# Big Data

**Paolo Ferragina & Dino Pedreschi**

Dipartimento di Informatica

Università di Pisa

KDD LAB   http://kdd.isti.cnr.it

ACUBE LAB   http://acube.di.unipi.it

# Siamo tutti pollicini digitali

- Plenty of digital breadcrumbs behind us
- La Vita Nova, e-magazine de Il Sole 24 Ore
- Fosca Giannotti, Dino Pedreschi
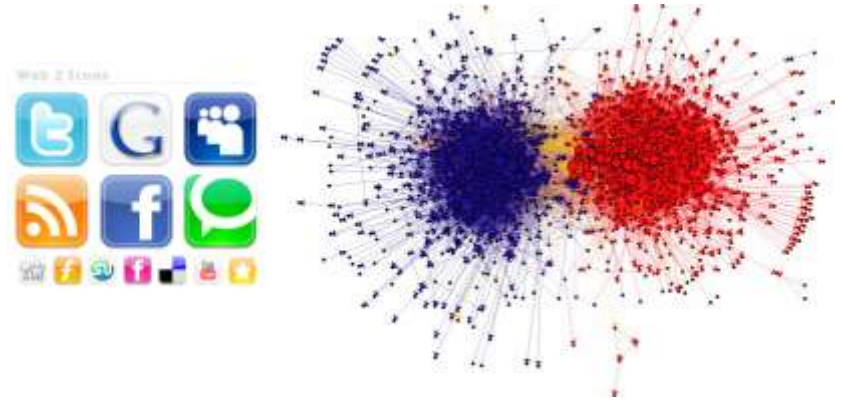- Dicembre 2012
- Everyone is becoming a «statistical entity»

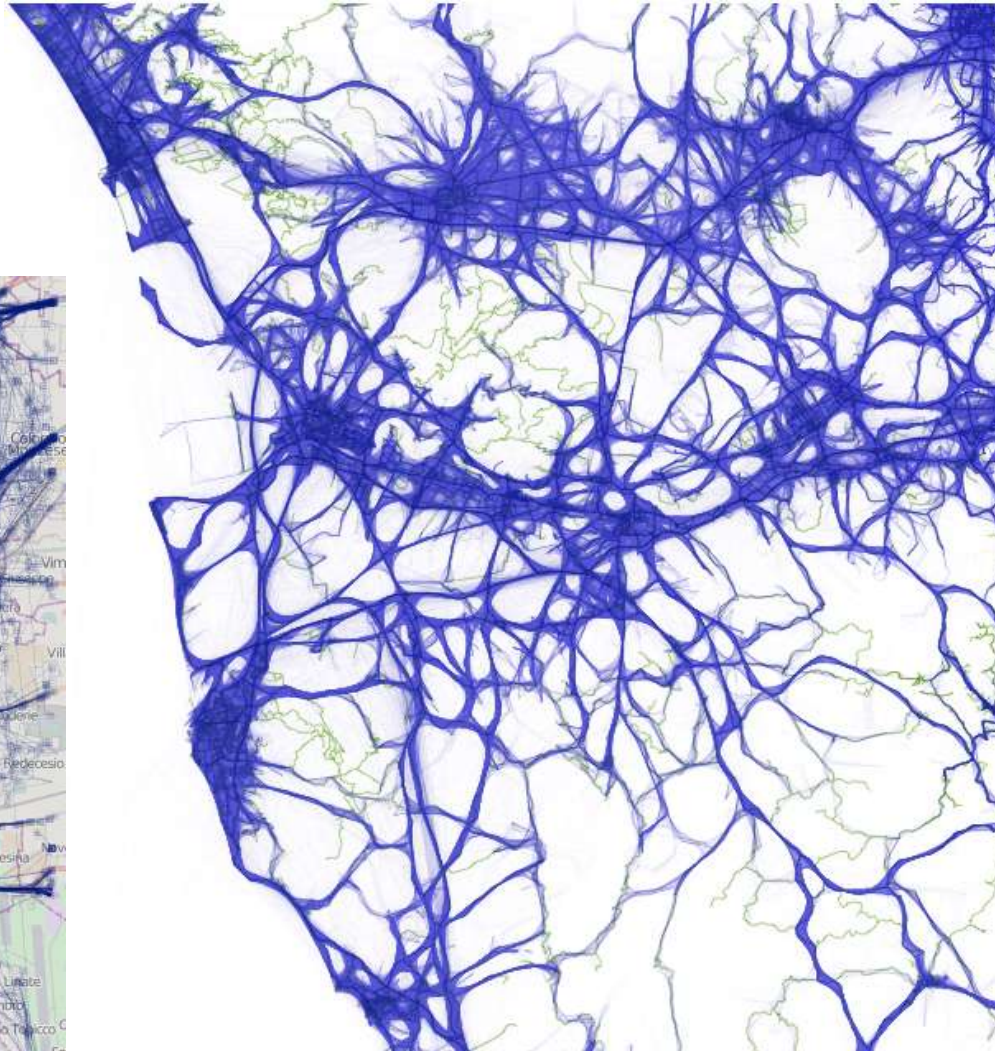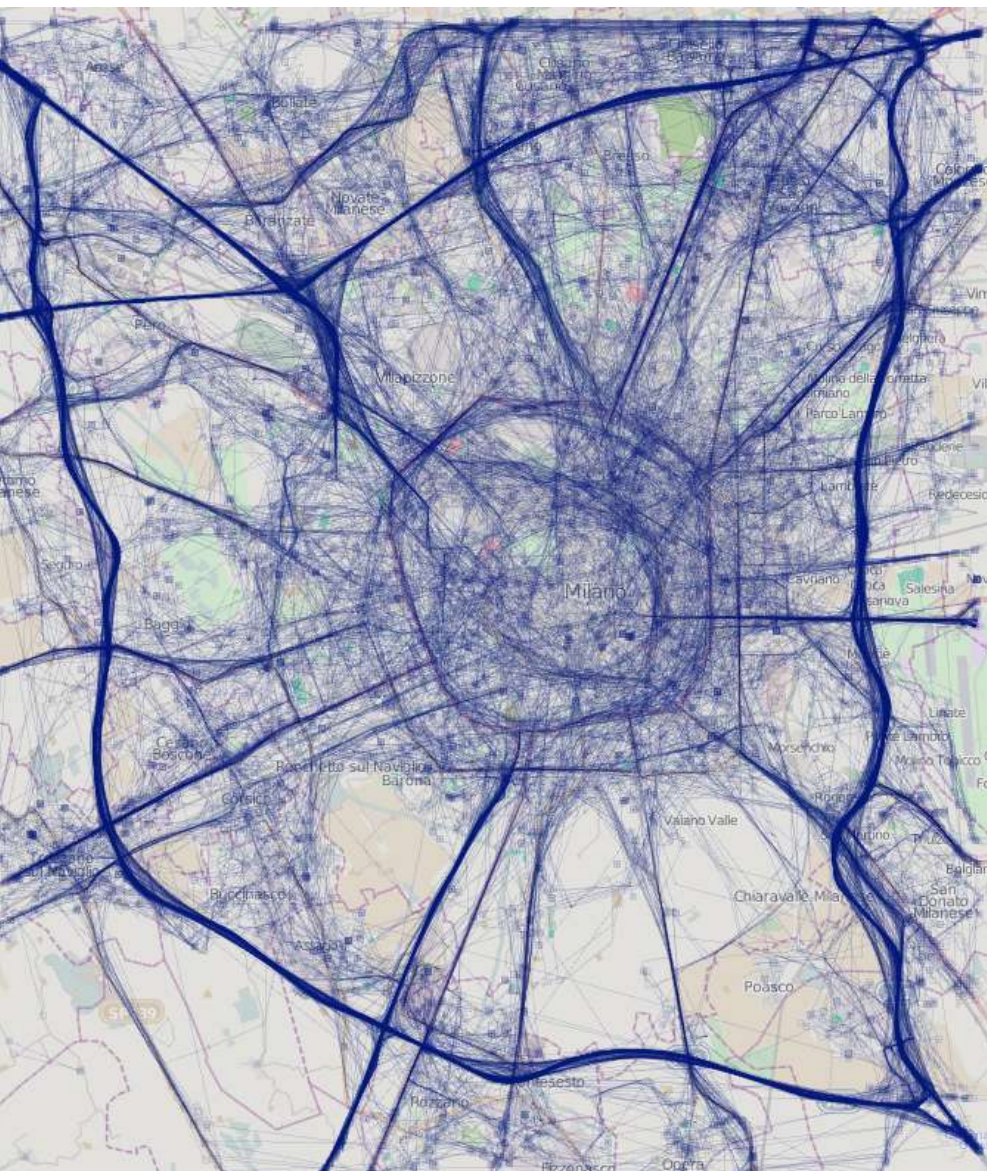# Big data "proxies" of social life
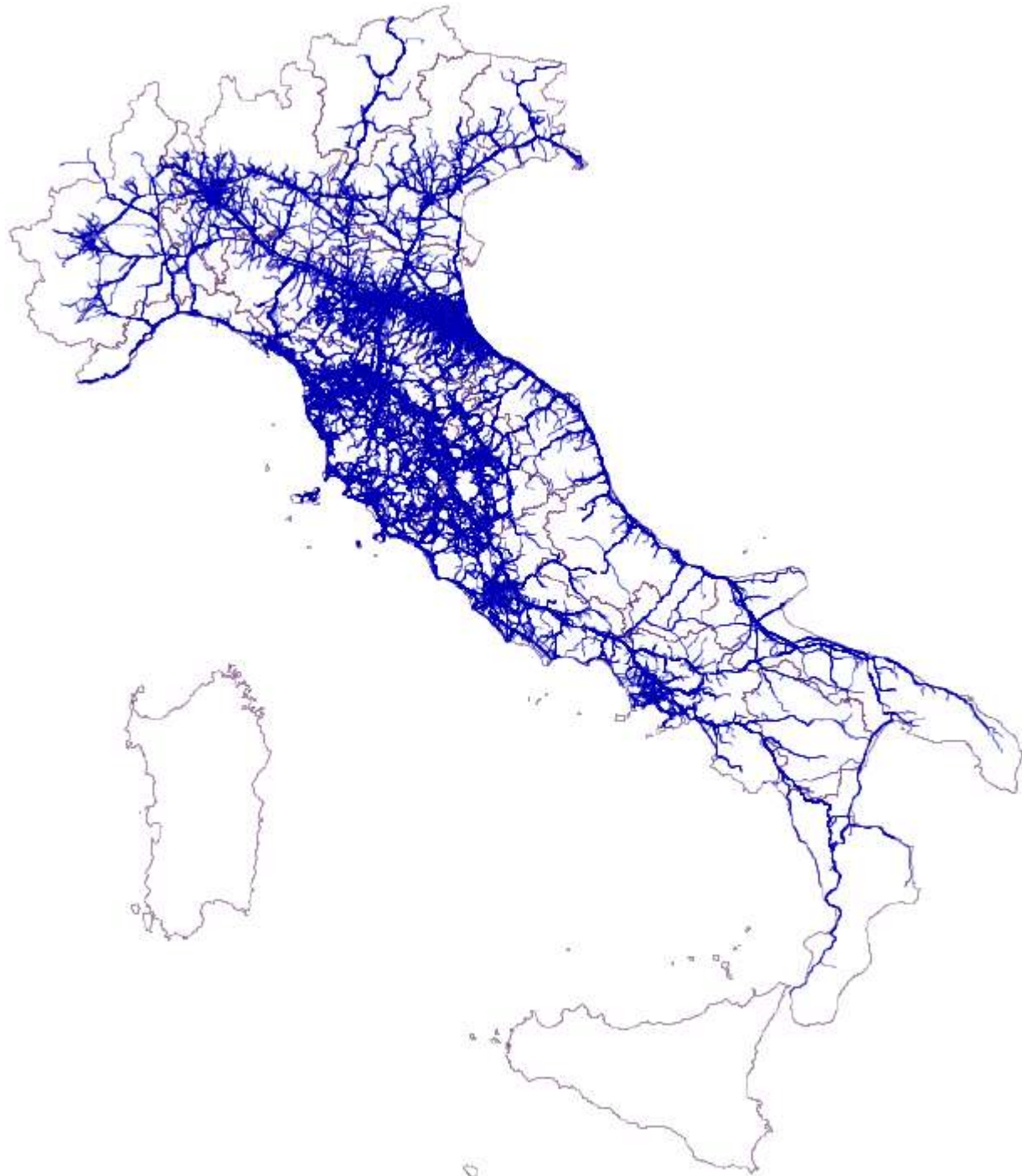
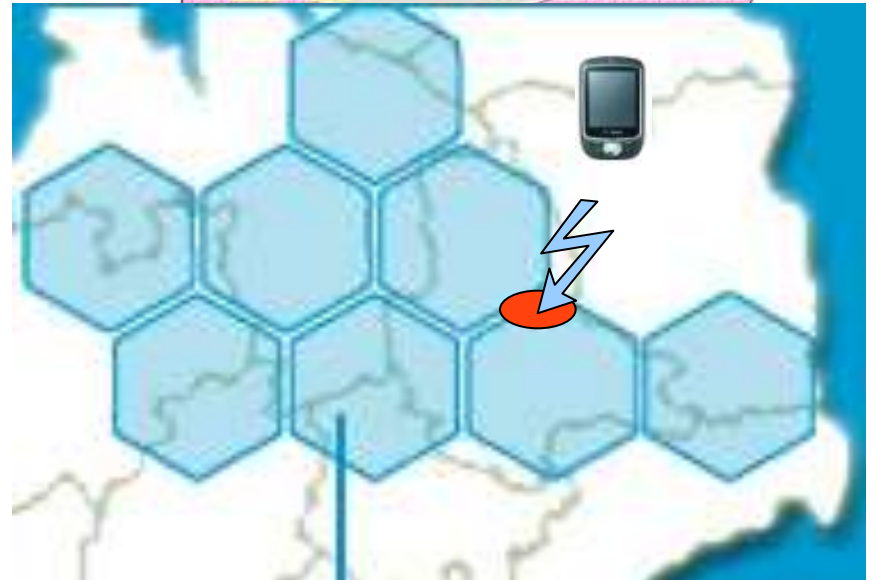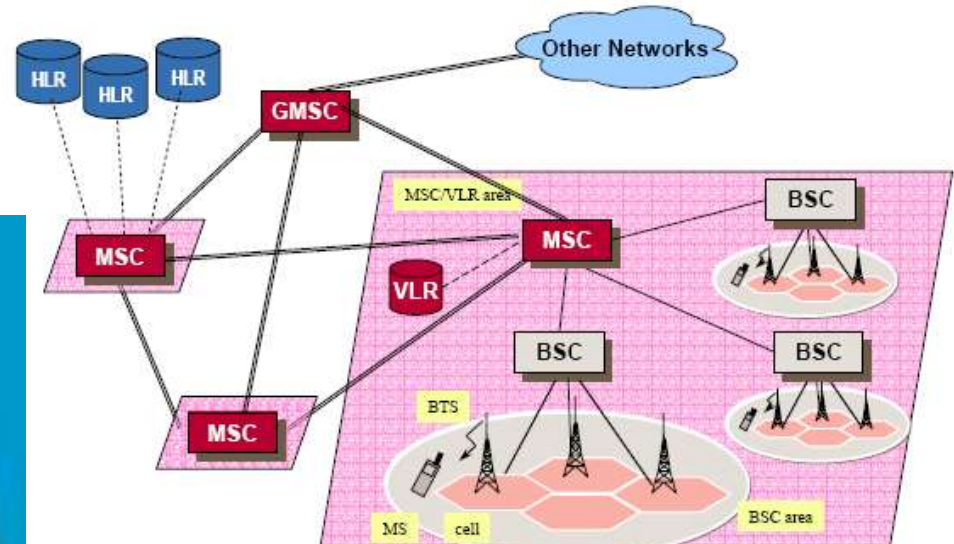Shopping patterns & lifestyle

Relationships & social ties

Movements

Desires, opinions, sentiments

# GSM data: CDR – call detail records

**Understanding human mobility
@ KDD LAB (Unipi + ISTI-CNR)**

*Cascina*

# Mobility atlas of many cities

**Dataset Temporal Distribution**

**Incoming** — **Inner** — **Outgoing**

**Rog Distribution**

**Radius of Gyration**
- 20 km
- 40 km
- 61 km
- 81 km
- 102 km
- 122 km
- 143 km

**Standard Deviation**
- 8 km
- 17 km
- 26 km
- 35 km
- 44 km

1403 m

**Time at Home**

**Time at Work**

# Incoming Traffic (38.464 Trajectories)



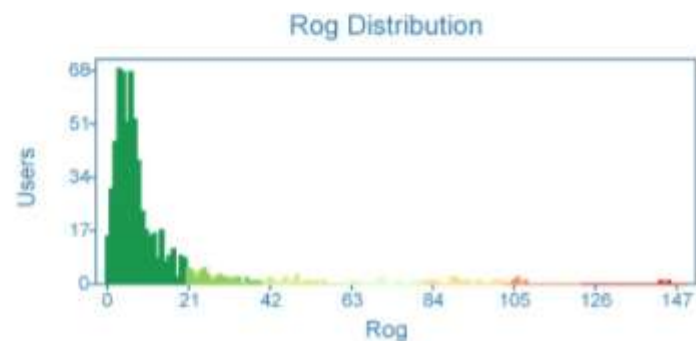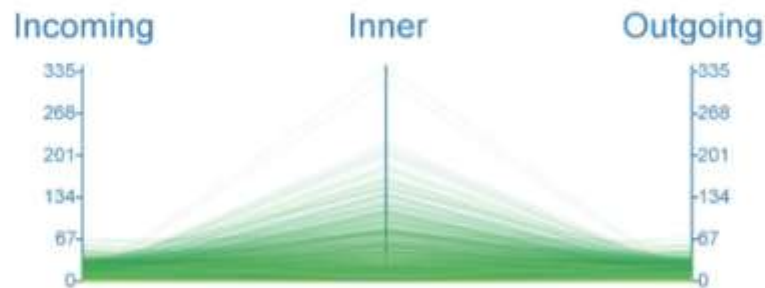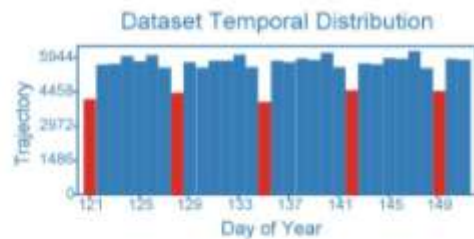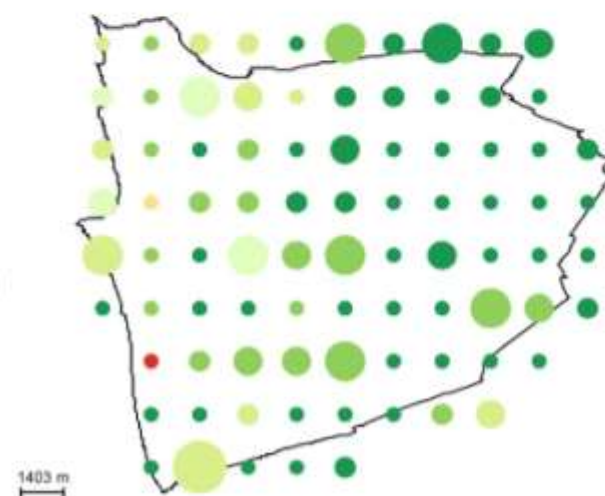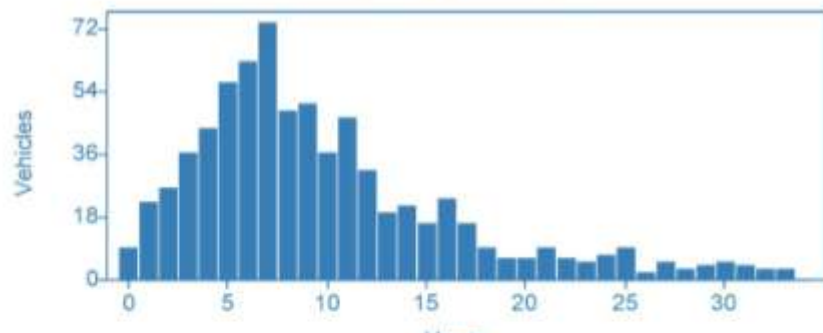| | City | Traj | Perc |
|---|---|---|---|
| NORD 32% | San Giuliano T.. | 4.816 | 62% |
| | Vecchiano | 1.425 | 94% |
| | Viareggio | 1.142 | 99% |
| | Lucca | 882 | 87% |
| | Camaiore | 358 | 94% |
| OVEST 0% | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| SUD 12% | Livorno | 2.843 | 92% |
| | Collesalvetti | 565 | 50% |
| | Rosignano Mari.. | 140 | 41% |
| | Fauglia | 137 | 19% |
| | Cecina | 124 | 45% |
| EST 54% | Cascina | 7.078 | 97% |
| | San Giuliano T.. | 2.881 | 37% |
| | Pontedera | 1.350 | 95% |
| | Calci | 795 | 79% |
| | Calcinaia | 693 | 92% |

## Incoming Temporal Matrix



## Regular VS Occasional



Regular
Occasional

# Outgoing Traffic (38.271 Trajectories)



| City | Traj | Perc |
|---|---|---|
| San Giuliano T.. | 4.842 | 62% |
| Vecchiano | 1.418 | 93% |

# Discover the borders of mobility

# Estimate O-D matrix from phone data



Figure 12: Mobile phone movements in Ivory Coast and Abidjan.

Mirco Nanni, Roberto Trasarti, et al.:
MP4-A Project: Mobility Planning for Africa. "Data for Development" Orange challenge, 2013

**proactive car pooling**

# Car Pooling



Trajectory matching



Carpooling Network



Carpooling Communities

# Carpooling potential

**PISA**

16467 users, 357137 trips

20% are systematic trips

40% of them are matching trips

142.740.060 saved Kms

**FIRENZE**

34864 users, 1040872 trips

9% are systematic trips

58% of them are matching trips

220.802.240 saved Kms



Pisa: no-carpooler 44, only-driver 22, driver-passenger 19, only-passenger 15

Firenze: driver-passenger 31, no-carpooler 28, only-passenger 26, only-driver 14

Legend:
- no-carpooler
- only-passenger
- only-driver
- driver-passenger

**Electrifiability**

# Electrifiability



In Tuscany **75% users** have a daily mobility covered at 100% by an electrical car (home charging only)

In Pisa 90.5% of daily trips are electrifiable:
562.061 km electrificable

# City users Sociometer

# Mobile phone socio-meter

**Analyze individual call habits to recognize profiles**

–Resident

–Commuters

–Visitors/Tourists

# City user profile **quantification**



- 🔴 Resident profile
- 🔵 Commuter profile
- 🟡 Visitor profile

**Classification outcome**

- 🟥 Residents 20%
- 🟦 Commuters 9%
- 🟨 Visitors 45%
- 🟪 Unclassified 26%

# Monitoring big events

**Presence of Visitors GSM - Pisa -  Historical Center**





**Presence of Visitors GSM - Lucca Comics 2012**

# Enabling technologies for Big Data analytics

# M-Atlas

An analytic platform to extract, store, combine different kinds of models to build mobility knowledge discovery processes.



WIND, ORANGE, GOUDAPPEL, Skype, OCTOTelematics, Telecom Italia, Toyota, ENEL, ISTAT IBM, Local administrations

One of best EU-FET results, invited for exhibition at Strasburg Parliament

# Privacy-by-Design in Big Data Analytics



**Cell Phone**

**Cell Phone**

**Cell Phone**

**Service Provider**

**Anonymization**

**Mining and Analytical Engine**

**InfoMobility**

**Socio-economic indicators**

**Health services**

**Anonymization is not trivial**

**De-identification is not enough**

**Text annotation and scalable analytics @ ACUBE LAB (Unipi)**

# Topic-based annotation

"Diego Maradona won against Mexico"

*Mexico's football team*

*Ex-Argentina's player*

Find anchors and annotate them with articles drawn from Wikipedia!

# Wikipedia is a rich source of instances



## Steve Jobs

From Wikipedia, the free encyclopedia

*For the biography, see Steve Jobs (book).*

**Steven Paul "Steve" Jobs** (/ˈdʒɒbz/; February 24, 1955 – October 5, 2011)[5][6] was an Arab-American[7] entrepreneur[8] and inventor,[9] who was the co-founder, chairman, and CEO of Apple Inc. Through Apple, he was widely recognized as a charismatic pioneer of the personal computer revolution and for his influential career in the computer and consumer electronics fields, transforming "one industry after another, from computers and smartphones to music and movies..."[12] Jobs also co-founded and served as chief executive of Pixar Animation Studios; he became a member of the board of directors of The Walt Disney Company in 2006, when Disney acquired Pixar. Jobs was among the first to see the commercial potential of Xerox PARC's mouse-driven graphical user interface, which led to the creation of the Apple Lisa and, one year later, the Macintosh. He also played a role in introducing the LaserWriter, one of the first widely available laser printers, to the market.[13]

After a power struggle with the board of directors in 1985, Jobs left Apple and founded NeXT, a computer platform development company specializing in the higher-education and business markets. In 1986, he acquired the computer graphics division of Lucasfilm, which was spun off as Pixar.[14] He was credited in *Toy Story* (1995) as an executive producer. He served as CEO and majority shareholder until Disney's purchase of Pixar in 2006.[15] In 1996, after Apple had failed to deliver its operating system, Copland, Gil Amelio turned to NeXT Computer, and the NeXTSTEP platform became the foundation for the Mac OS X.[16] Jobs returned to Apple as an advisor, and took control of the company as an interim CEO. Jobs brought Apple from near bankruptcy to profitability by 1998.[17][18][19]

**Steve Jobs**

Jobs holding a white iPhone 4 at Worldwide Developers Conference 2010

**Born**        Steven Paul Jobs

## PARC (company)

From Wikipedia, the free encyclopedia
(Redirected from PARC User Interface)

# Why is it a difficult problem?

"Diego Maradona won against Mexico"

- Ex-Argentina's coach
- ~~His nephew~~
- ~~Maradona Stadium~~
- ~~Maradona Movie~~
- ...

- ~~Mexico nation~~
- ~~Mexico state~~
- Mexico football team
- ~~Mexico baseball team~~
- ...

*Don't annotate!*

**TAGME** is a powerful tool that is able to identify *on-the-fly* meaningful short-phrases (called "spots") in an unstructured text and link them to a pertinent Wikipedia page in a fast and effective way. This annotation process has implications which go far beyond the enrichment of the text with explanatory links because it concerns with the *contextualization* and, in some way, the *understanding* of the text.
Try **TAGME** now!

You can play with the demo interface below or check the TAGME RESTful API we are currently supporting.

Currently **TAGME** is available in English and in Italian and it is based on Wikipedia snapshots of July, 2012.

**NEWS!** As of August 2012, new RESTful functions are available and new advanced parameters can be used. For instance, you can compute semantic relatedness between topics identified by TAGME, or enable the special parser for Twitter messages. Check the RESTful API page for further details.

Developed by Paolo Ferragina and Ugo Scaiella at A³ Lab
Dipartimento di Informatica, University of Pisa.

**Input Text**

Italiano English

On this day 24 years ago Maradona scored his infamous "Hand of God" goal against England in the quarter-final of the 1986

Many links

Few links

Reset

TAGME!

## Tagged text | Topics

On this day 24 years ago Maradona scored his infamous "Hand of God" goal against England in the quarter-final of the 1986

**Less links**

## Tagged text | Topics

On this day 24 years ago Maradona scored his infamous "Hand of God" goal against England in the quarter-final of the 1986

**More links**

Original paper:
IEEE Software 2012

Clustering appl.
ACM WSDM 2012

**Details on…**

**http://acube.di.unipi.it/tagme**

**Regional project with:**

net7

SPAZIODATI

Classification appl.
ECIR 2012

- Nodes $\cong$ users, entities **(~ 1 bil)**

- Edges *explicit* = friend, follower… **(~ 10 bil)**

- Edges *implicit* = similarity, click… **(» 100 bil)**

- Textual Data = post, tweet, news… **(> Ptb)**

The Knowledge Graph

Learn more about one of the key breakthroughs behind the future of search.

See it in action

Discover answers to questions you never thought to ask, powered by the Knowledge Graph.

# New searching algorithms

The goal is:

- Minimize the occupied space
- Maximize the **substring-search** throughput

2007 → 

Under US-patenting

We were the first to show
how to search *bzip*-ed data

> 600 citations



Total citations

Citations per year

89

0

2000  2001        2005        2009        2013

Scholar articles   Opportunistic data structures with applications
P Ferragina, G Manzini - Foundations of Computer Science, 2000. Proceedings. ..., 2000
Cited by 435 - Related articles - All 29 versions

2003 → 

# Collaborations



4 submitted US patents: Yahoo and NY Univ.
2 accepted US patents: Rutgers Univ. and AT&T-Lucent

# SoBigData

*Bootstrap Workshop*

# Towards a European Laboratory

## on Big Data Analytics and Social Mining

**18 July** 2013, h **10:00-17:30**
Auditorium of the National Research Council
*Area della Ricerca* **CNR**, Via Moruzzi 1, **Pisa**

## Program

**10:00: Reception and refreshments**

**10:30-11:30: Setting the stage**
**Welcome address** D. Laforenza (President CNR Research Campus), N. De Francesco (Univ. Pisa, Vice-Rector)
**Opening remarks** C. Montani (ISTI-CNR, Director), F. Turini (Dip. Informatica Univ. Pisa, Director).
**Towards a Euro Lab on Big Data Analytics & Social Mining** M. Conti (DIITET-CNR, Director)
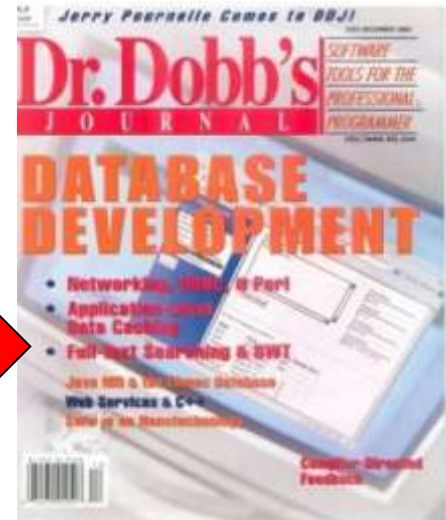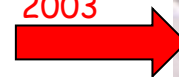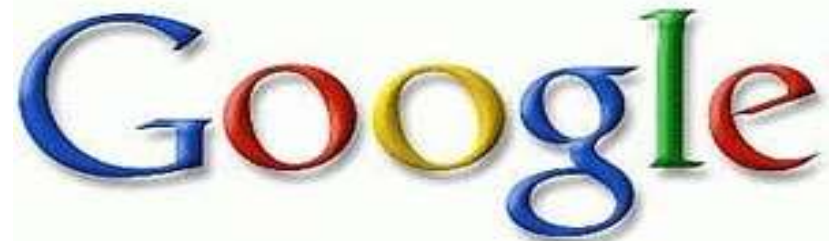**Big Data Analytics & Social Mining for Science and Society** F. Giannotti (ISTI-CNR)
**Democratizing big data: the ethical challenges of social mining** D. Pedreschi (Univ. Pisa)

**11:30-12:00: Keynote**
**Big data big insights: the coming age of computational social science**
D. Lazer, Professor of Political Science and Computer and Information Science, Northeastern University, Boston, MA. Director,
Program on Networked Governance, Harvard University

**12:00-13:30: Panel – Big data & social mining: new models for participation and policy making**
S. Targetti (Regione Toscana, Vice-President), F. Accordino (EC), F. Sestini (EC),E. Baldacci(ISTAT), C. Comella (Garante Privacy),
F. Marzano (Stati Innovazione)
Moderator: L. De Biase (Il Sole 24 Ore)

**13:30-14:30: Lunch break**

**14:30-16:00: SoBigData pills**
**Monitoring trend and engagement with social media mining** M. Tesconi (IIT-CNR)
**Exploring the structure of society**
A. Passarella (IIT-CNR)
**Sentiment quantification and opinion surveys**
F. Sebastiani (ISTI-CNR)
**Good answers for difficult questions** R. Perego (ISTI-CNR)
**Understanding human mobility** C. Renso (ISTI-CNR | Univ. Pisa)
**Big data in finance and economics**
F. Lillo (Scuola Normale Superiore), G. Caldarelli (IMT Lucca)
**Big data and official statistics – monitoring poverty/well-being at any scale**
M. Pratesi (Univ. Pisa), F. Maggino (Univ. Firenze)
**Do you need a big computer or a great algorithm ?** P. Ferragina (Univ. Pisa)

**16:00-17:30: Panel – Big data & social mining: new models for social innovation and business**
R. Soru (Tiscali), O. Cicchetti (Telecom Italia), T. Martino (Octotelematics), G. Gigliucci (ENEL Ricerca),
A. Di Benedetto (CNA, Ass. Giovani Industriali)
Moderator: L. De Biase (Il Sole 24 Ore)

**17:30: Conclusion**

**BIG DATA**

**Registration: www.sobigdata.eu/registration**