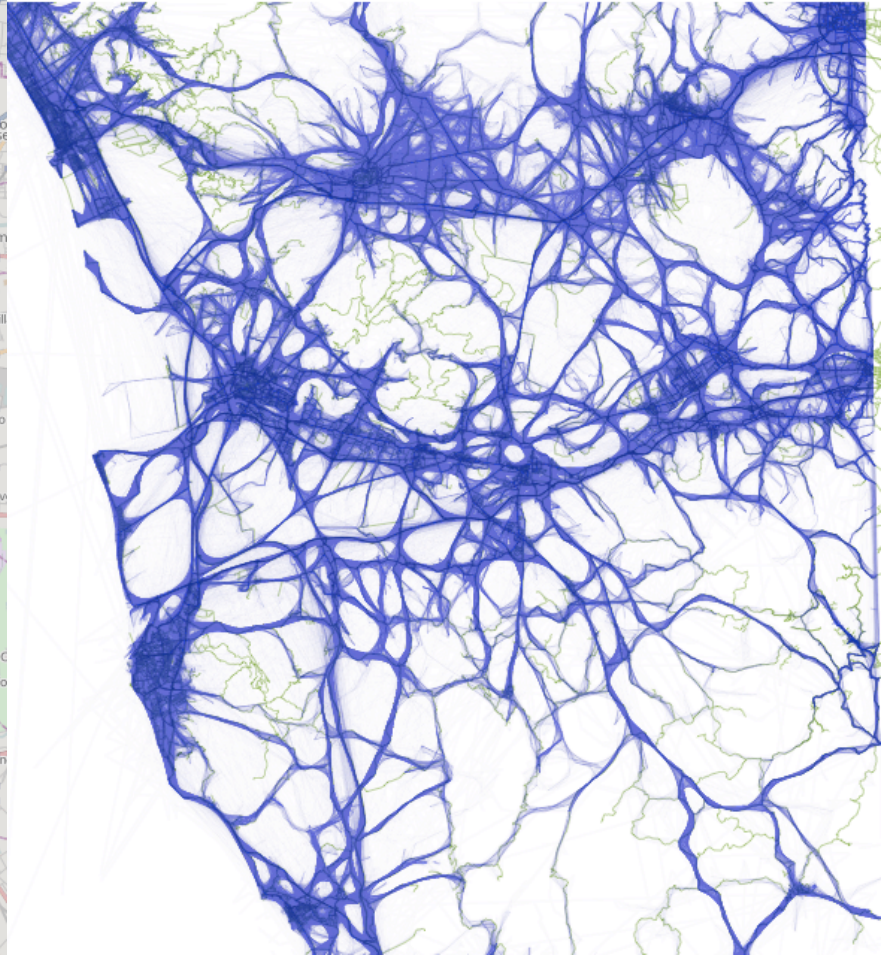
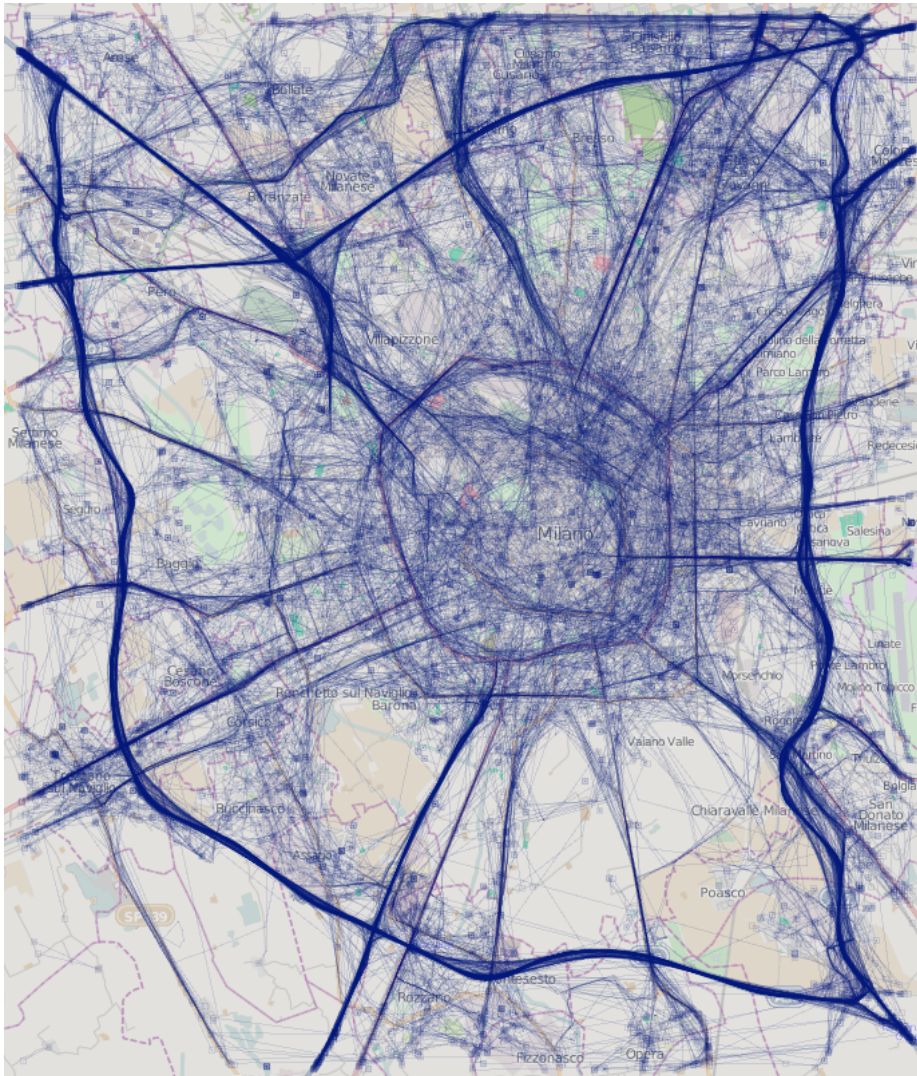
James H. Fowler, Nicholas A. Christakis.
Dynamic Spread of Happiness in a Large Social Network:
Longitudinal Analysis Over 20 Years in the Framingham Heart Study
British Medical Journal 337 (4 December 2008)



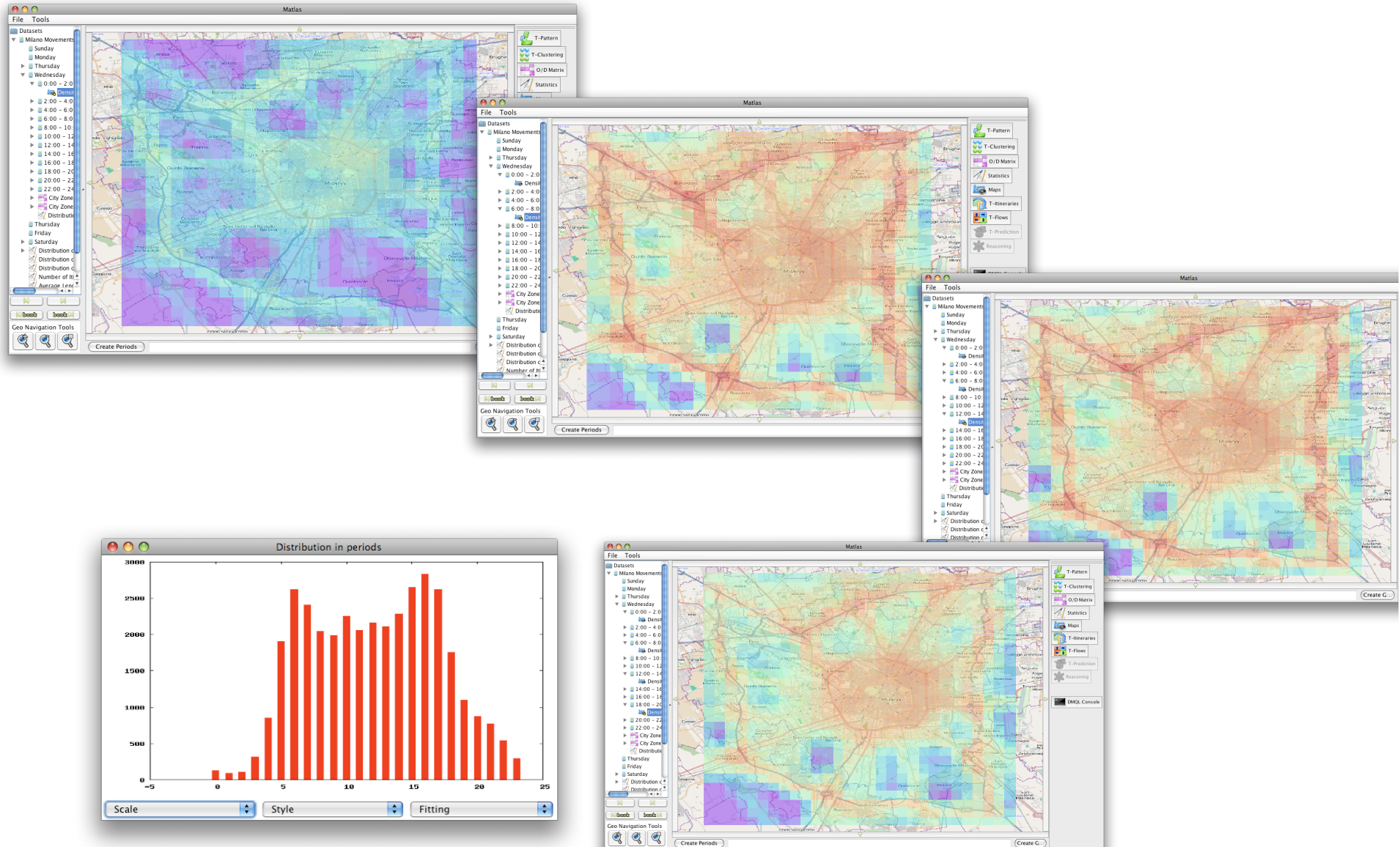


Mobility Analytics for Science of the Cities

Big Data for Smart Cities



How do people move during the day?



Big Data for Urban Mobility Atlas

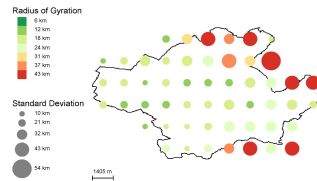
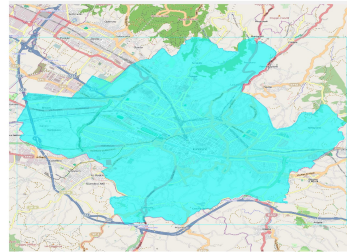
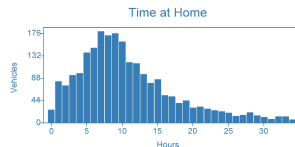
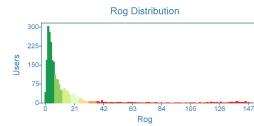
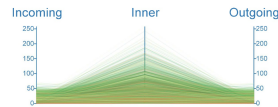
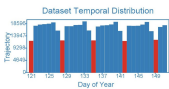
Firenze

Surface area: 106 km²

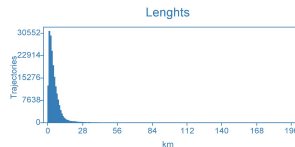
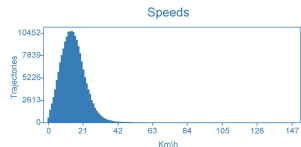
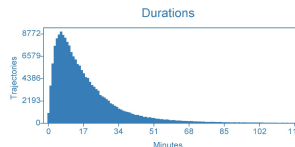
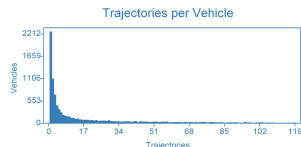
Coordinates: 43,78 11,24

Vehicles: 32.752

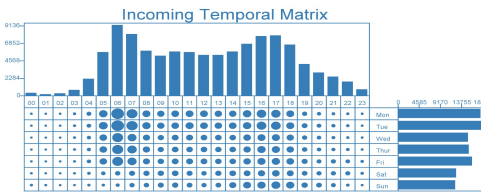
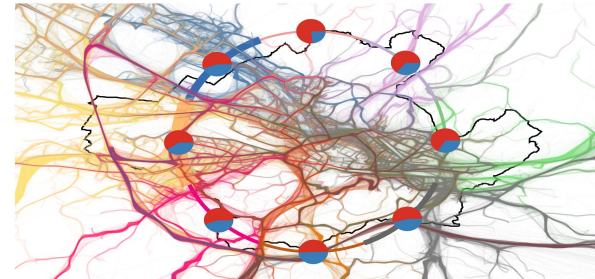
From: 2011-05-01 To: 2011-05-31



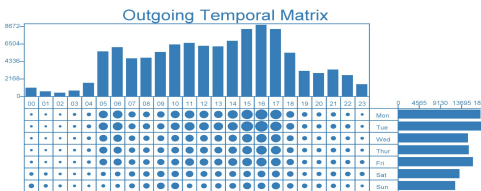
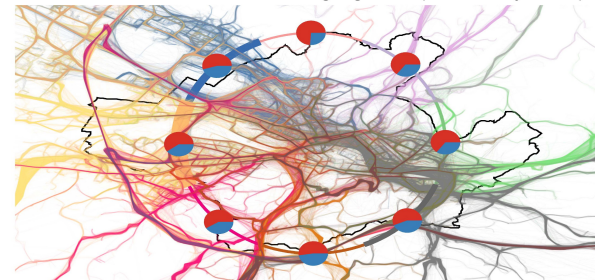
Inner Traffic (145.034 Trajectories)



Incoming Traffic (105.881 Trajectories)



Outgoing Traffic (106.418 Trajectories)



	City	Traj	Perc
NORD 13%	Sesto Fiorent.	5.123	85%
	Catanzaro	1.434	44%
	Vaglia	758	81%
	Campi Bisenzio	436	6%
	Borgo San Lore.	380	46%
OVEST 50%	Scandicci	13.447	98%
	Campi Bisenzio	6.058	93%
	Prato	6.048	94%
	Sesto Fiorent.	4.824	34%
	Lastra a Signa	2.342	99%
SUD 16%	Impuneta	3.983	87%
	San Casciano I.	1.838	75%
	Figine Valder	1.190	81%
	Grove in Chian.	898	36%
	Tavernette Val.	744	93%
EST 19%	Bagno a Ripol.	7.314	93%
	Fiesole	3.970	95%
	Pontassieve	2.787	97%
	Grove in Chian.	1.516	63%
	Rignano sull'A.	774	92%

Regular VS Occasional

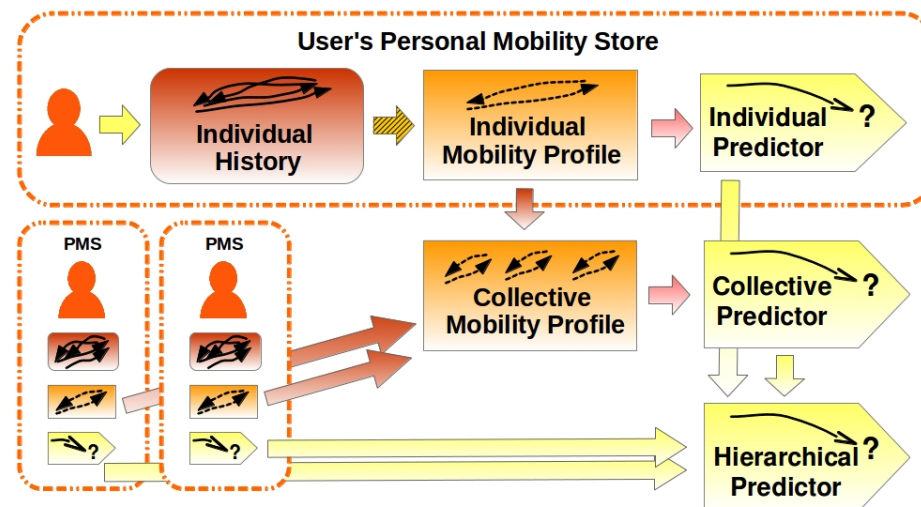
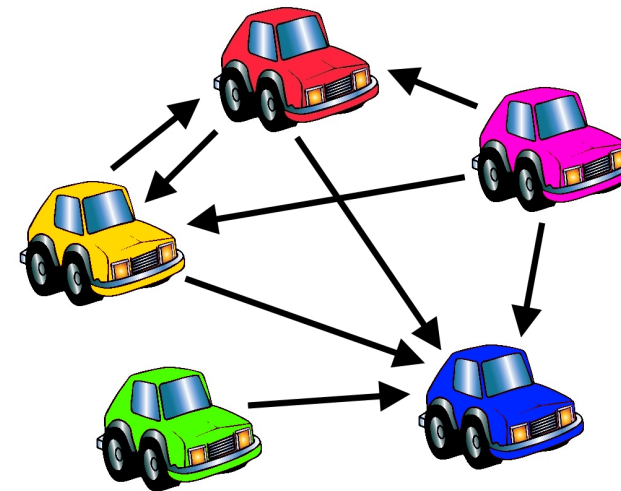
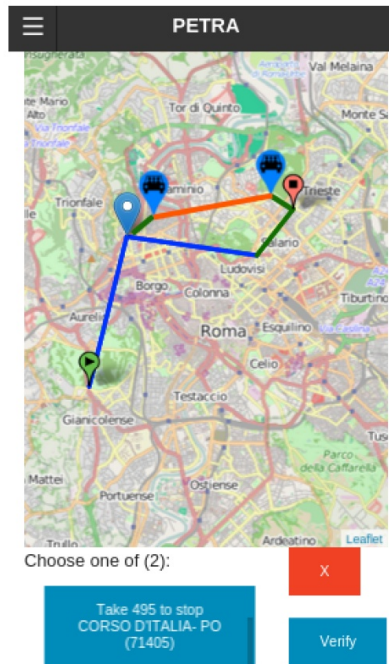
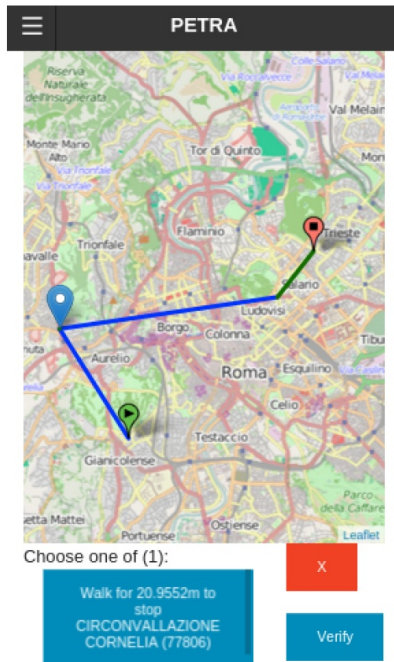


	City	Traj	Perc
NORD 12%	Sesto Fiorent.	8.358	60%
	Catanzaro	1.235	36%
	Vaglia	845	87%
	Campi Bisenzio	487	7%
	Borgo San Lore.	456	54%
OVEST 52%	Scandicci	13.439	98%
	Prato	6.166	95%
	Campi Bisenzio	5.845	92%
	Sesto Fiorent.	5.521	39%
	Lastra a Signa	2.423	98%
SUD 14%	Impuneta	3.985	85%
	San Casciano I.	1.801	72%
	Figine Valder	1.155	77%
	Tavernette Val.	742	92%
	Grove in Chian.	739	29%
EST 20%	Bagno a Ripol.	7.701	96%
	Fiesole	3.982	94%
	Pontassieve	2.806	98%
	Grove in Chian.	1.670	67%
	Rignano sull'A.	818	96%

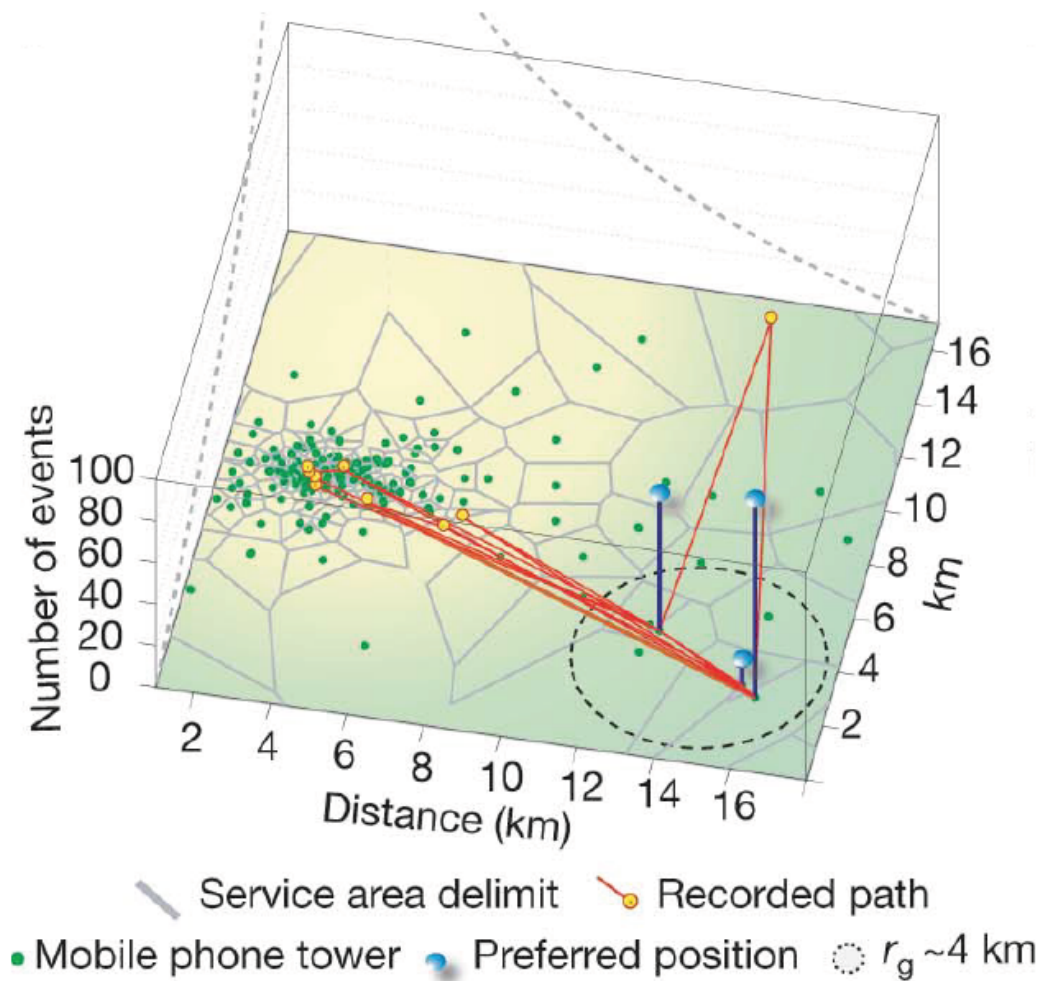
Regular VS Occasional



Personal mobility assistant



Focus on country-wide CDR data



when
you
call

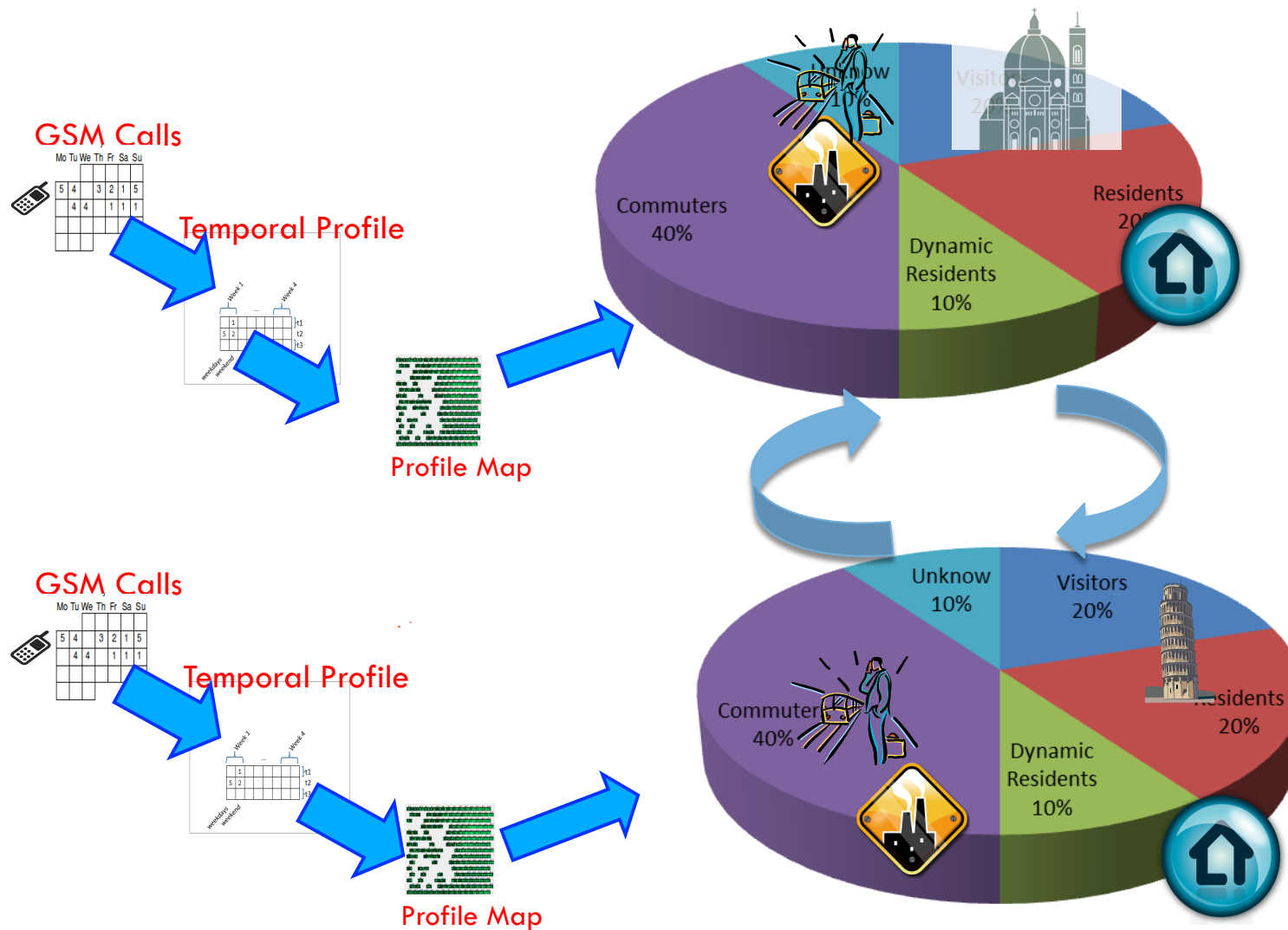


where
you
call

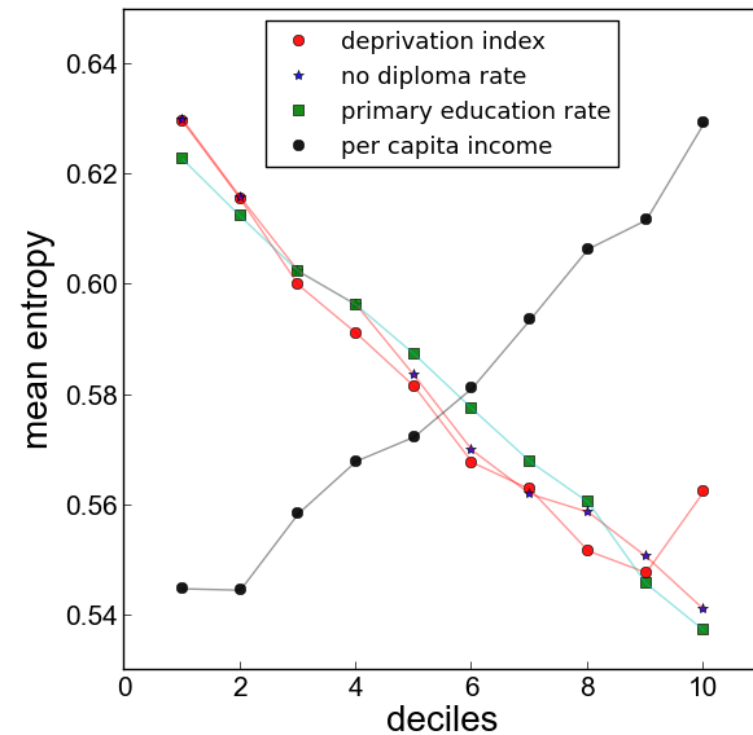
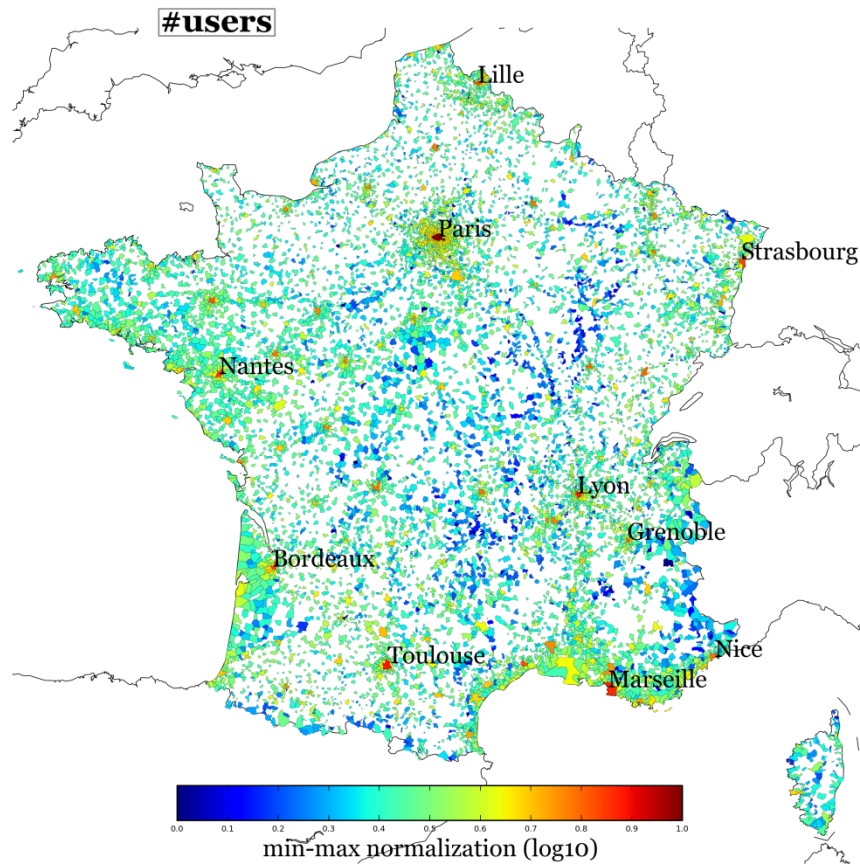


who
you
call

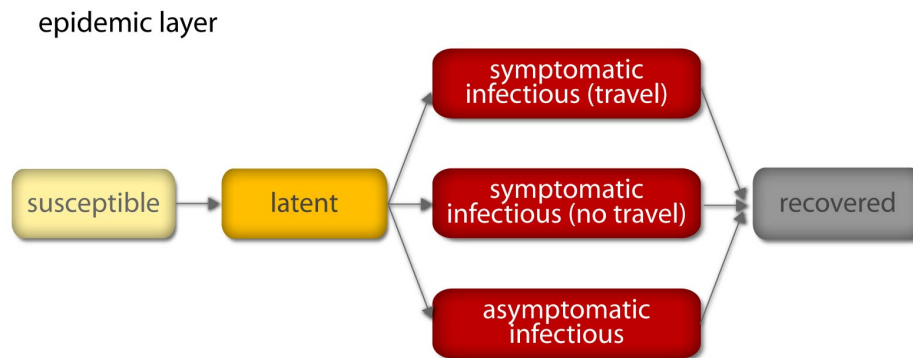
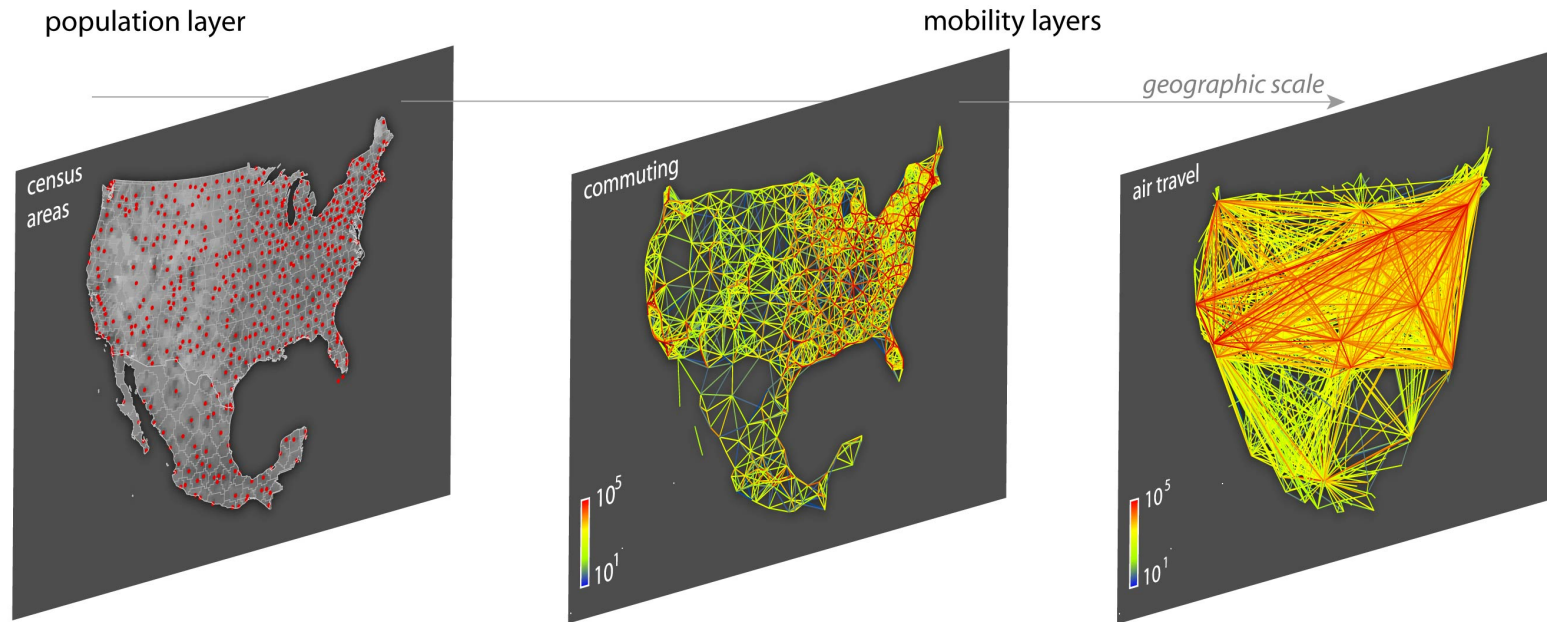
Big Data in Official Statistics: Persons & Places



Big Data: mobility diversity and wellbeing

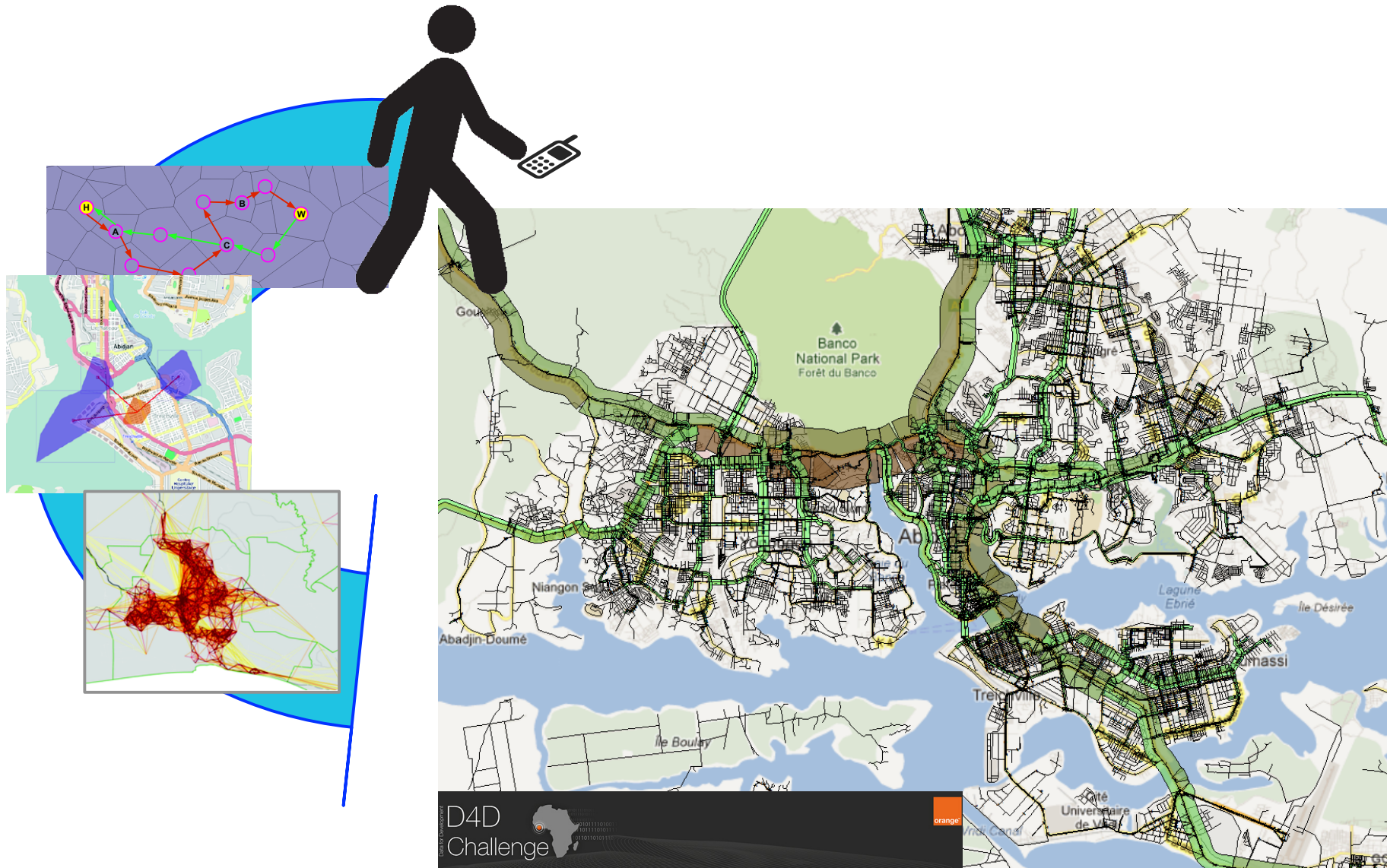


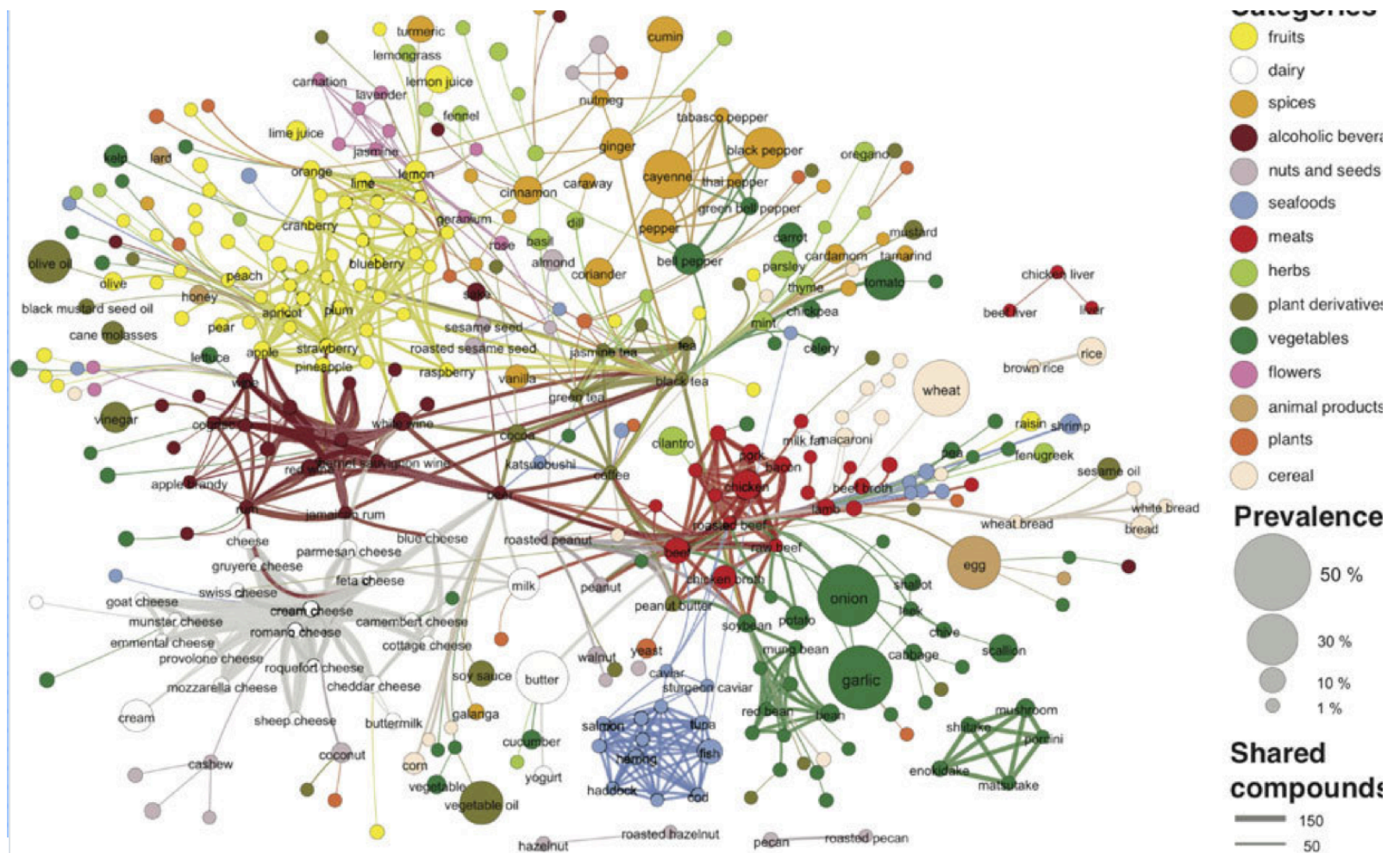
Big Data and epidemics



Parameter	Value	Description
β	from R_0	transmission probability
ε^{-1}	1.9 [1.1-2.5] d	average latency period
μ^{-1}	3 [3-5] d	average infectious period
p_t	50%	probability of traveling for infectious individuals
p_a	33%	probability of being asymptomatic
r_β	50%	relative infectiousness of asymptomatic infectious individuals

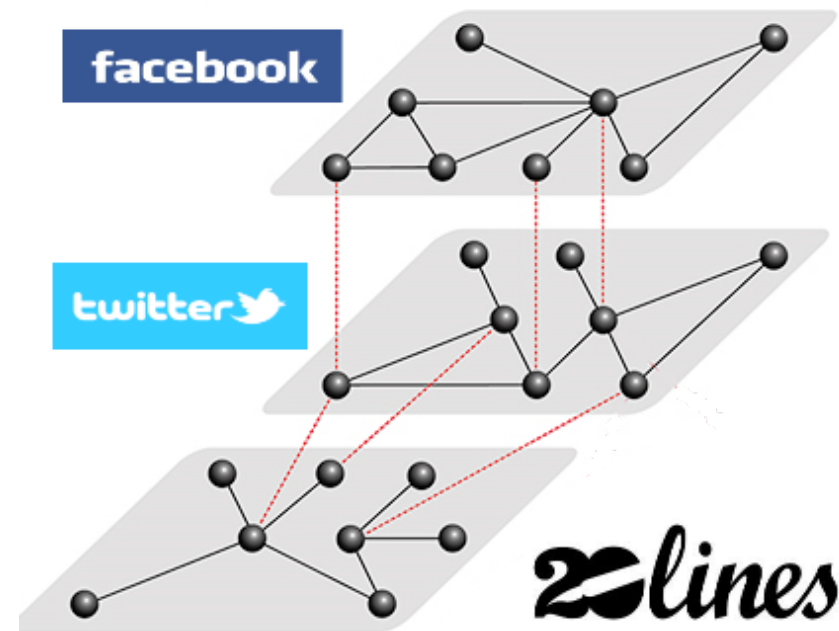
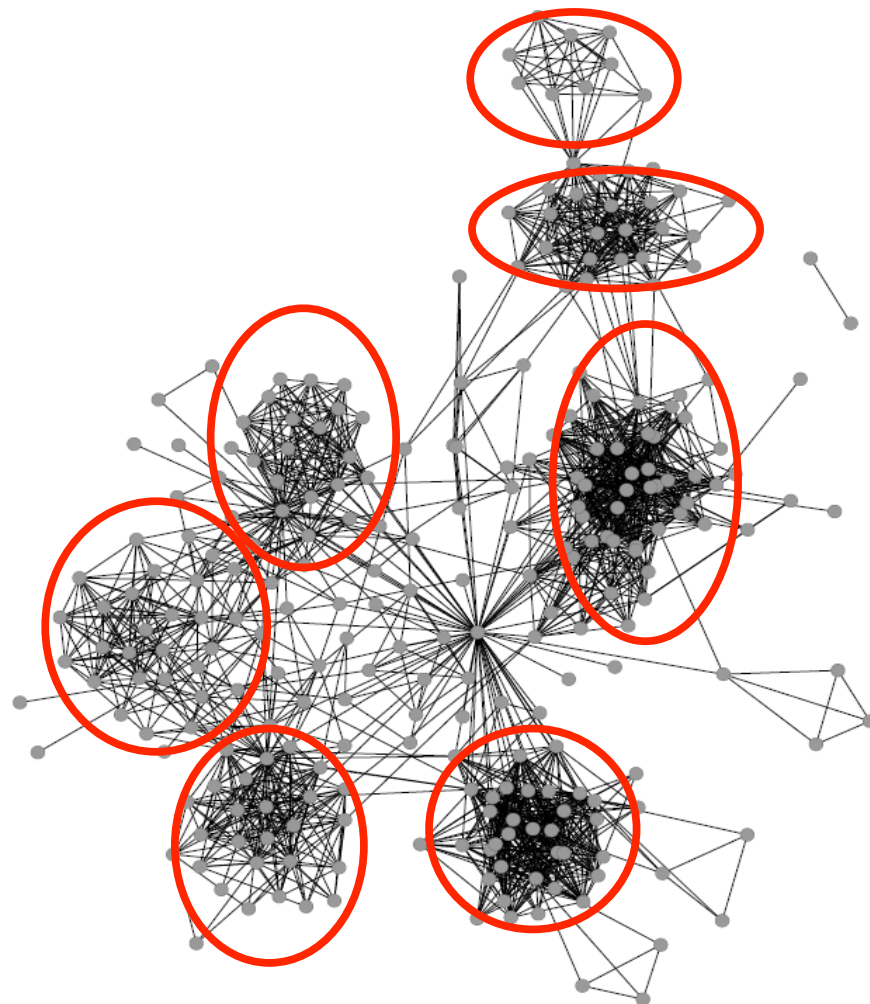
Big Data for Developing Countries



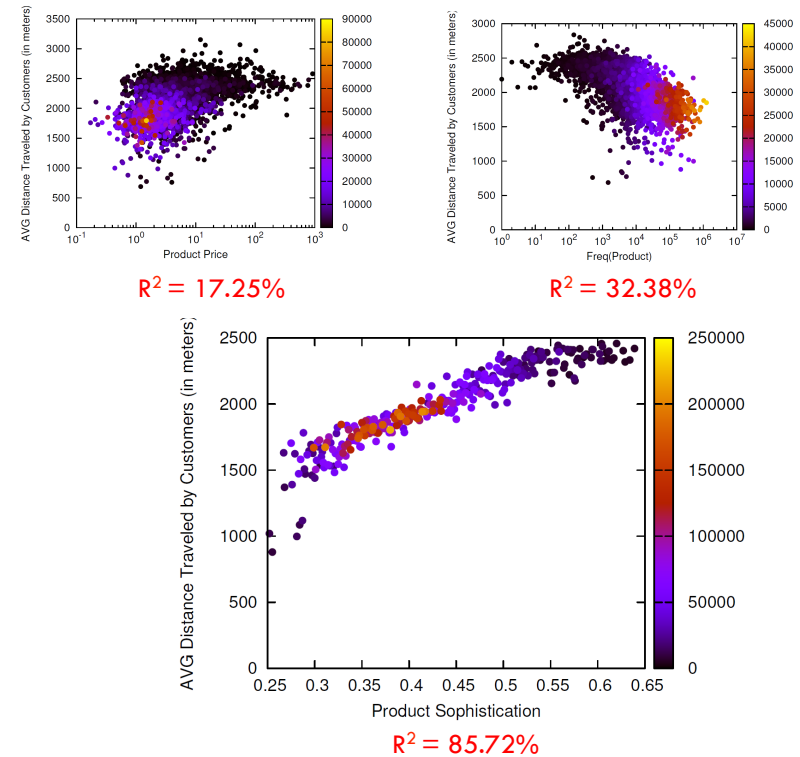


Social Network Analysis

Community Discovery, Evolution, Diffusion, Multidimensionality,...



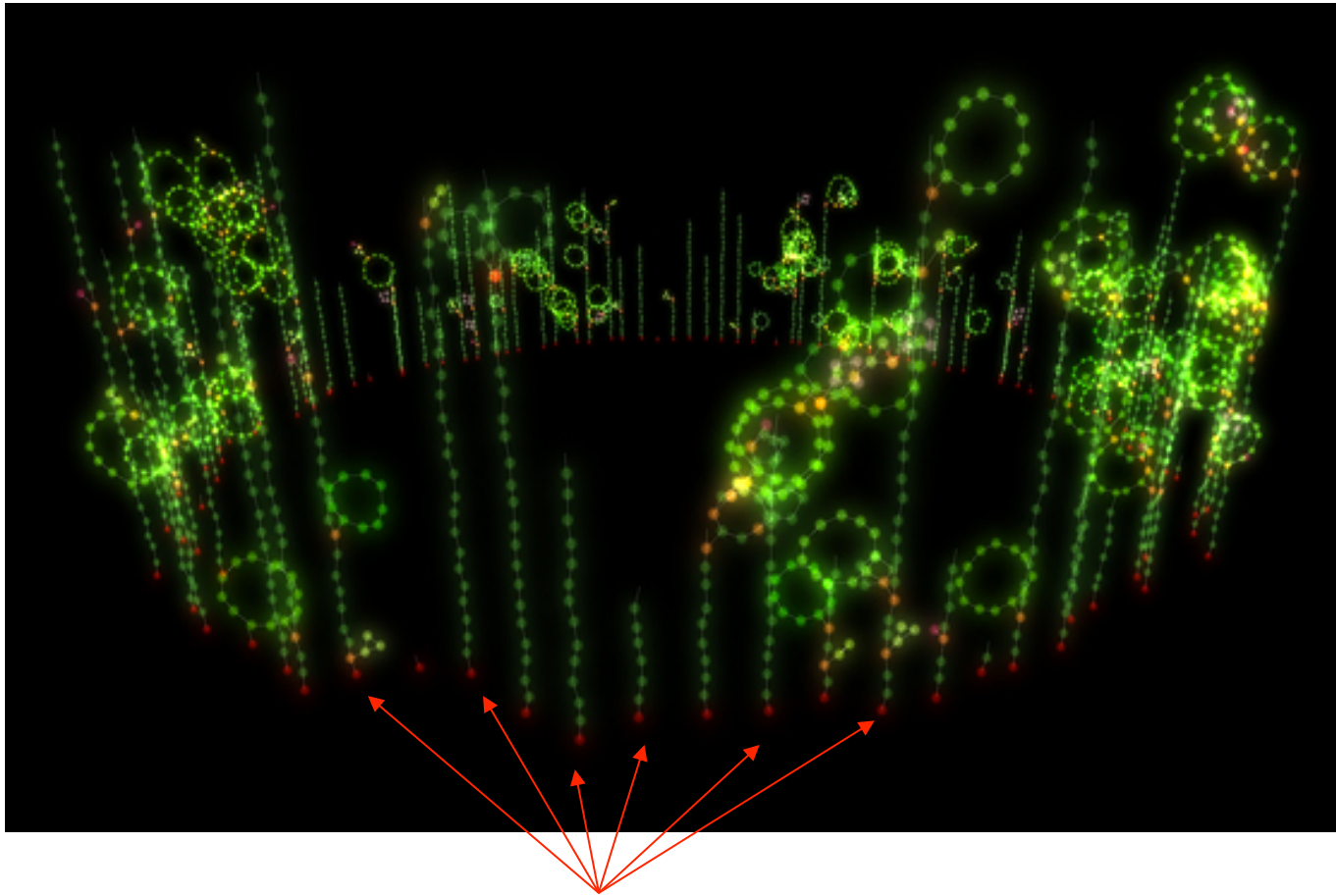
Retail Market as Complex system



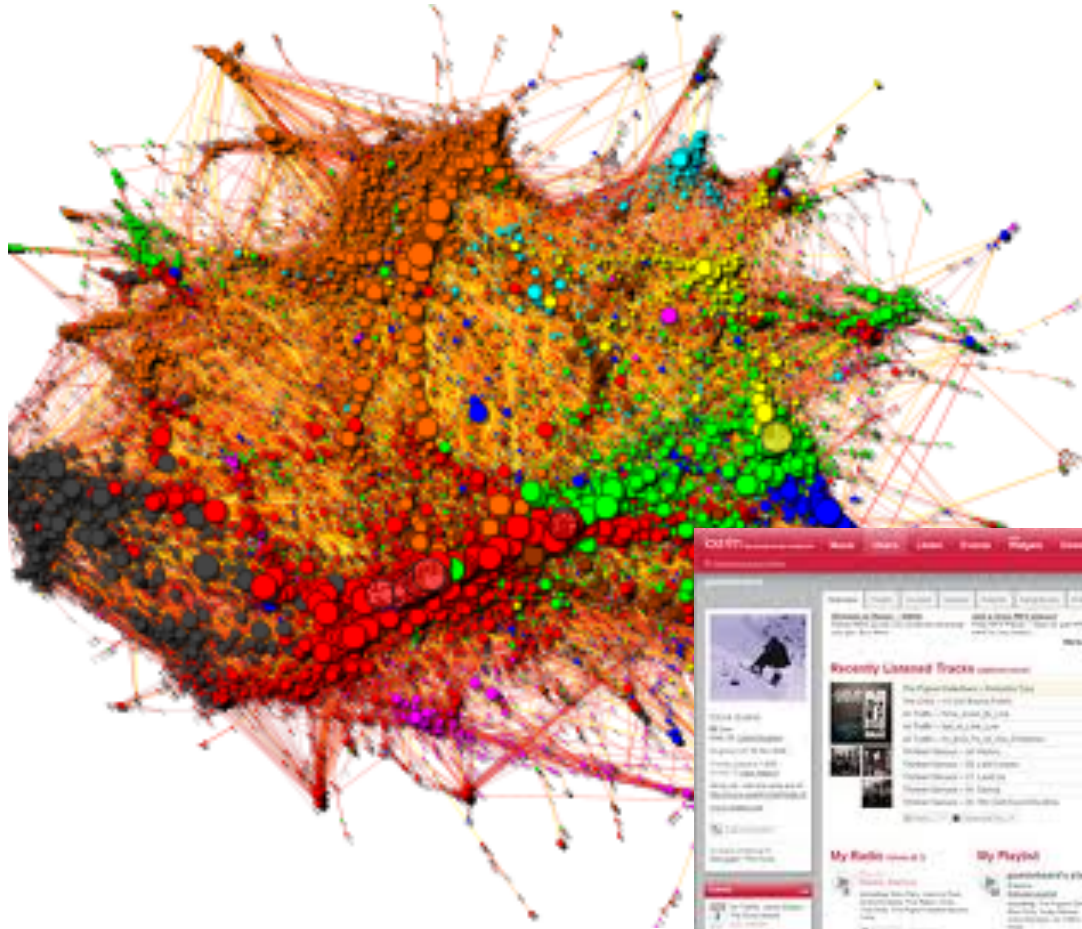
THE PARADOX OF SOCIAL INFLUENCE



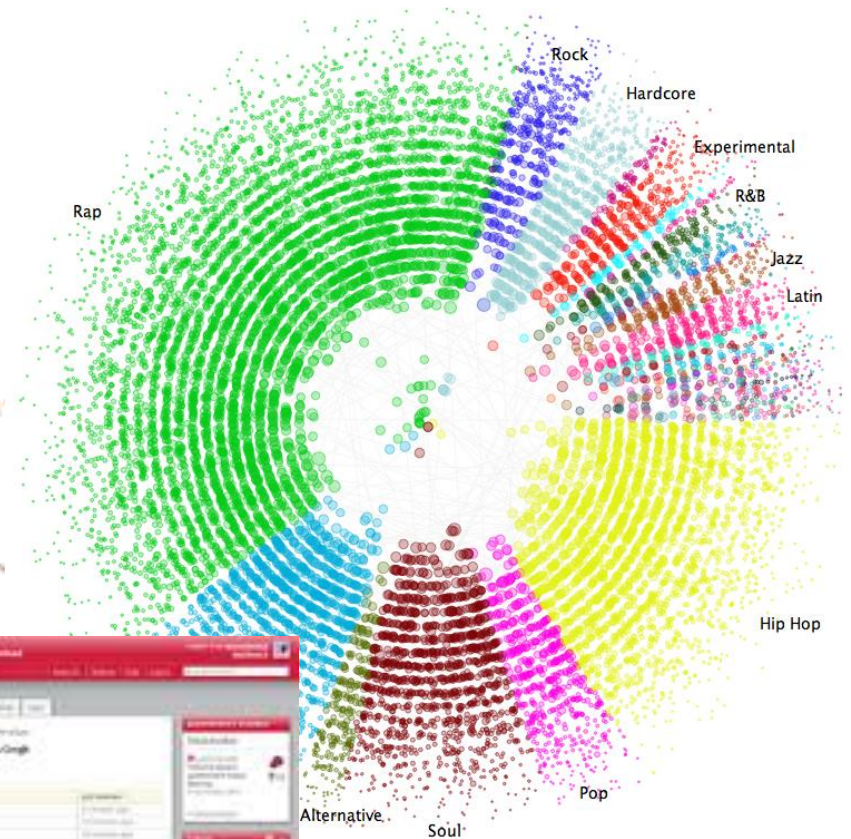
Social Influence: Leaders



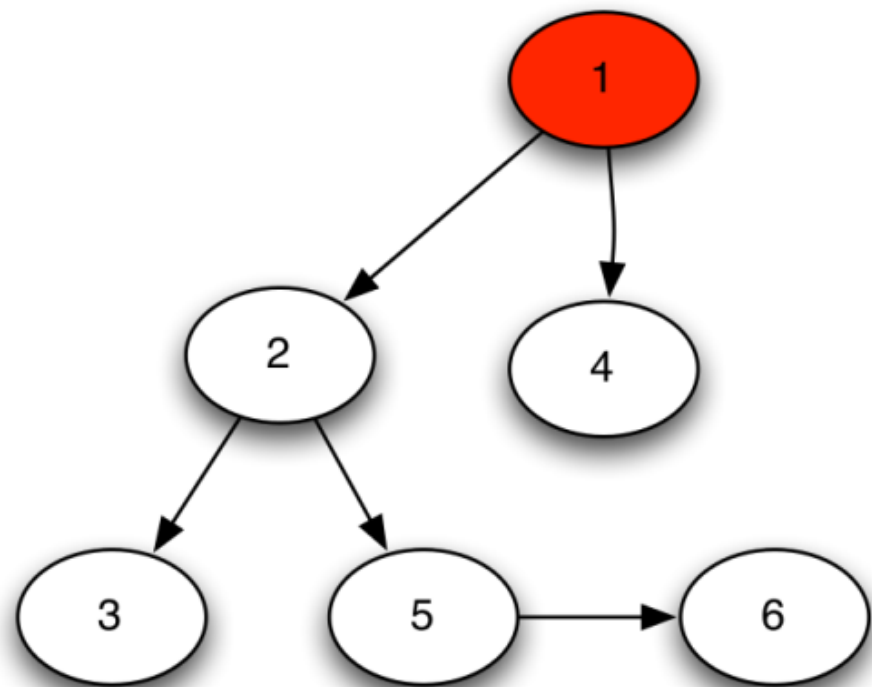
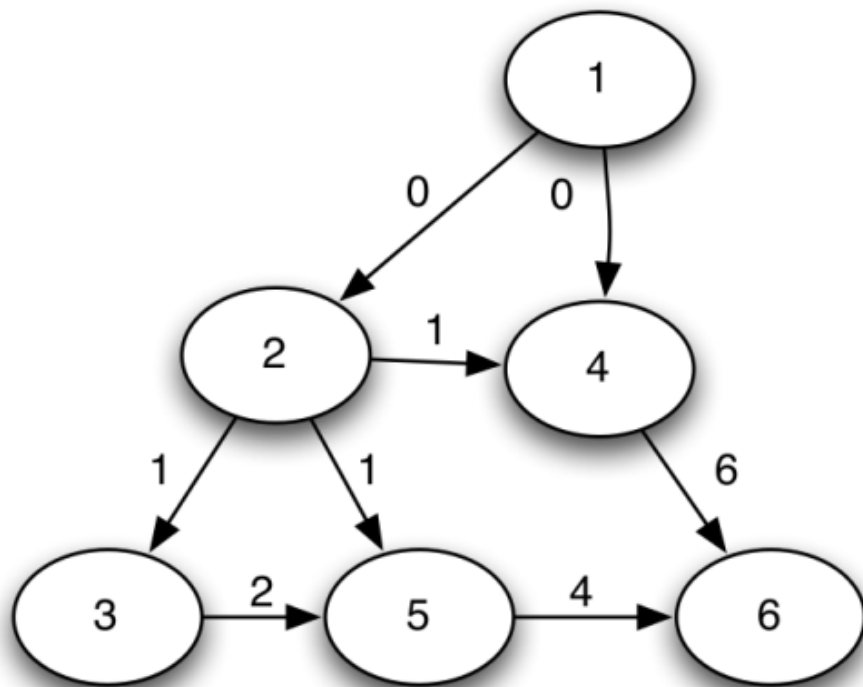
Ask to LAST.FM



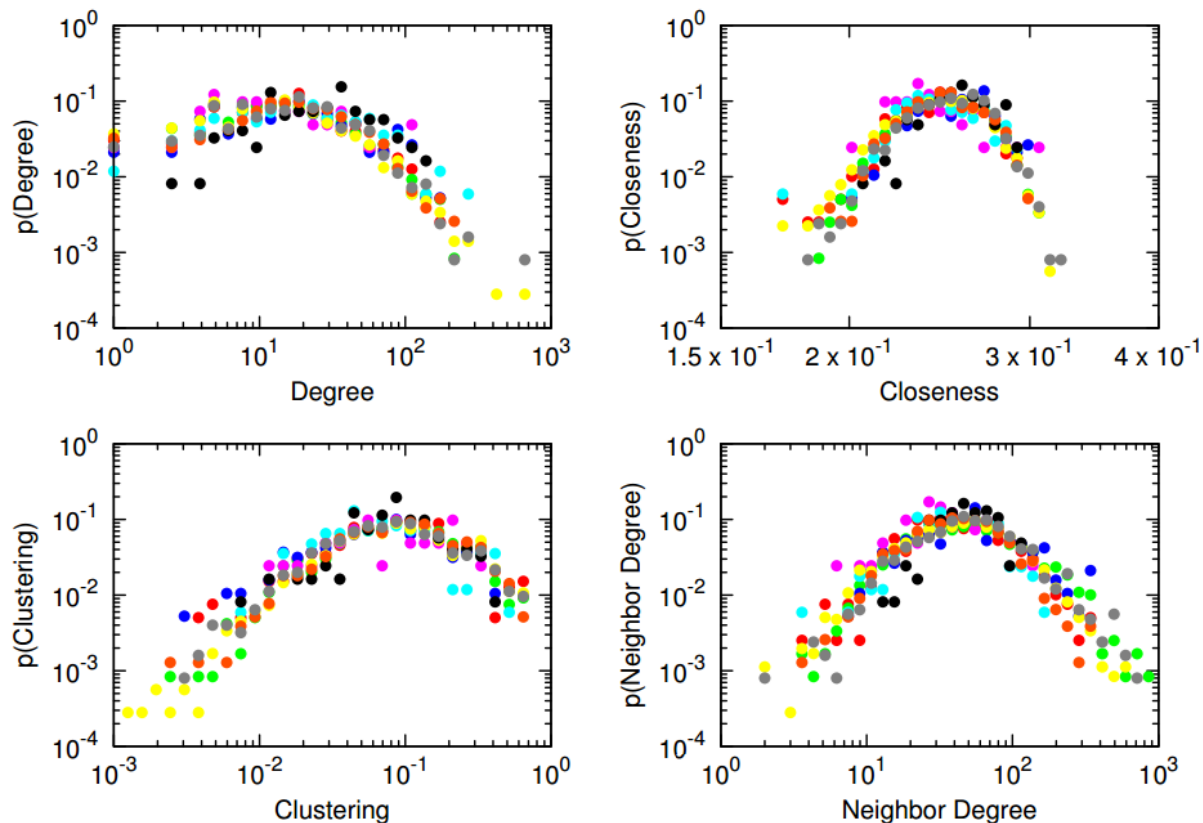
80.000 utenti, 4000.000 connessioni



Leader finding



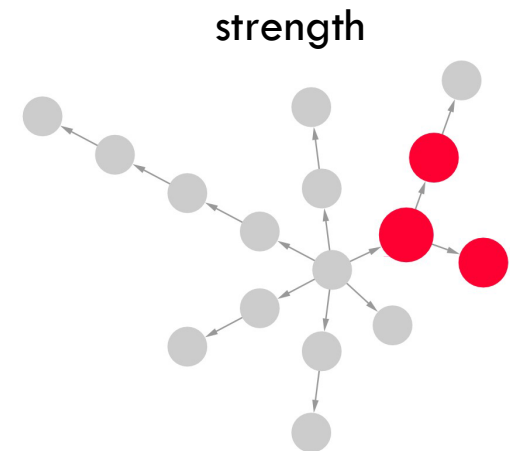
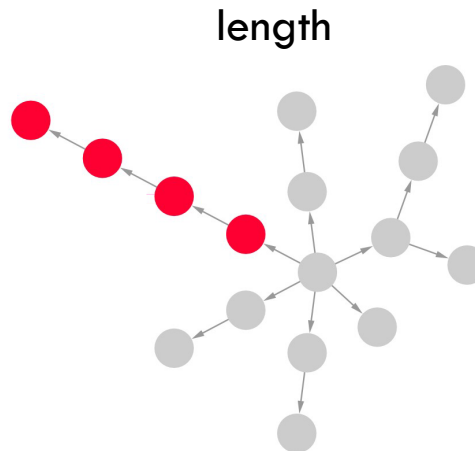
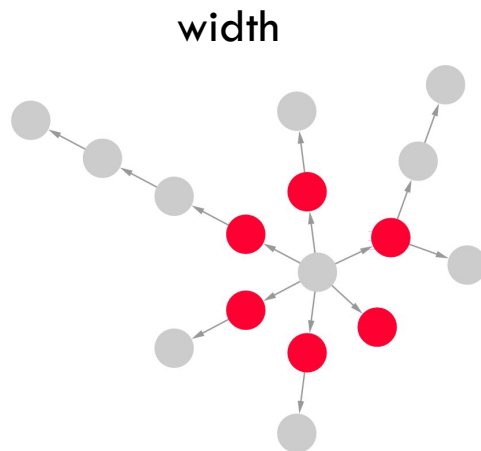
From BigData...true influenzer are not leaders



... abbiamo scoperto che i leader teorici, quelli che avrebbero in teoria il potere di influenzare la rete sociale, non hanno una grande influenza pratica sulla rete.

What is Social Prominence?

- It has been observed that a small set of users in a Social Network is able to anticipate (or influence) the behavior of the entire network
- We detected 3 possible scenarios:





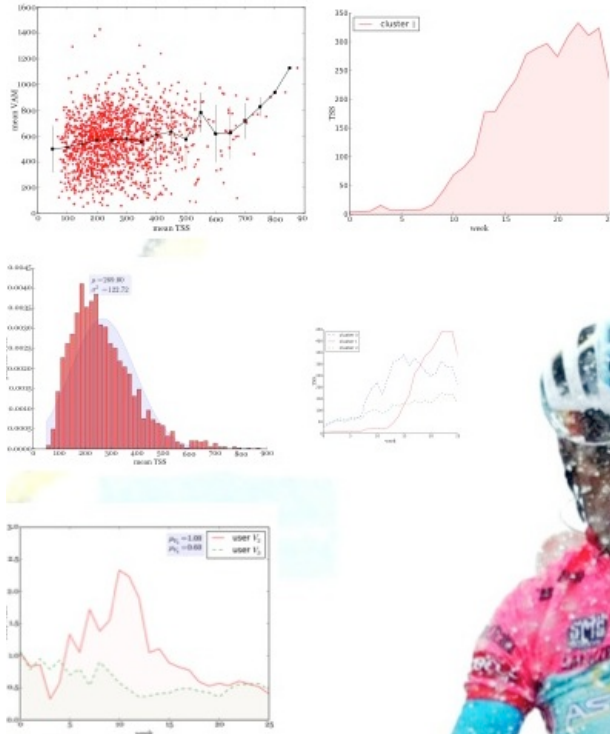
No limits to creativity

If data are available, then any phenomenon becomes measurable, quantifiable and possibly predictable ... including human behaviour

Big Data: the way of **Success**

The patterns of success in
cycling:

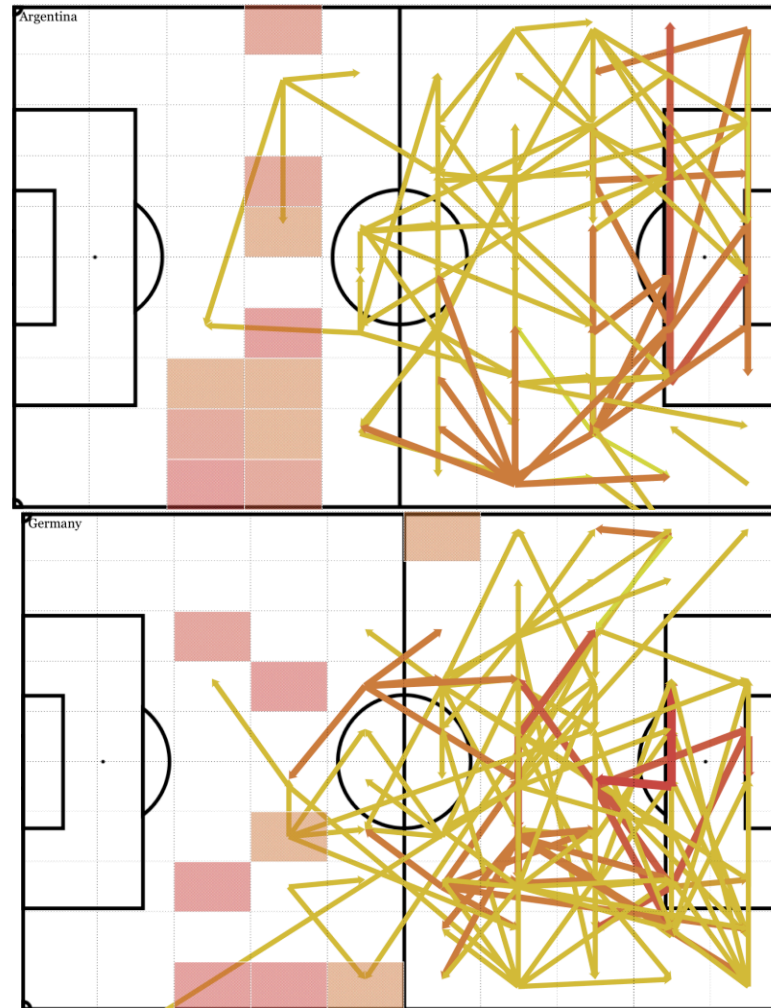
- data from Strava.com
- How you train is fundamental
- A confirmation of the
“overcompensation” theory



The patterns of success in Sports

“Football is a simple game: 22 men chase a ball for 90 minutes and at the end, the Germans always win”

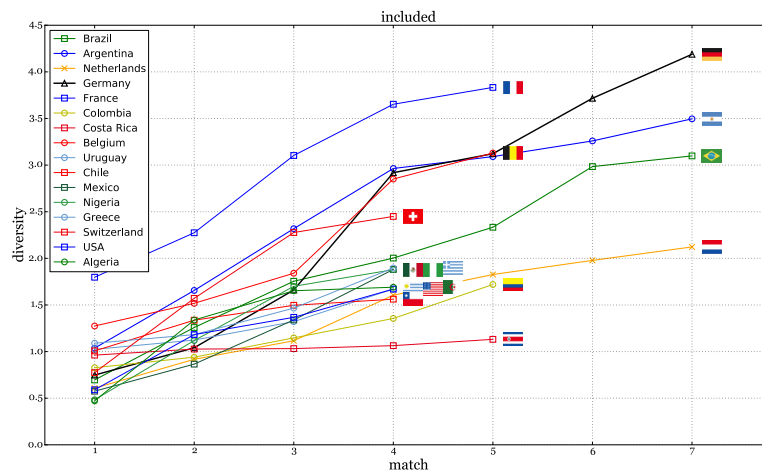
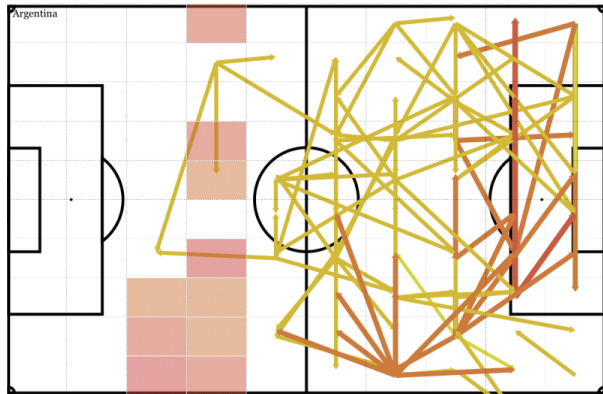
-- Gary Lieneker (after Italy 1990 Final)



Big Data: the way of **Success**

The patterns of success in football:

- detailed data on every match (trajectories, passes, goals, ...)
- a network approach to study the strategy of teams
- a data mining approach to study the performance of players



According to our models the final will be Germany-Argentina. Are our data-driven models correct? Let's see what happens!!! [#WorldCup2014](#)

9:00 PM - 8 Lug 2014 📍 Pisa, Italia

1 RETWEET 2 FAVORITES



A heatmap of a football pitch where color intensity represents player density. Red and orange areas indicate high concentration of players, while green and blue areas indicate lower density. The pitch is divided into three vertical sections by white lines. A semi-transparent dark grey box is centered over the middle section, containing a list of event data.

Data from Opta:

All events during the match

```
...  
<tackle,15.4,41.1,112>  
<pass,25.0,67.1,113>  
<pass,65.0,87.1,115>  
<assist,82.1,35.8,120>  
<goal attempt,82.1,35.8,121>  
.....
```


Big Data Analytics & Social Mining



Who's got the benefits of big data so far?

- A few **latifundists of data**

- **GAFA**

Google



facebook

amazon.com.

- Profiling for behavioral advertising and target marketing

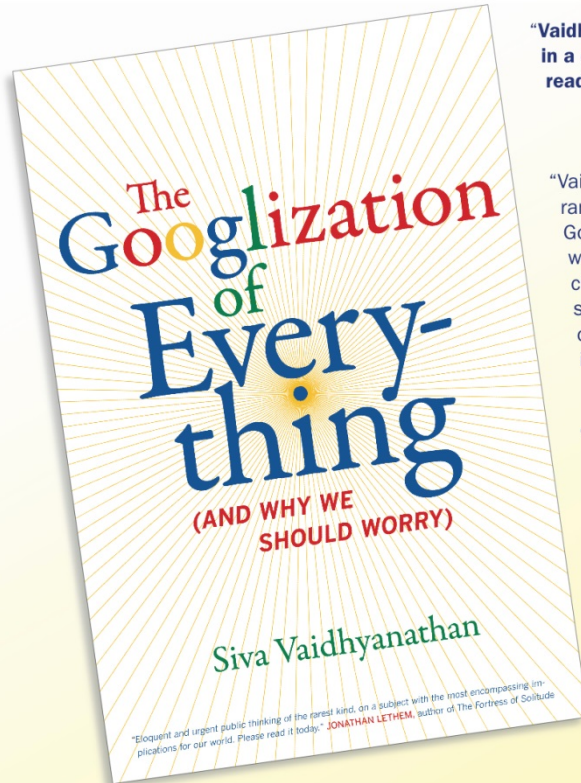
- **NSA**



- Profiling for discovering potential threats to homeland security
- Mass surveillance

**“Finely written and engaging....
A book for anyone who has used Google.”**

—Toby Miller, author of *Makeover Nation*



**“Vaidhyanathan is everything you could want
in a cultural critic: funny, fantastically
readable, and insightful as hell.”**

—Cory Doctorow, author of *For the Win*
and co-editor of *Boing Boing*

**“Vaidhyanathan’s lively, thoughtful, and wide-
ranging book makes clear, in detail, how
Google is reshaping the way we live and
work. He finds much to admire, but also
challenges us to not only use Google’s
services, but to go beyond them to
create a new and genuinely democratic
information order.”**

—Anthony Grafton, author of *Codex in Crisis*

**“Thoughtfully examines the insiders
influence of Google on our society....
As Vaidhyanathan points out, we
must be cautious about embracing
Google’s mission and not accept
uncritically that Google has our
best interests in mind.”**

—Publishers Weekly, Starred Review

**“A critically important book because it’s really about the
Googlization of All of Us.... A brilliant meditation on technology,
information, and consumer inertia, as well as an ambitious
challenge to change how, where, why, and what we Google.”**

—Dahlia Lithwick, senior editor and writer, *Slate Magazine*



At bookstores or www.ucpress.edu/go/googlization



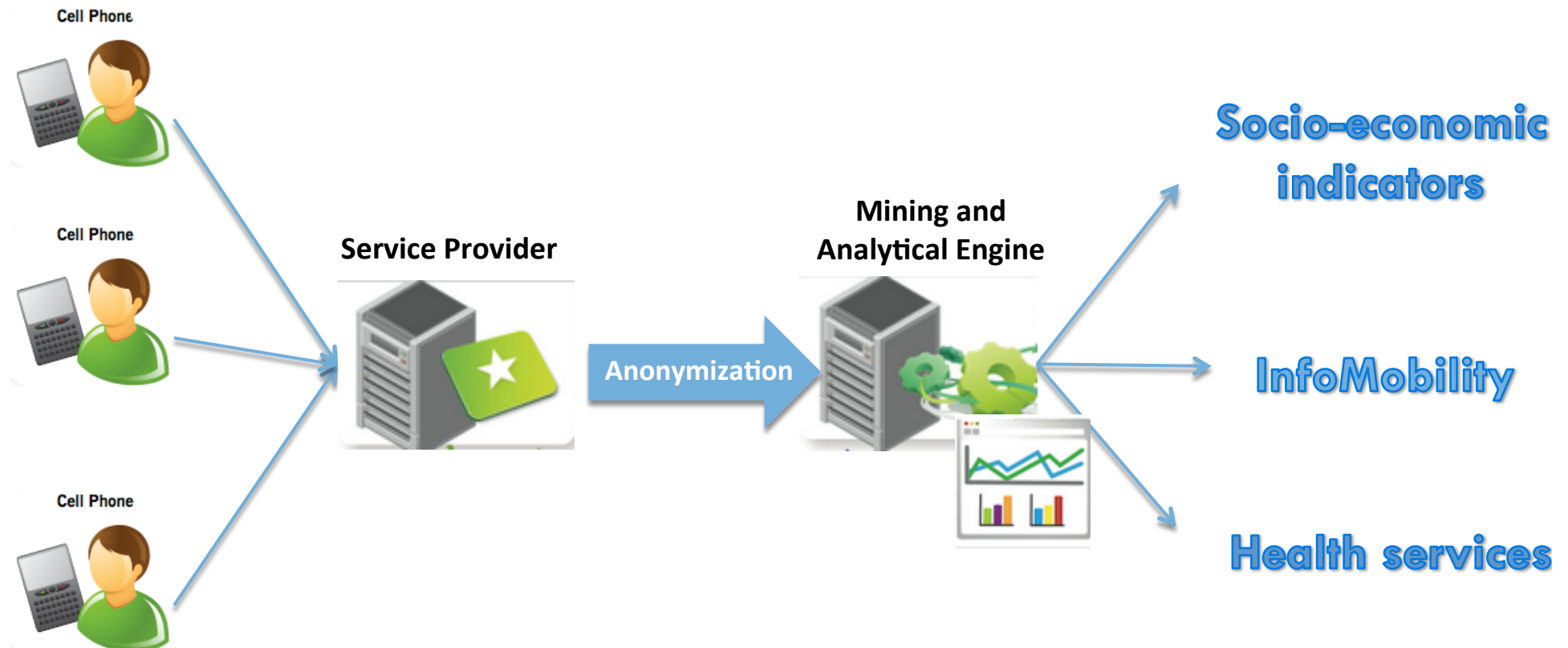
UNIVERSITY OF CALIFORNIA PRESS

We are not Google’s
customers,
we are its products.

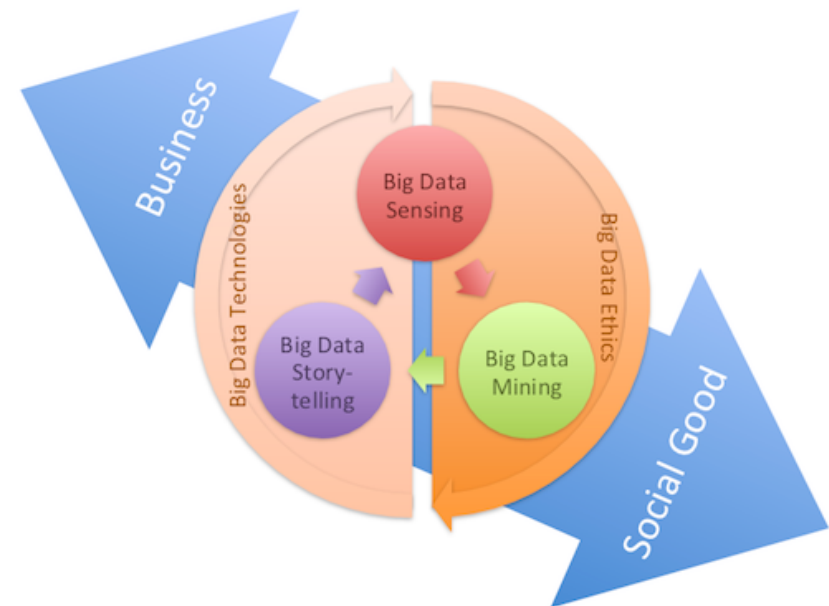
We – our fancies,
fetishes, predilections,
and preferences – are
what Google sells to
advertisers.

Privacy-by-design in big data analytics and social mining

Anna Monreale^{1,2*}, Salvatore Rinzivillo², Francesca Pratesi^{1,2}, Fosca Giannotti² and Dino Pedreschi¹



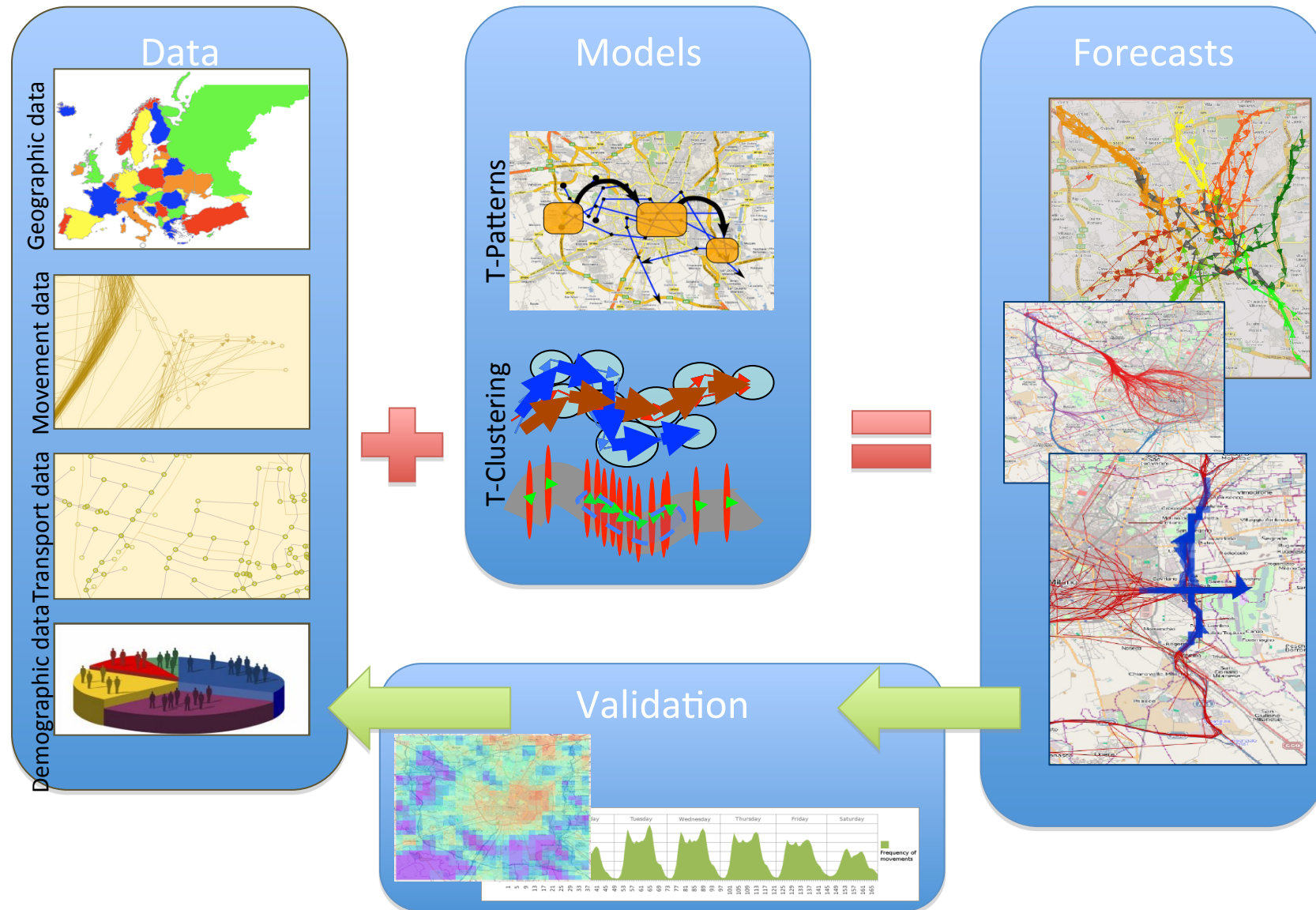
The modern data scientist!!!



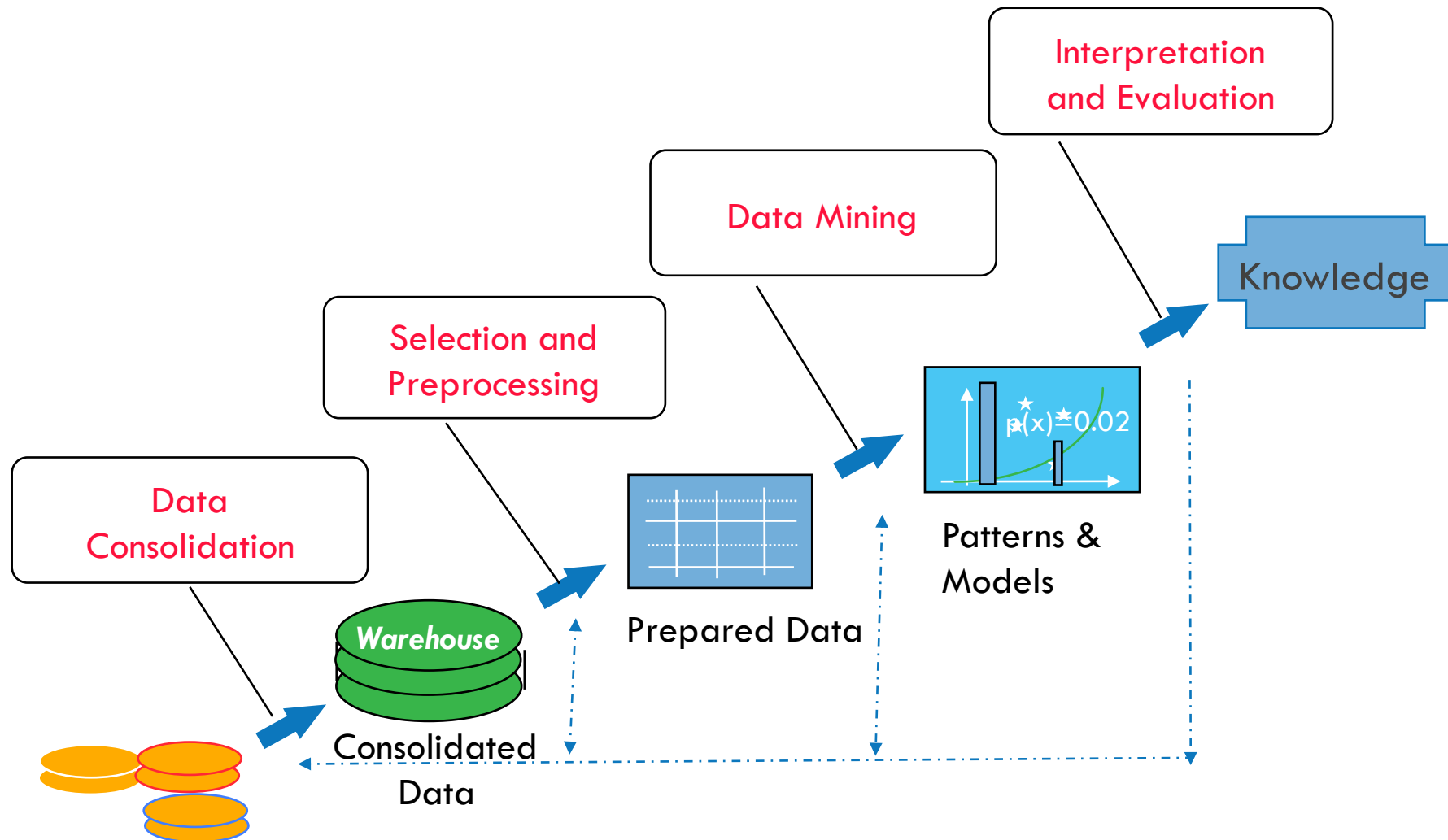


How to develop a big data analytics project

From DATA to KNOWLEDGE

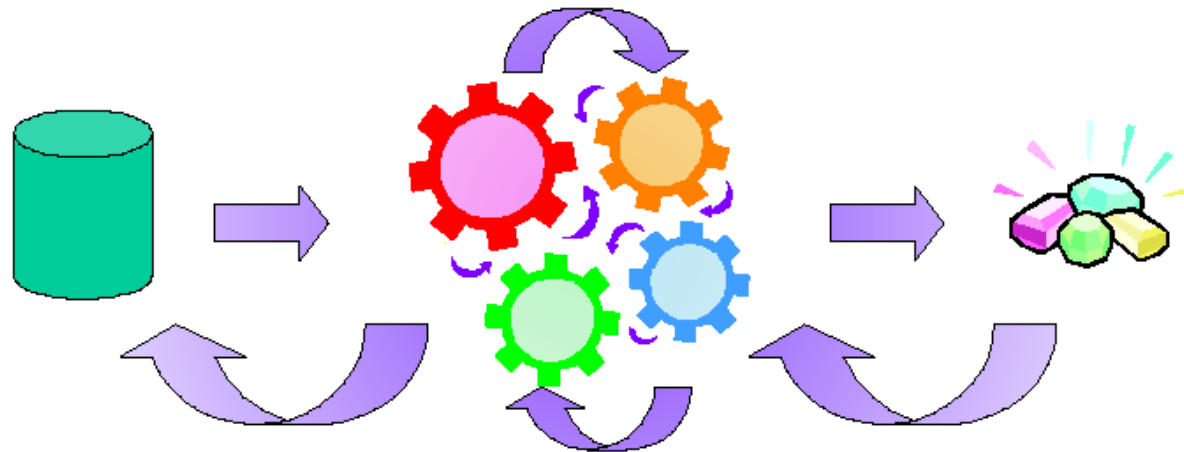


The KDD process

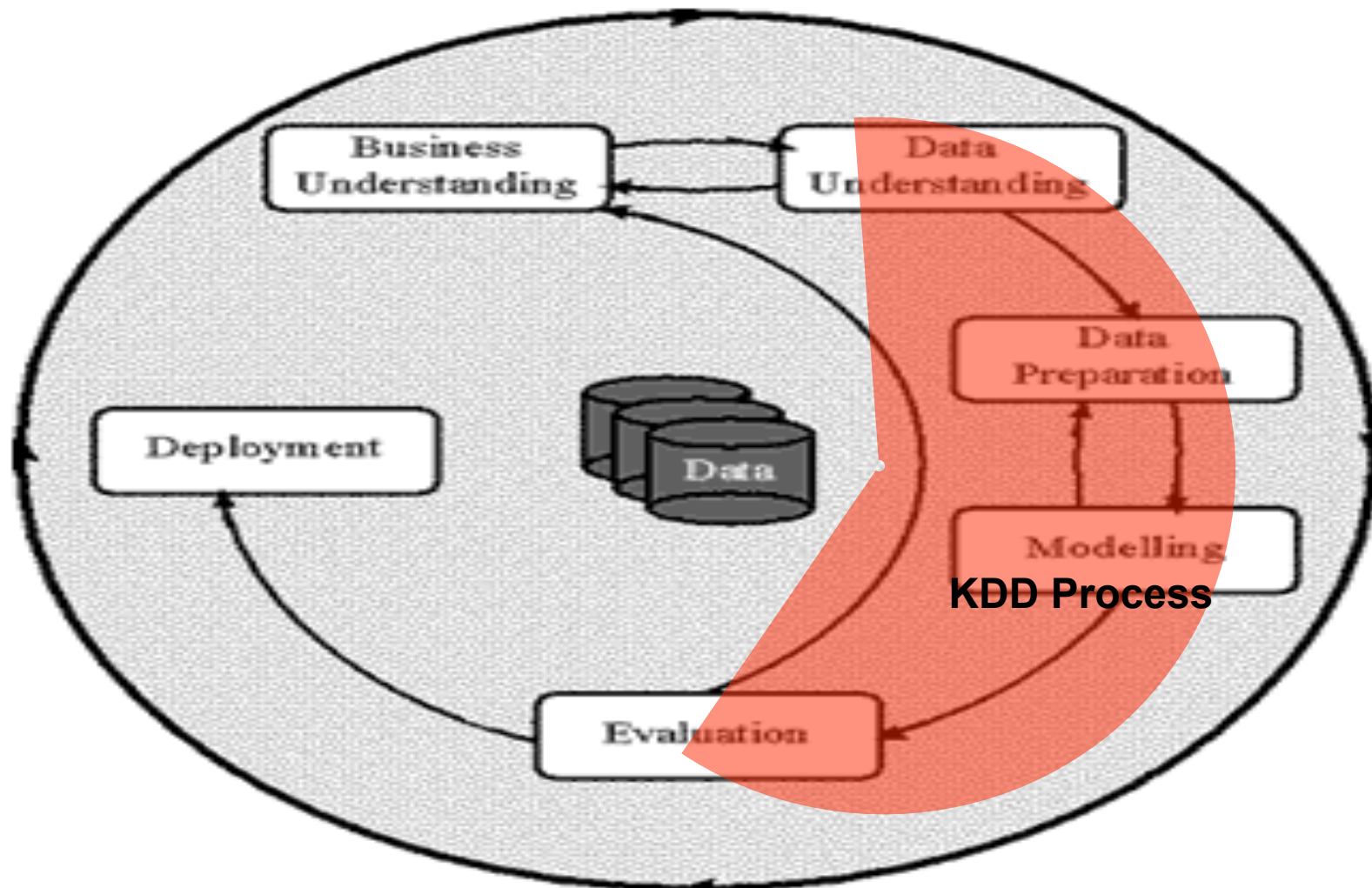


The KDD Process in Practice

- KDD is an Iterative Process
 - ▣ art + engineering rather than science



CRISP-DM: The life cycle of a data mining project





Big Data Number

Social Media Timeline

Big Data

2011

Google+

2010

Instagram

Pinterest

2009

foursquare

2008

2007

bad00

tumblr.

reddit

2006

NETLOG

slideshare

lost.fm

2005

Twitter

2004

YouTube

digg

flickr

facebook

2003

LinkedIn

myspace

2002

friendster

delicious

2001

Wikipedia

SOCIAL MEDIA TIMELINE

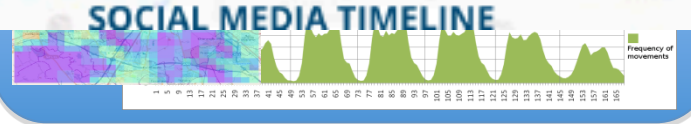
Frequency of movements

Statistics:

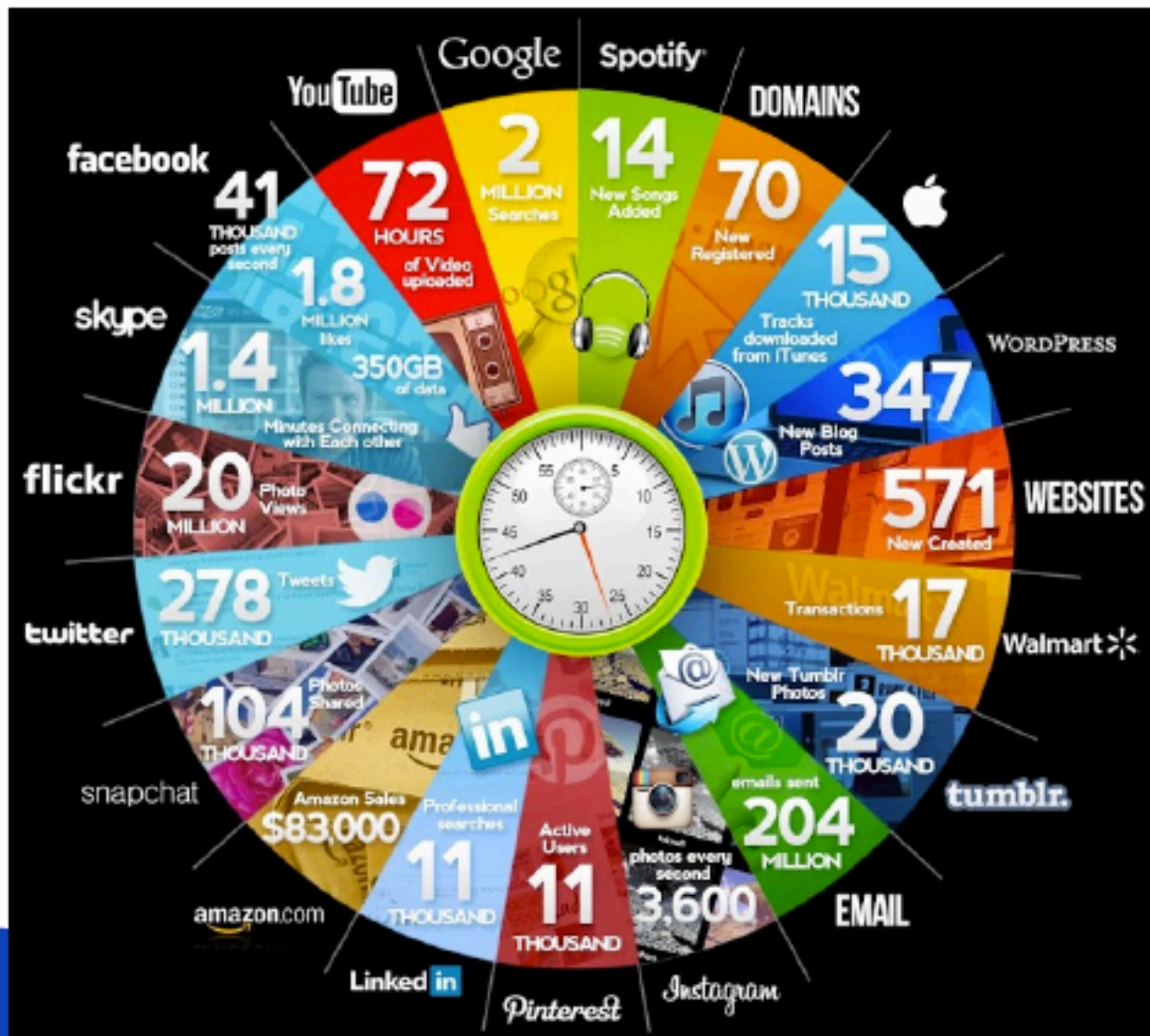
- 271 million active users
- 580 million tweets per day
- 2.1 billion queries per day
- 125 countries
- 35 different languages

Facebook Statistics:

- 1.2 billion of active users
- 150 billion of FB friend connections
- 50 million of FB pages
- 4.5 billion daily FB likes
- 213 countries
- 70 languages



Every minute in Social Media



Data....

- 1,200,000,000,000,000,000,000,000 bytes of data
- Facebook - 1,150 million users
- Gmail – 425 million users
- Skype – 300 million users
- Tweeter – 500 million users (200M active)
- WhatsApp – 300+ million users
- Youtube – 1,000 million users (4 B daily views)
- Instagram - 150 million users

Sources:

<http://expandedramblings.com/index.php/resource-how-many-people-use-the-top-social-media/> September 15, 2013

Data....

- Waze – 50 million users
- Amazon – 209 million users
- Ebay - 120 million users
- Paypal - 132 million users
- Google searches – ~12 billion (monthly, US alone)

Sources:

<http://expandedramblings.com/index.php/resource-how-many-people-use-the-top-social-media/> September 15, 2013

Big Data and Vs

- **Volume and complexity** of data is increasing. “complexity”: it refers to the context of data (creation, provenance, relations) in which it exists and which must be considered when interpreting or re-using the data.
- **Velocity** with which data is being created and characterised is changing
- **Variety** of data in all respects and the challenges of combining variety
- Veracity related to aspects such as trust in dealing with data, i.e. statistical significance.
- Value
- Privacy

With the datafication comes big data, which is often described using the four Vs:

- Volume
- Velocity
- Variety
- Veracity

Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Volume...

...refers to the vast amounts of data generated every second. We are not talking Terabytes but Zettabytes or Brontobytes. If we take all the data generated in the world between the beginning of time and 2008, the same amount of data will soon be generated every minute. New big data tools use distributed systems so that we can store and analyse data across databases that are dotted around anywhere in the world.

Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Velocity...

...refers to the speed at which new data is generated and the speed at which data moves around. Just think of social media messages going viral in seconds. Technology allows us now to analyse the data while it is being generated (sometimes referred to as in-memory analytics), without ever putting it into databases.

Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Variety...

refers to the different types of data we can now use. In the past we only focused on structured data that neatly fitted into tables or relational databases, such as financial data. In fact, 80% of the world's data is unstructured (text, images, video, voice, etc.) With big data technology we can now analyse and bring together data of different types such as messages, social media conversations, photos, sensor data, video or voice recordings.

Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Veracity...

...refers to the messiness or trustworthiness of the data. With many forms of big data quality and accuracy are less controllable (just think of Twitter posts with hash tags, abbreviations, typos and colloquial speech as well as the reliability and accuracy of content) but technology now allows us to work with this type of data.

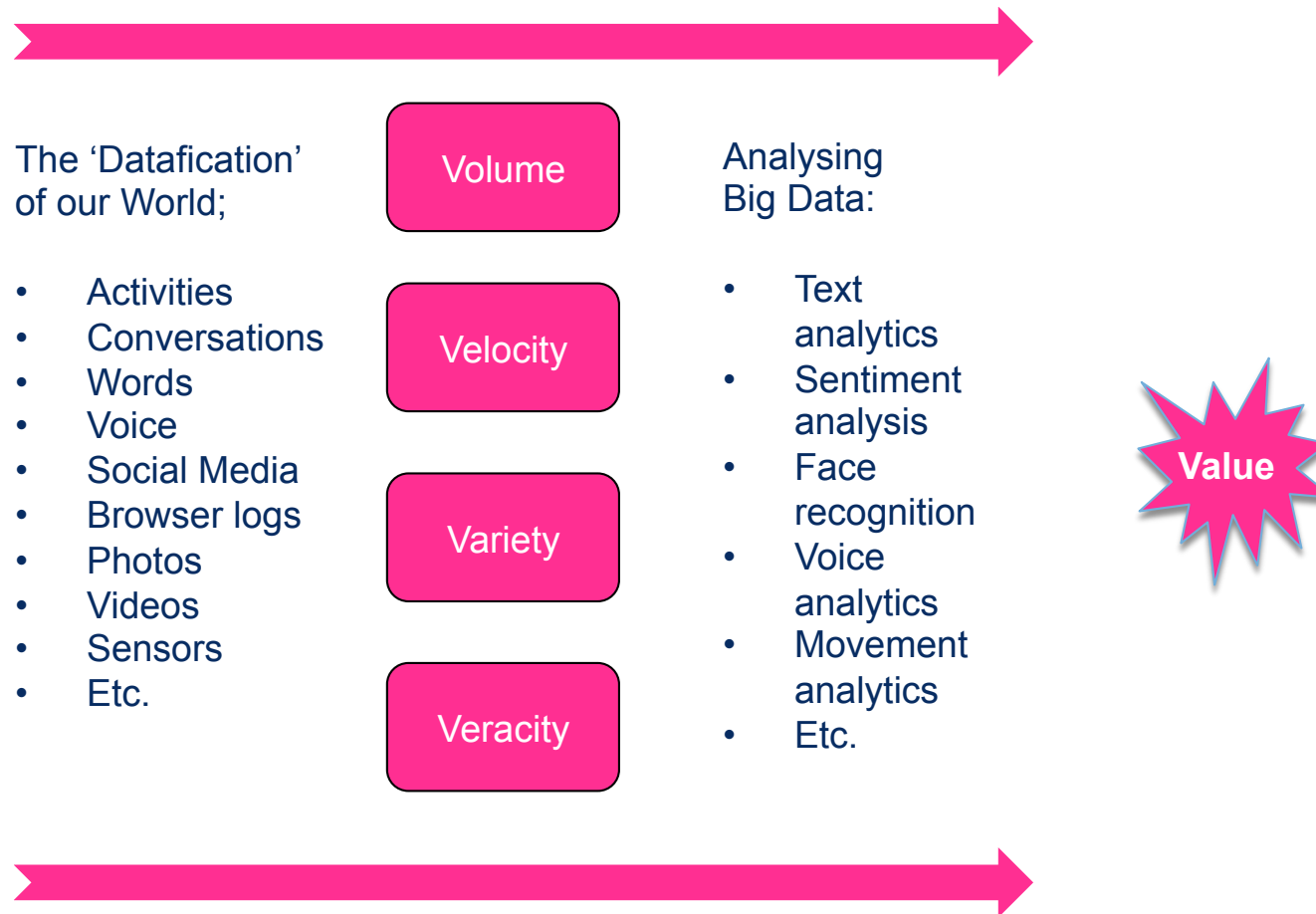
Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Turning Big Data into Value:

The datafication of our world gives us unprecedented amounts of data in terms of Volume, Velocity, Variety and Veracity. The latest technology such as cloud computing and distributed systems together with the latest software and analysis approaches allow us to leverage all types of data to gain insights and add value.

Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions

Turning Big Data into Value:



Bernad Marr Bigdata: using Smart BigData analytics and metrics
To make better decisions



First exercise: look for open datasets

Esempi – dati governativi

- [Data.gov](#)
- [US Census Bureau](#)
- [Open Data Portal dell'Unione Europea](#)
- [Data.gov.uk](#)
- [the CIA World Factbook](#)
- [Healthdata.gov](#)
- [Dati.gov.it](#)

About DwB

[Work Packages](#)[Deliverables](#)[Participants](#)

Metadata Services

[CIMES for European National
Official Statistics Microdata](#)[MISSY for Integrated European
Official Statistics Microdata](#)

Access Services

[Legal Frameworks](#)[National Accreditation & Access](#)[Transnational Access](#)[Guides](#)[Routines for Integrated
Microdata](#)[Synthetic data tools](#)

Activities and Events

[European Data Access Forums](#)[Users' Conferences](#)[Training Events and Material](#)[Staff Visits to Research Data](#)

Data without Boundaries – DwB

The Data without Boundaries – DwB – project came to a formal end on 30 April 2015. The project had a mission to support equal and easy access to the rich resources of official microdata for the European Research Area, within a structured framework where responsibilities and liability would be equally shared. During its four-year lifespan the DwB worked towards preparing a comprehensive European service with better and friendly metadata, a more harmonized transnational accreditation and a secure infrastructure that would allow transnational access to the highly detailed and confidential microdata, both national and European, so that the European Union would be able to continuously produce cutting-edge research and reliable policy evaluations.

The resulting output is presented on this website as well as the tools and services that are maintained and developed by project partners beyond the project end. Please, use the navigation on the left to learn more. The main event to display and discuss the project work was the Second European Data Access Forum (EDAF) held in March 2015 in Luxembourg. All the presentation slides and audio files are available from the EDAF event page.

Bridging Three Communities



The DwB Formal End

The project came to a formal end on 30 April 2015. [\[More...\]](#)

NEW Deliverable Available!

D5.5 (Final report & recommendations for the continuation of services for European OS microdata) is now publicly available online [\[More...\]](#)

New Service: Visualisation Tool

Legal frameworks for official statistics microdata access demonstrated by Visualisation Tool. [\[More...\]](#)

Microdata Information System MISSY

Systematically structured metadata for official statistics. [\[More...\]](#)

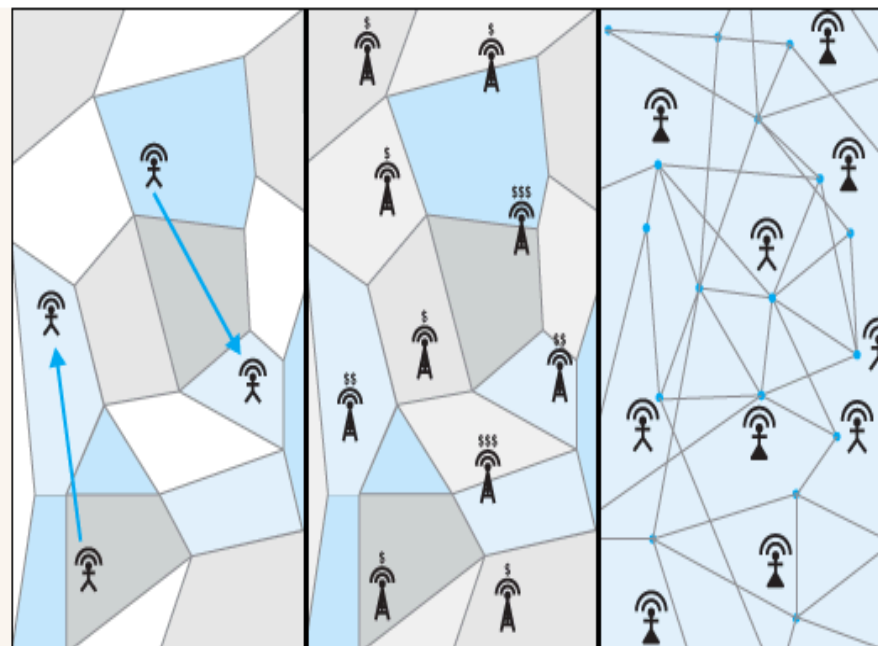
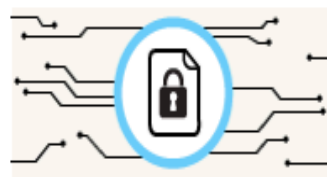
New Service: CIMES

Overview of metadata of official microdata for research purposes. [\[More...\]](#)

[HOME](#)[RESEARCH](#)[LABS](#)[BLOG](#)[MULTIMEDIA](#)[ABOUT](#)[CONTACT](#)

MOBILE PHONE NETWORK DATA FOR DEVELOPMENT

A synthesis of a growing body of research on mobile phone data analysis in development or humanitarian contexts. [Read More /](#)

**SUBSCRIBE TO OUR NEWSLETTER****GO****RESEARCH PROJECTS****PRIVACY PRINCIPLES****ROADMAP****PULSE LABS**

NEWS

[Thoughts on the Google Flu Trends debate](#)

TWITTER



4 Apr

A WORLD THAT COUNTS

MOBILISING THE DATA REVOLUTION FOR SUSTAINABLE DEVELOPMENT

