

Crowdsourcing mobile networks from the **macac** experiment

Katia Jaffrès-Runser

University of Toulouse, INPT-ENSEEIHT,
IRIT lab, IRT Team

Ecole des sciences avancées de Luchon

Networks and Data Mining, Session II

July 1st, 2015

2 The smartphone phenomenon

- Multiple sensing and communication capabilities
 - Sensors, camera, GPS, microphone
 - 3G, WiFi, Bluetooth, etc.
 - Storage capabilities (several Gbytes)
 - Computing power



Mobile Traffic is growing constantly

- Increasing volume of mobile data between 2014-2018
 - “...worldwide mobile data traffic will increase nearly **11-fold** over the next four years and reach **an annual run rate of 190 exabytes** (10^{18}) by 2018...”
 - 54% of mobile connections will be ‘smart’ connections by 2018

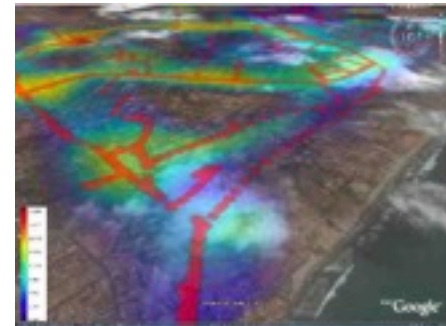
[Cisco VNI Global Mobile Data Traffic Forecast (2013-2018)]



+



=



In 2013, 4.1 billion
users worldwide

Next Big Networking Challenge: meet traffic demand !

1. If data is not delay sensitive:
 - e.g. Videos, Application / system updates, music, podcasts, etc.

Leverage opportunistic encounters to route or flood **delay tolerant** data hop by hop

Benefit: Reduce downloads from infrastructure wireless network

2. If several connectivity options exist:
 - e.g. 3G/4G, WiFi, Femto cells

Offload / Pre-fetch data using the 'best' available connectivity, at the best time and location

Benefit: Load balancing between available infrastructures

Crowdsourcing (part of) this huge network

- This huge network of users is constantly active.
 - The context each user is evolving in is changing
 - The content each user is consuming / sending is evolving as well
- To provide the next intelligent data communications, **we need to understand how this network evolves**
- How is this big dynamic network evolving?
 - Getting network traces
 - Model the interactions of this dynamic network to capture its evolution
- How to get network traces?
 - Network operator monitoring (cf. Marco's talk)
 - **Crowdsourcing using smartphone capabilities (this talk)**

Outline of this talk

1. Crowdsourcing using smartphone capabilities
 - Building a Mobile app for that
 - First statistics of Macaco Project
2. Classifying social interaction from such contact traces
 - RECAST algorithm

EU CHIST-ERA MACACO Project

Mobile context-Adaptive CAching for COntent-centric networking

www.macaco.inria.fr

INRIA (Paris), University of Toulouse, SUPSI (Lugano),
University College London, CNR-IEIT (Torino), UFMG (Brazil)



Crowdsourcing Mobile app

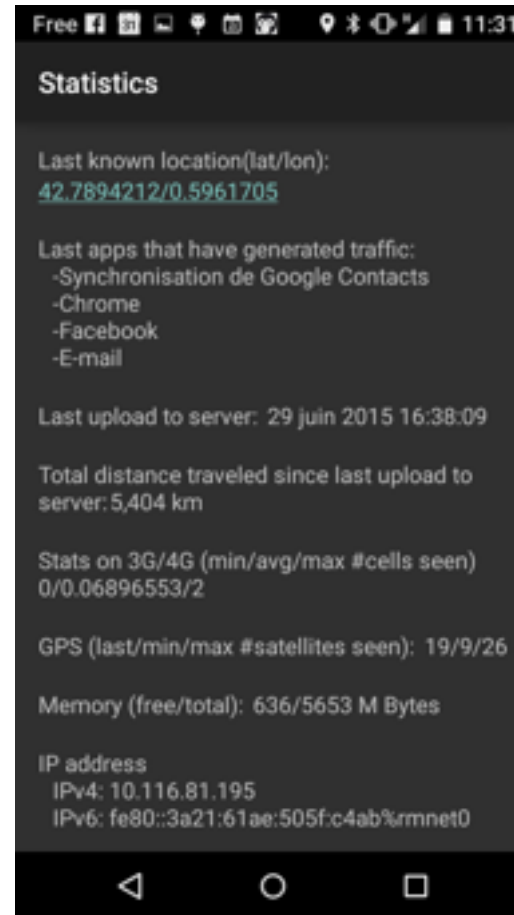
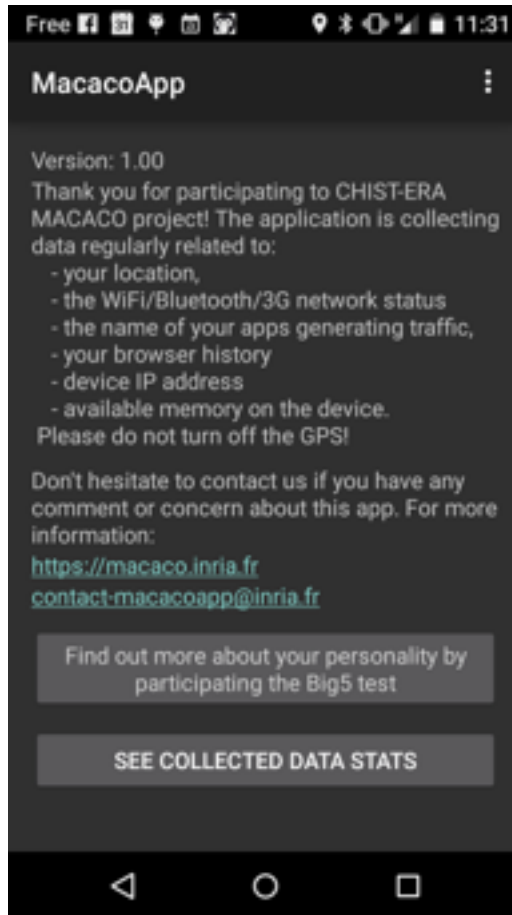
Goal : Sample user context and content data

- Runs in background on **volunteer** phone users
 - Monitors different sensors periodically (5 mins)
- Should be **seamless** with respect to regular phone usage
 - **Upload data** to our servers before memory is full
 - Full memory = no reactivity
 - But : does not ruin the 3G data plan !
Favor uploads on WiFi
- **Energy** constraint !!
 - Monitoring all sensors is costly

The macaco App

www.macaco.inria.fr

Available on Play Store



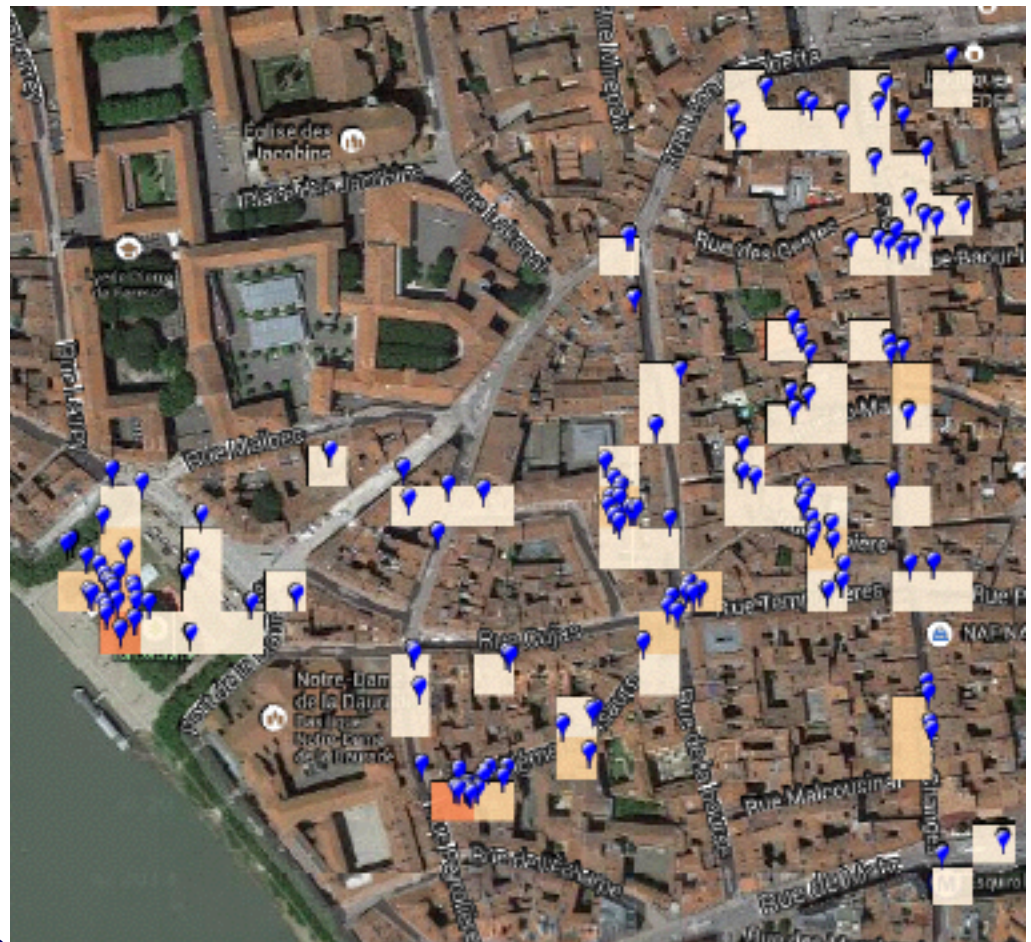
Measured data every 5 minutes :

- Context data

- Location (GPS, Internet)
- WiFi connectivity
- Bluetooth connectivity
- Cellular network towers
- Battery discharge
- Accelerometer
- Big 5 personality test

- Content data

- Name of applications that have generated traffic
- Browser history
- Name of applications run



Main issue: getting volunteers :-)

- **Privacy** issues (discussion with CNIL)
 - Keep data within project partners,
 - Have data anonymized (hashed IMEI - location)
 - Limit storage duration of non-anonymized data use
 - Option to remove its own data from the collection
- **Energy efficient app design**
 - Keep the volunteers using the app
- Provide **a motivation** for participating
 - Added value of the app (e.g. visualize its own data, game, ...)
 - Financial retribution (voucher)
 - Lottery
 - For the greater good :-) ...

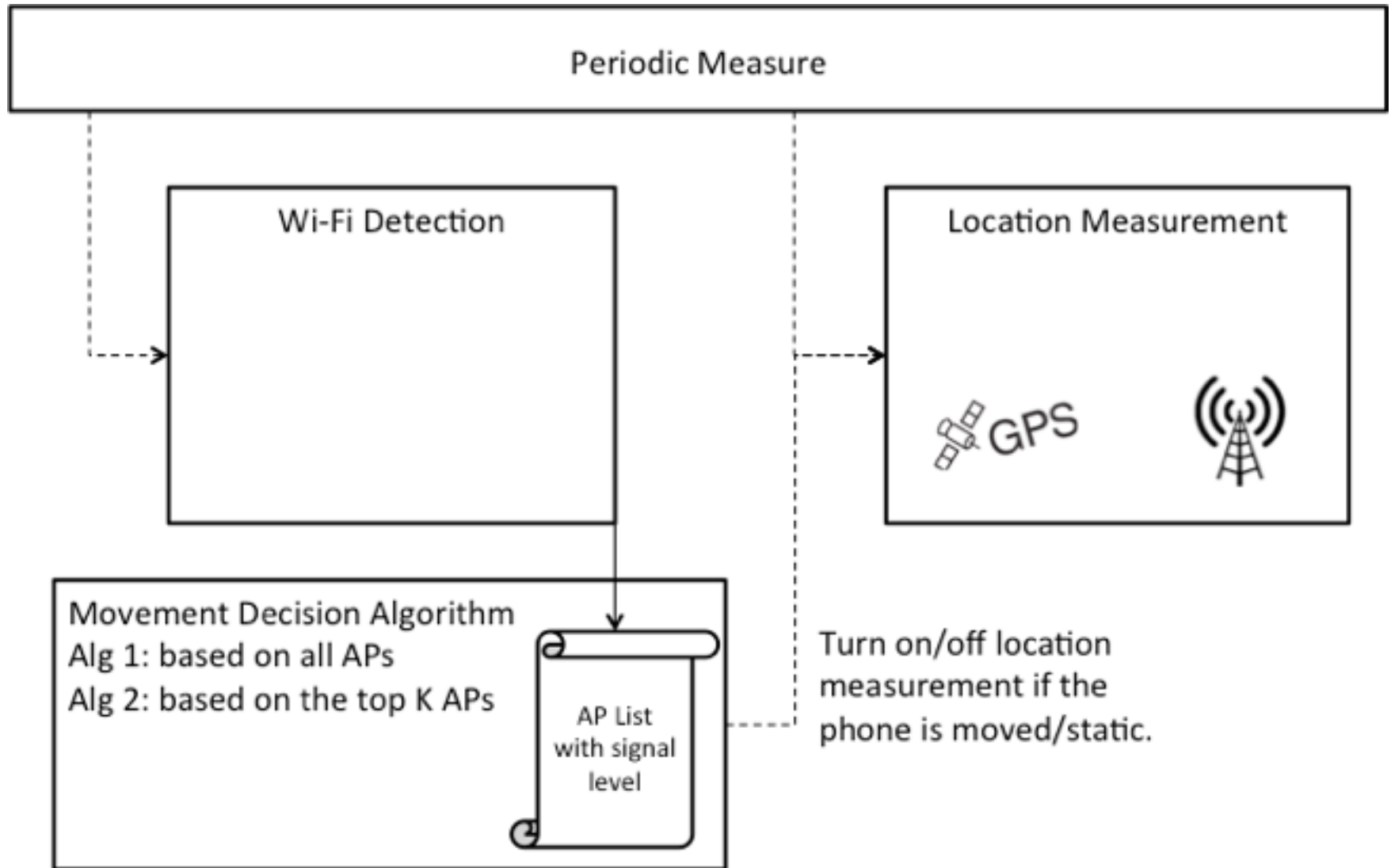
Energy aware design

- Energy hungry sensors:
 - **GPS** localisation
 - Unavailable indoors
 - Useless if no motion -> DETECT MOVEMENT
 - **Bluetooth** scan
 - Use Low-Energy bluetooth
 - Useful to detect available opportunistic communications
 - **Accelerometer**
 - Reduce the sampling duration and interval

Movement detection algorithm

- Main idea
 - if (Movement detected)
then trigger GPS measurement
- Two options:
 - Use accelerometer / gyroscope sensors : only works if the user is moving during the sampling duration + additional energy
 - Leverage for 'free' the **wireless networking context**
- Wireless networking context:
 - List of received signal strength (RSSI) for all APs measured at current location

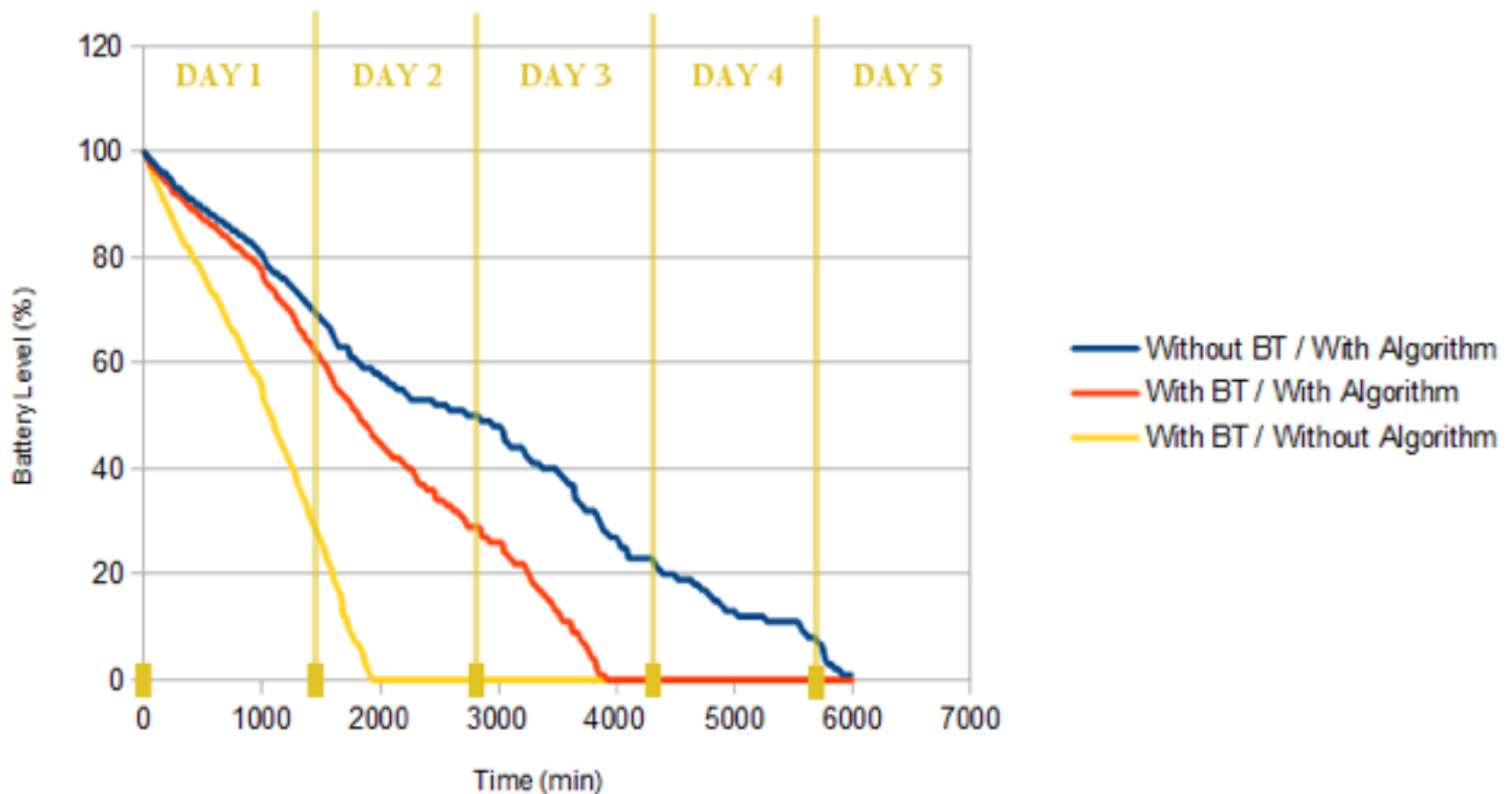
Motion detection algorithm



Energy depletion with movement detection

% remaining battery if the phone stands still

- w./w.o. movement detection
- w./w.o. bluetooth measurements



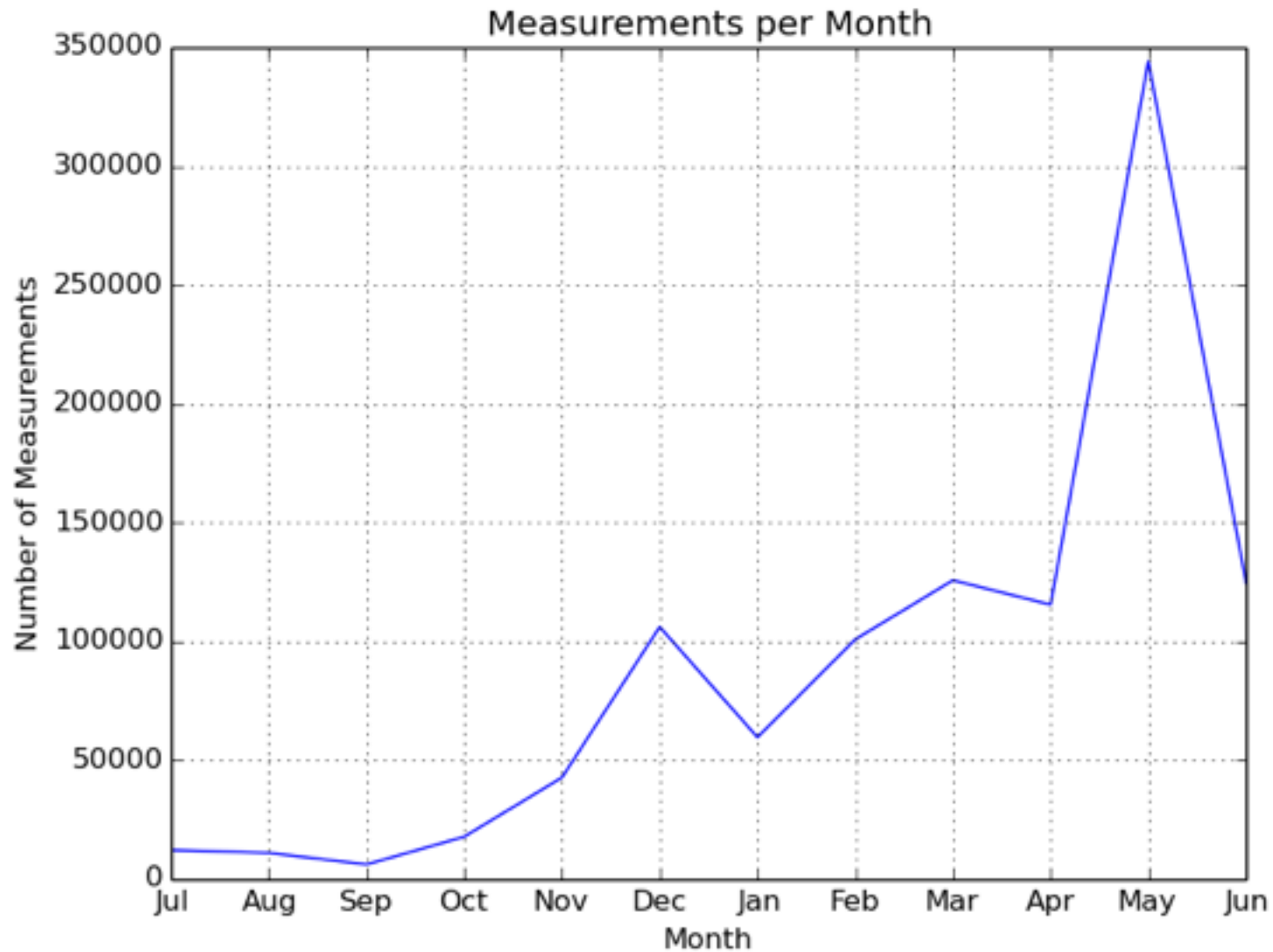
First Macaco data statistics

- Collected with MacacoApp
 - Up to now, for one year (2014 July – 2015 June)
 - 57 devices over one year
 - 1,069,083 Measurements

- Top contributors:

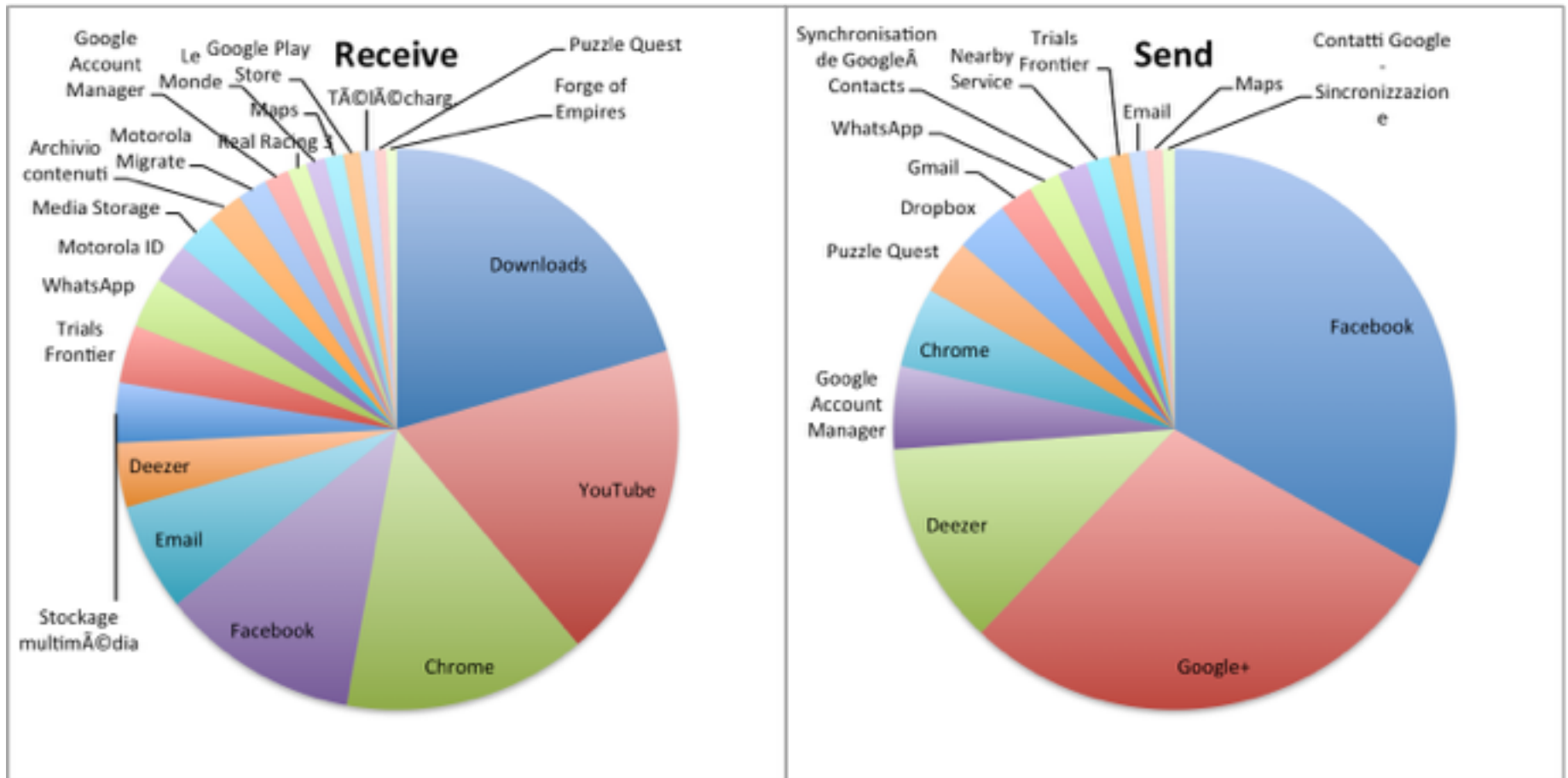
Hash(IMEI)	Period	# measurements
203a...	2014-11-04 - 2015-06-22	187879
bacd...	2014-08-27 - 2015-06-22	145619
f1d9...	2014-08-06 - 2015-06-20	126215
46bd...	2014-08-19 - 2015-06-13	119634
4517...	2012-01-01 - 2015-06-22	65812
e6d2...	2015-05-05 - 2015-06-22	59997
008f...	2015-05-07 - 2015-06-22	55059

First Macaco data statistics



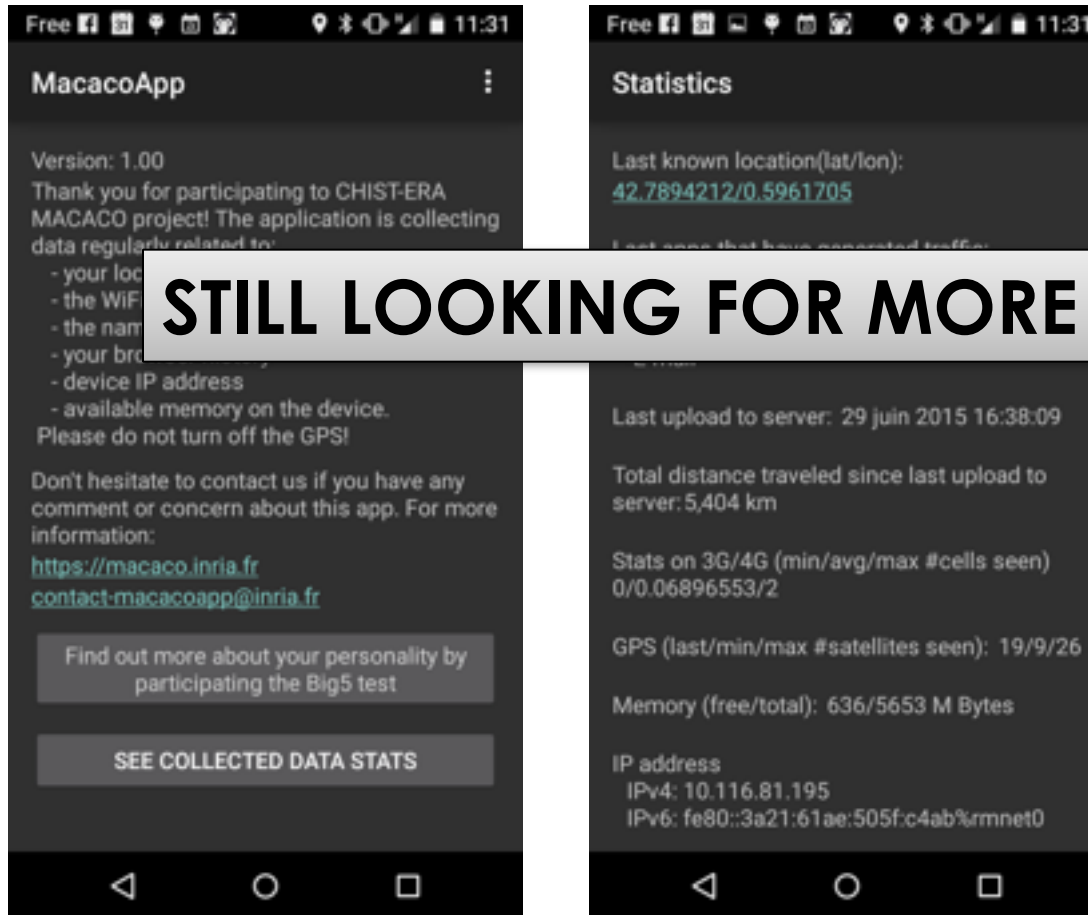
First Macaco data statistics

- Total traffic download: 55534 MB
- Total traffic upload: 10679 MB



Mobile context-Adaptive CAching for COntent-centric networking

www.macaco.inria.fr

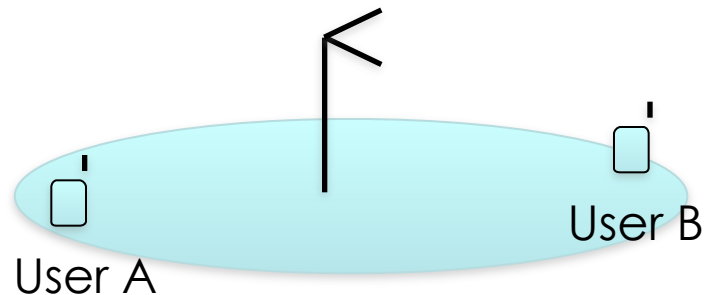


STILL LOOKING FOR MORE VOLUNTEERS :-)

 App
Available on Play Store

How to exploit such datasets?

- Other open datasets exist (cf. Crowdad <http://crowdad.cs.dartmouth.edu/>)
- Different types of temporal contact measurements
 - Measure a direct link between User A and B (e.g. Bluetooth, WiFi Direct connectivity)
 - Assume a link exists between User A and User B if they are connected to the same WiFi access point



- Measure location of users (GPS): if users are close enough, assume they are connected
- MACACO : adds the content dimension to the context

Example open data sets

Data collection to build *contact traces*

- ▶ Log the contact time and duration of a node to an access point
- ▶ Log the GPS coordinates of mobile nodes regularly

Derive a time-varying contact graph

Dataset	Local	# entities	Duration	Type	Avg. # encounters/ node/day
Dartmouth ¹	campus	1156	2 months	Individuals	145.6
USC ²	campus	4558	2 months	Individuals	23.8
San Francisco ³	City	551	1 month	Cabs	834.7

- ▶ Dartmouth and USC collect connection dates/durations to WiFi APs,
- ▶ San Francisco collects GPS locations of taxi cabs.

¹T. Henderson et al. "The changing usage of a mature campus-wide wireless network," in Proc. of ACM MobiCom 2004

²W. jen Hsu et al. "Impact: Investigation of mobile-user patterns across university campuses using wlan trace analysis," CoRR, vol. abs/cs/0508009, 2005

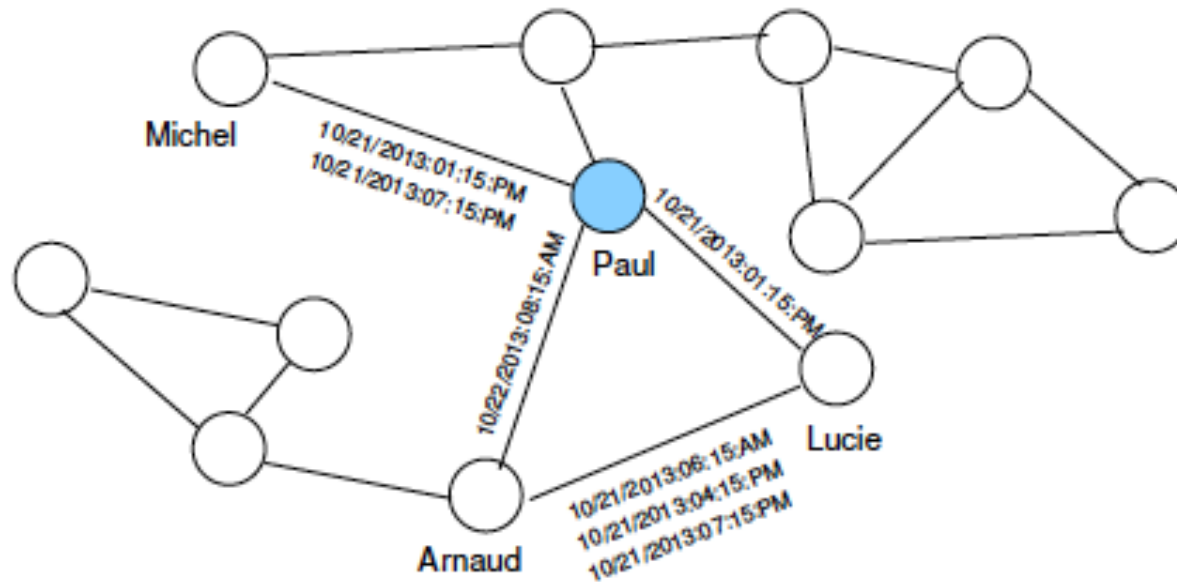
³A. Rojas et al. "Experimental validation of the random waypoint mobility model through a real world mobility trace for large geographical areas," in Proc. of the 8th ACM MSWiM 2005

Rationale and related initiatives

Characterize **interactions**, i.e. edges of contact graph

- ▶ Regularity of contacts : How often did Arnaud and Paul meet per day? during the whole trace?

Miklas et al.⁴ determine whether 2 nodes are *friends* or *strangers* using an empirical threshold (friends encounter 10 times or more within 14 weeks).

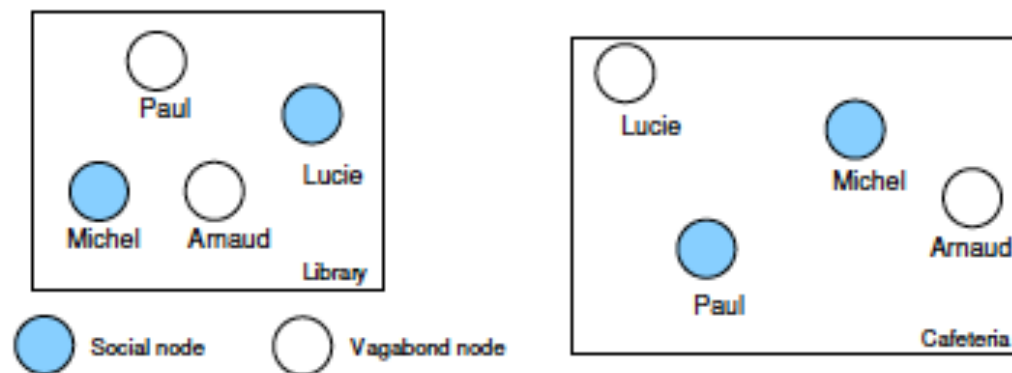


⁴A. G. Miklas et al., "Exploiting social interactions in mobile systems," in *Proceedings of the UbiComp '07*

Rationale and related initiatives

Characterize **node's** behavior, i.e. vertices of contact graph

Using localization information, Zyba et al.⁵ differentiate *social* from *vagabond* nodes. Socials appear regularly in a given area while vagabonds visit an area rarely and unpredictably.



- ▶ Monitor the total appearance and regularity of appearance

Paul is social at the cafeteria but vagabond at the library: a per node/per area approach → *geographical dependency*

⁵G. Zyba, G. Voelker, S. Ioannidis, and C. Diot, "Dissemination in opportunistic mobile ad-hoc networks: The power of the crowd," in *Infocom'11*

RECAST classifier ^[1]

- Characterizes the interactions of nodes based on their probability to originate from a random or social behavior
- Identify different kinds of social interactions (friends, acquaintances, bridges or random)
- No geographical dependency, i.e., is of general validity

Together with

Pedro O. Vaz de Melo, Antonio Loureiro – UMFG Brazil

Aline Viana - Inria, Marco Fiore - IIT-CNR Italy

Frédéric Le Mouël – INSA Lyon

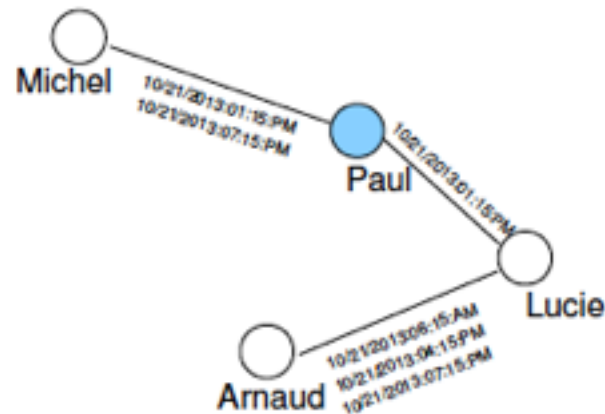
[1] RECAST: Telling Apart Social and Random Relationships in Dynamic Networks, P. Olmo Vaz de Melo, A. Viana, M. Fiore, K. Jaffrès-Runser, F. Le Moüel and A. A. F. Loureiro, 16th ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems (ACM MSWim 2013), Barcelona, Spain, 3-8 November 2013.

Graphs extracted from contact traces

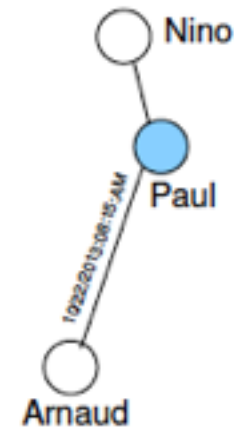
Two possible representations

1. δ event graph: $\mathcal{G}_k(\mathcal{V}_k, \mathcal{E}_k)$

There is an edge in \mathcal{E}_k if contact within $\delta = 1$ day for instance.



Day 1 event graph $\mathcal{G}_1(\mathcal{V}_1, \mathcal{E}_1)$



Day 2 event graph $\mathcal{G}_2(\mathcal{V}_2, \mathcal{E}_2)$

2. Accumulative graph $G_t(\mathcal{V}_t, \mathcal{E}_t)$

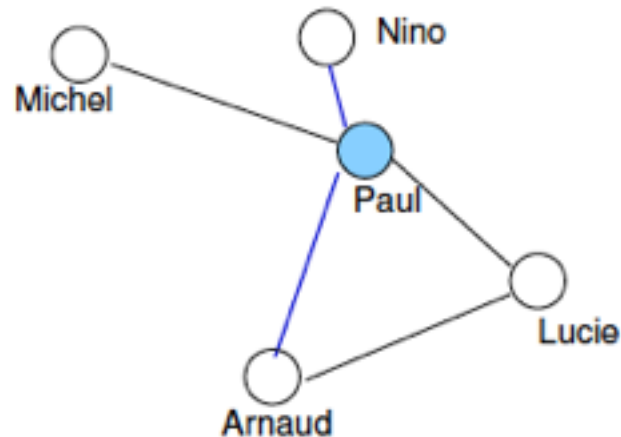
Graphs extracted from contact traces

Two possible representations

1. δ event graph: $\mathcal{G}_k(\mathcal{V}_k, \mathcal{E}_k)$

There is an edge in \mathcal{E}_k if contact within $\delta = 1$ day for instance.

2. Accumulative graph $G_t(V_t, E_t)$: $G_t = \{\mathcal{G}_1 \cup \mathcal{G}_2 \cup \dots \cup \mathcal{G}_t\}$



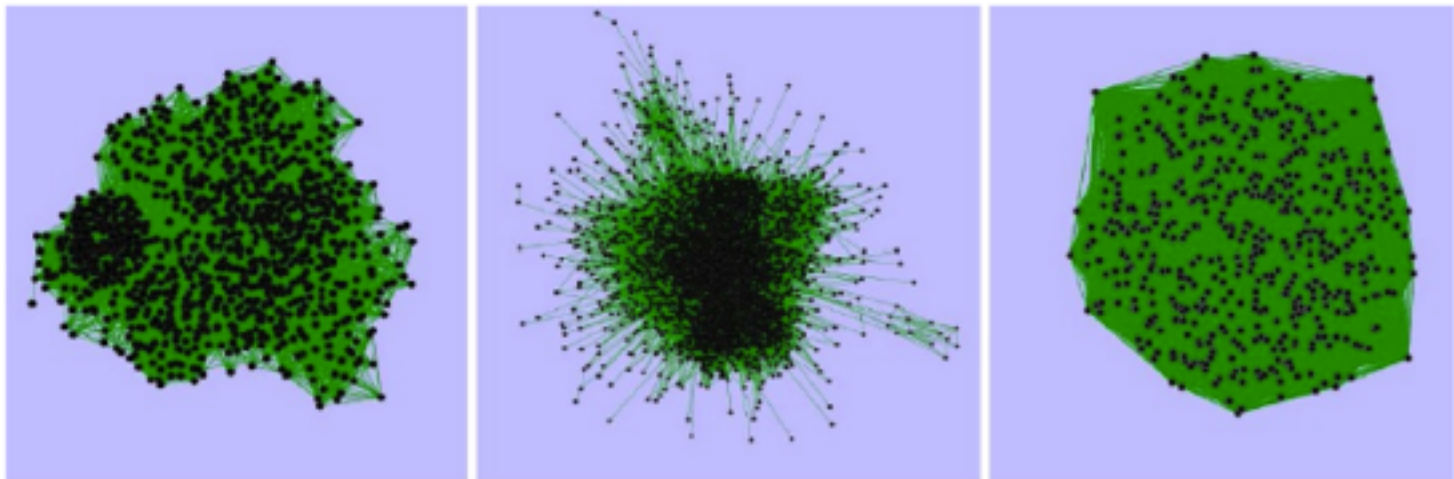
$G_2(V_2, E_2)$ Accumulative graph up to Day 2

Accumulates all event graphs up to time step t .

Graphs extracted from contact traces

Example accumulative graph G_t for $t = 2$ weeks

For $\delta = 1$ day and using force-direct layout algorithm for plotting



(a) Dartmouth

(b) USC

(c) San Francisco

Seems difficult to extract any knowledge from these social graphs:
→ gathers all social AND random interaction!

Social graph and its random counterpart

Random graph equivalent of G

Calculate a **random graph** G^R from a graph $G(V, E)$:

- ▶ Keep same number of vertices and edges,
- ▶ Randomly assign edges to keep the same node degree distribution using *RND* algorithm⁶:

An edge is set between nodes of degree d_i and d_j with probability

$$p_{ij} = (d_i \times d_j) / \sum_{k=1}^{|V|} d_k$$

Random accumulative graph G_t^R

Random accumulative graph derived from event graphs $\{\mathcal{G}_i\}_{i \in [1, \dots, t]}$

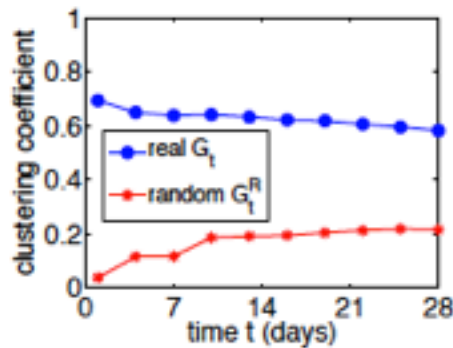
$$G_t^R = \{RND(\mathcal{G}_1) \cup RND(\mathcal{G}_2) \cup \dots \cup RND(\mathcal{G}_t)\}$$

⁶F. Chung and L. Lu, "Connected Components in Random Graphs with Given Expected Degree Sequences," *Annals of Combinatorics*. Nov. 2002

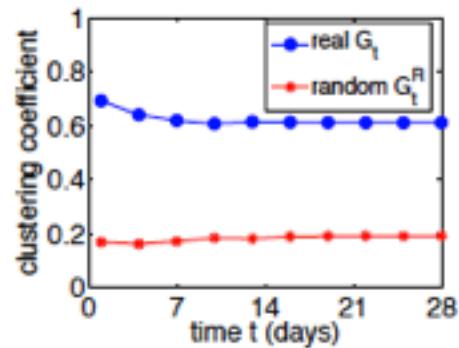
Comparison social vs. random graphs

Network clustering coefficient can identify a network with an elevated number of clusters (i.e. communities).

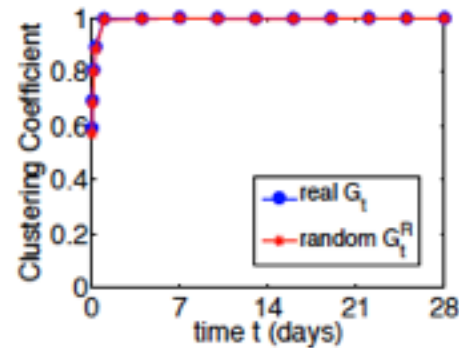
- ▶ If $\bar{c}c(G) \gg \bar{c}c(G^R)$, parts of the decisions of the nodes of G are NOT random



(a) Dartmouth



(b) USC



(c) San Francisco

- ▶ Dartmouth / USC traces have an order of magnitude higher $\bar{c}c$ than $G^R \rightarrow$ social decisions
- ▶ San Francisco: each individual taxi in the trace encounters most of the other taxis \rightarrow closer to a random behavior

Social network features: Regularity and Similarity

Social nodes' behavior tend to

- ▶ repeat on a regular basis (because of daily activities for instance)
→ Regularity
- ▶ build persistent communities and generate common acquaintances
→ Similarity

Mathematical metrics

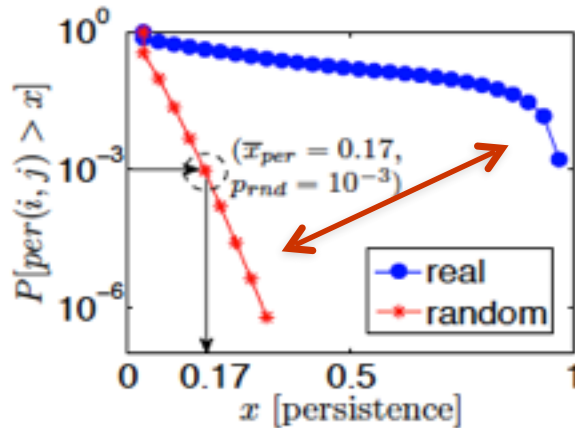
- ▶ **Edge persistence** $per(i, j)$ ⁷ :
Percentage of time steps an edge exists over the past discrete time steps in the event graphs $\{\mathcal{G}_i\}_{i \in [1, \dots, t]}$
- ▶ **Topological overlap** $to(i, j)$ ⁸ :
Ratio of neighbors shared by two nodes calculated for the accumulative graph G_t .

⁷N. Eagle et al., "From the Cover: Inferring friendship network structure by using mobile phone data," Proceedings of the National Academy of Sciences, Sept. 2009

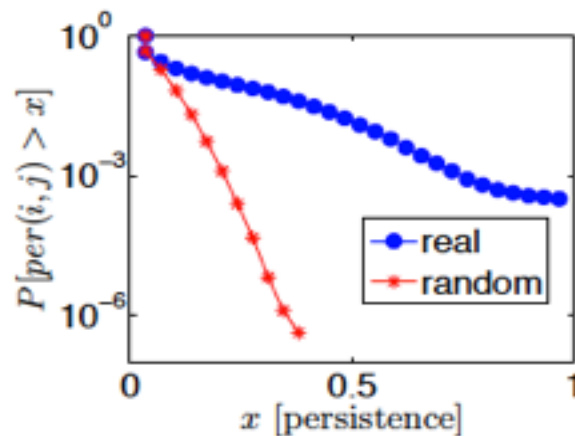
⁸J. P. Onnela et al., "Structure and tie strengths in mobile communication networks", Proc. of the National Academy of Sciences, May 2007

CCDF of edge persistence after 4 weeks

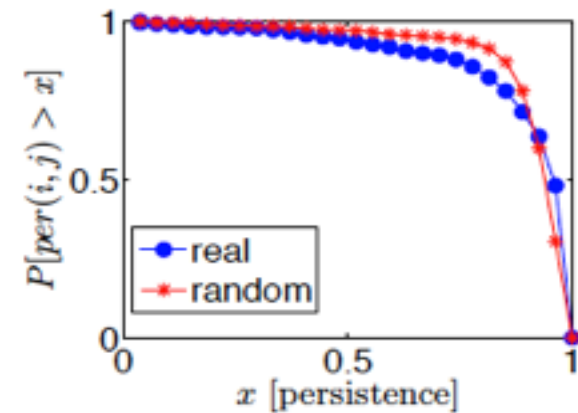
Individuals tend to see each other regularly



(d) Dartmouth



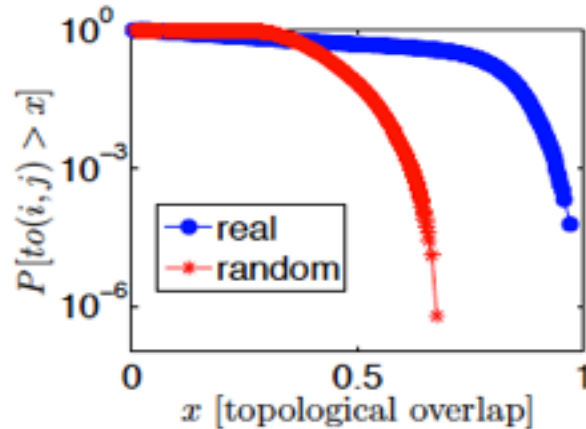
Encounters occur almost in a random fashion



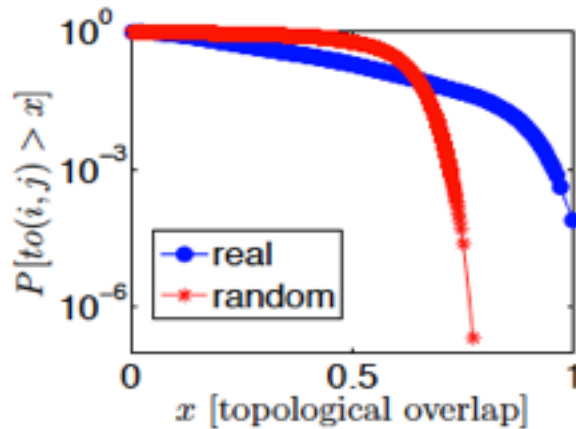
(f) San Francisco

CCFD of topological overlap after 4 weeks

Individuals of G_t have common neighbors

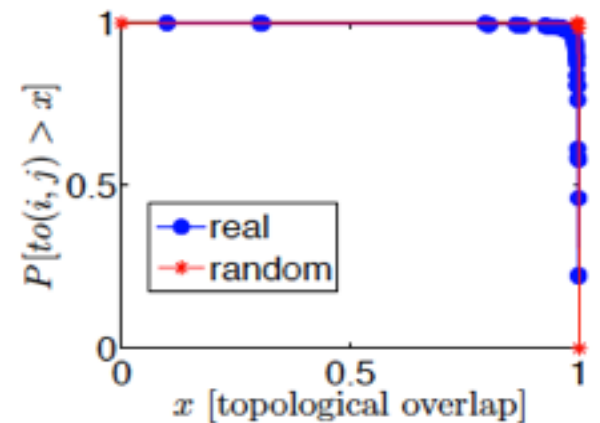


(g) Dartmouth



(h) USC

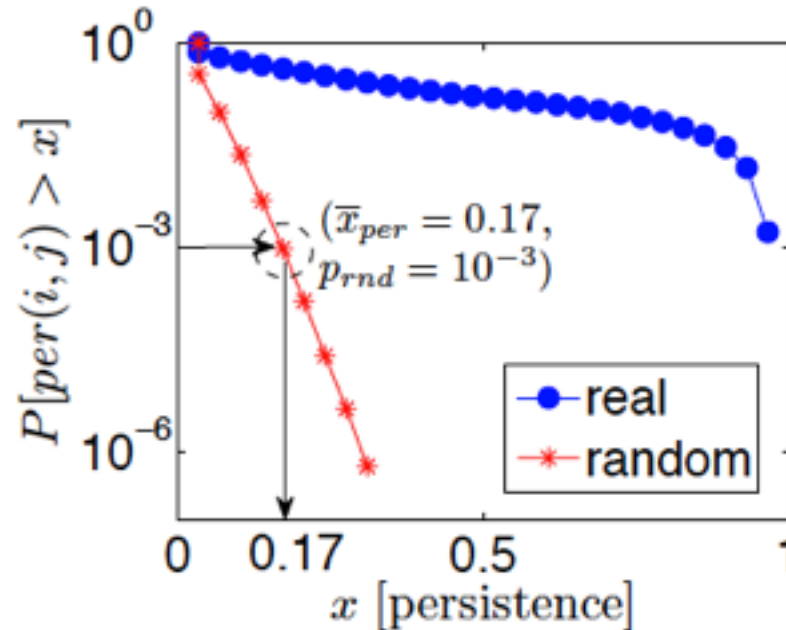
Common neighbors occur in a random fashion



(i) San Francisco

Social vs. Random Edges

In the **random network**, we only have a probability of 10^{-3} to have edges with a persistence of more than $\bar{x}_{per} = 0.17$.



→ Thus, in the **social graph** G_t :

- ▶ edges with $\text{per}(i, j) > \bar{x}_{per}$ can be classified as *social edges*
- ▶ edges with $\text{per}(i, j) < \bar{x}_{per}$ can be classified as *random edges*

Note that there is a p_{rnd} chance that a social edge is actually random (mis-classification)

RECAST classification algorithm

Only parameter of RECAST: p_{rnd} , the mis-classification error bound.

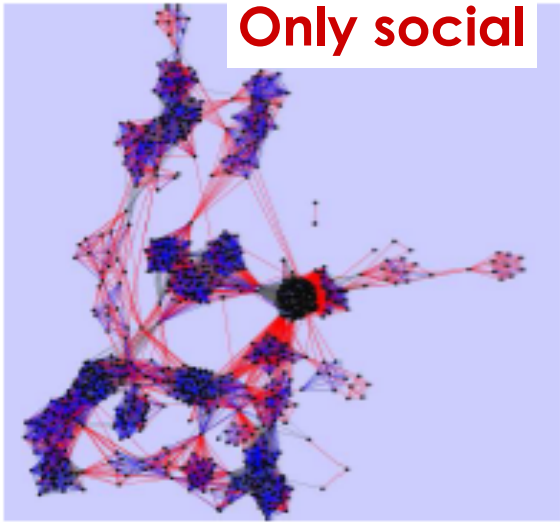
Main steps

- ▶ Calculate the $per(i,j)$ and $to(i,j)$ for each edge
- ▶ Knowing p_{rnd} , calculate \bar{x}_{per} and \bar{x}_{to} from CCDF's
- ▶ For each edge,
 - ▶ if $per(i,j) > \bar{x}_{per} \rightarrow (i,j)$ is **social** for edge persistence
else (i,j) is random for edge persistence
 - ▶ if $to(i,j) > \bar{x}_{to} \rightarrow (i,j)$ is **social** for topological overlap
else (i,j) is random for topological overlap
- ▶ Classify edges into classes of relationships according to:

Class	Edge persistence	Topological overlap
Friends	social	social
Acquaintances	random	social
Bridges	social	random
Random	random	random

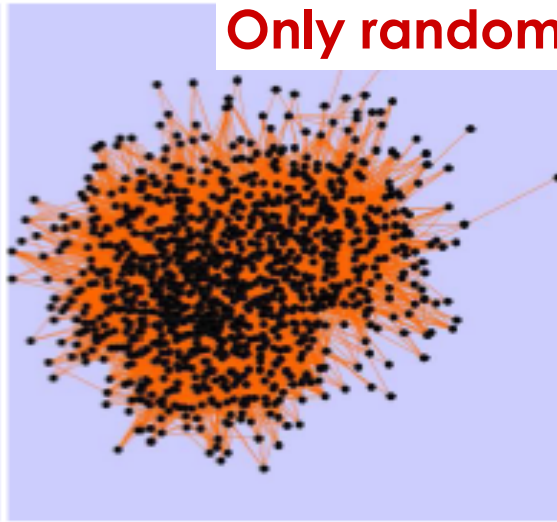
Classification after 2 weeks

Only social



(a) Dartmouth, only social edges

Only random

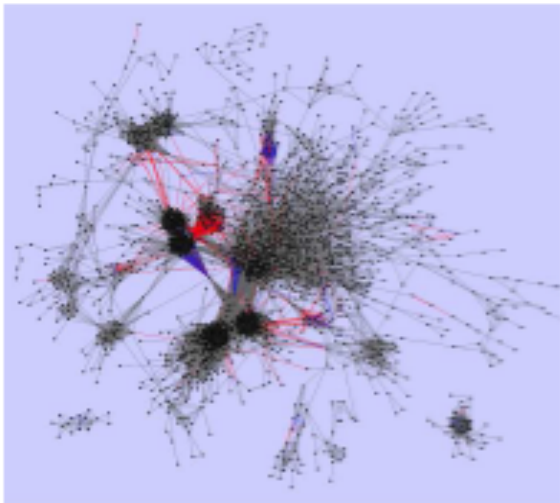


(b) Dartmouth, only random edges

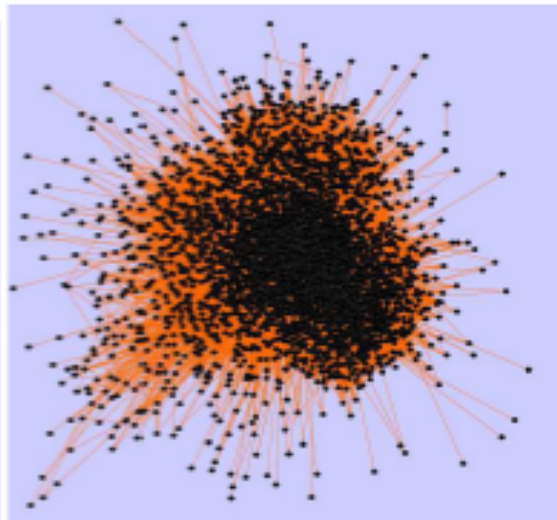
Friends edges are in blue
Bridges edges are in red
Acquaintance edges are in gray
Random edges are in orange

- **Social-edges network**
Complex structure of *Friendship* communities, linked to each other by *Bridges* and *Acquaintanceship*

- **Random-edges network** No structure appears, looking like random graphs



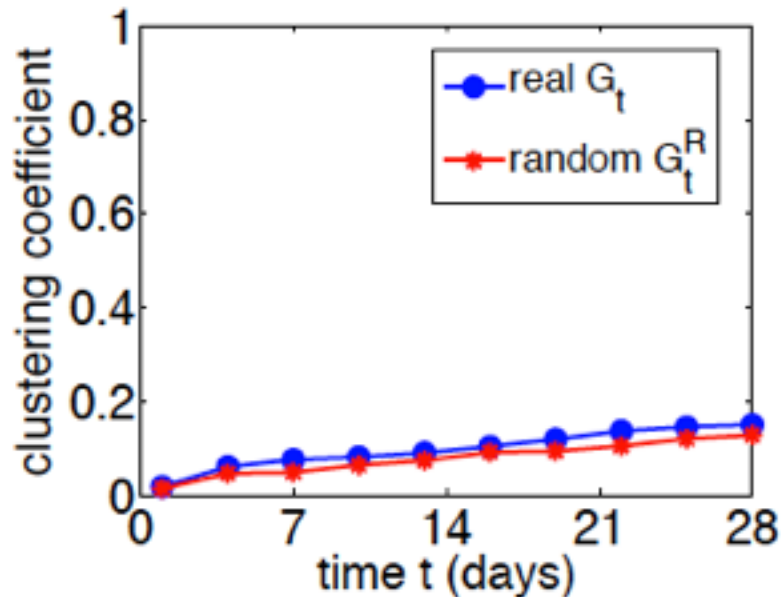
(c) USC, only social edges



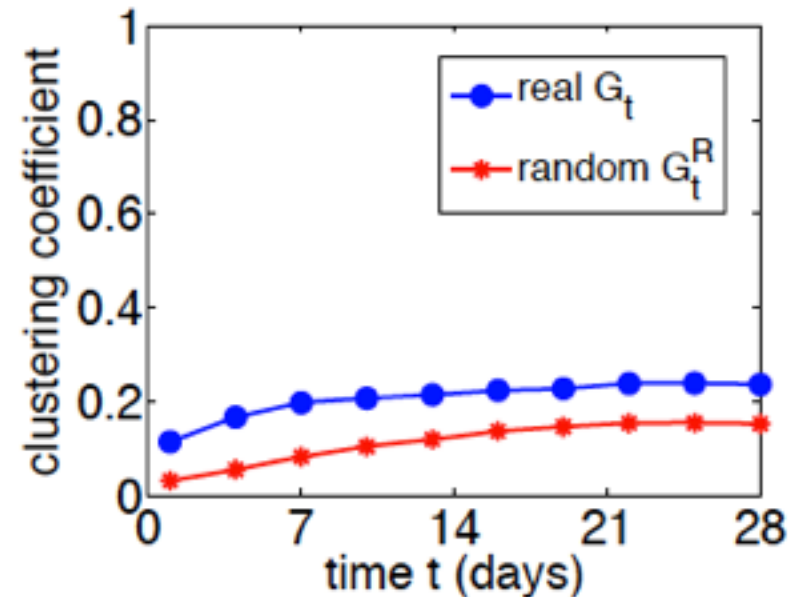
(d) USC, only random edges

Cluster coefficient analysis for **random edges only**

Dartmouth



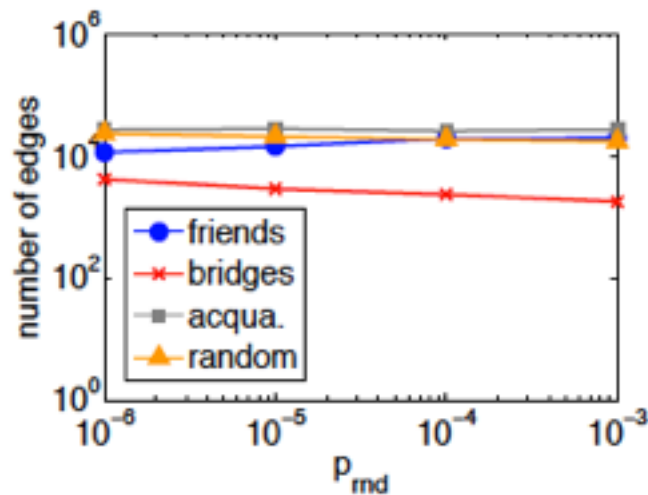
USC



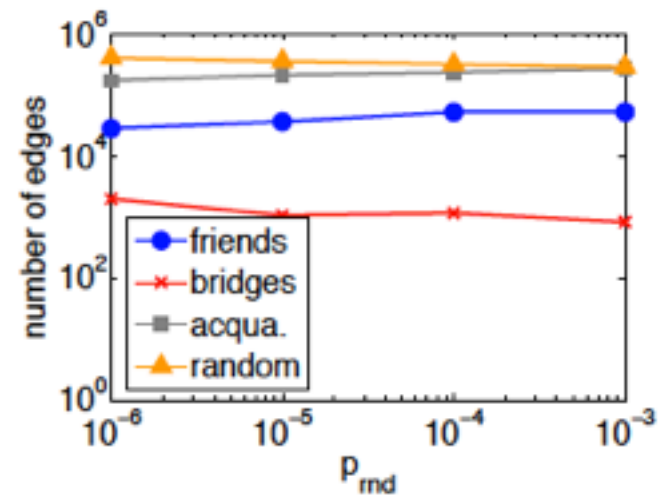
Validates the efficiency of RECAST to identify random edges for Dartmouth and USC

Impact of p_{rnd}

Number of edges of a each class that appear in the first 4 weeks vs. p_{rnd}



Dartmouth



USC

RECAST is not sensitive to p_{rnd} !

Forwarding using relationship information

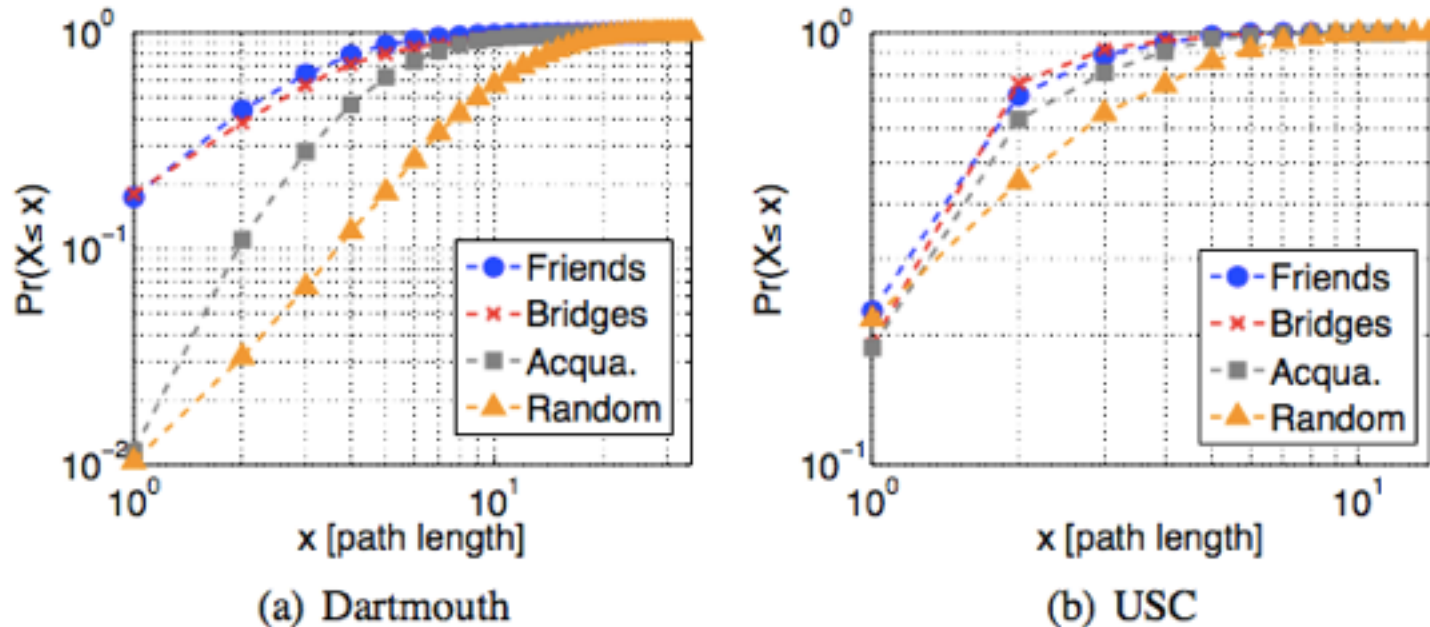
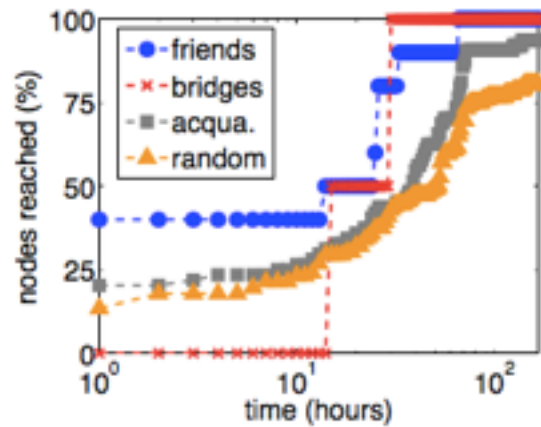
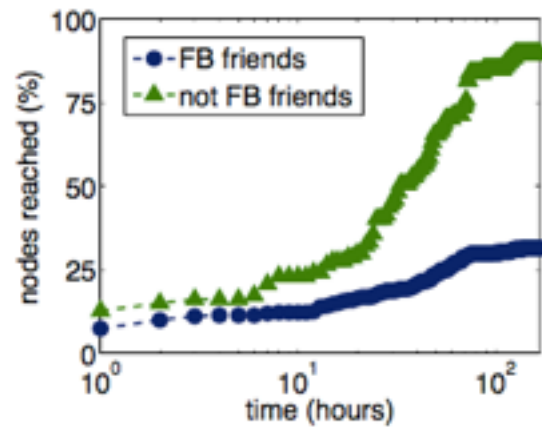


Figure 11: The histogram of the path lengths of messages between users i and j who share a determined class of relationship.

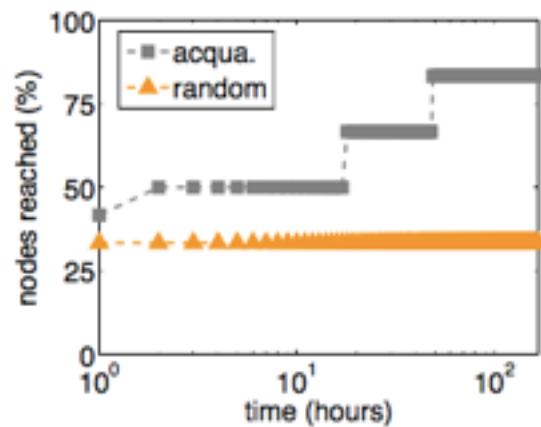
Forwarding with recast or FB data



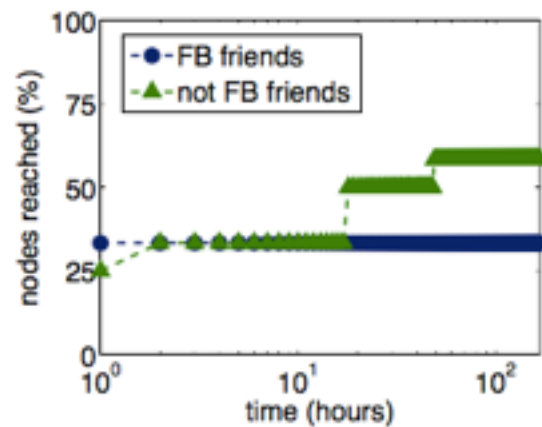
(a) Sassy: RECAST



(b) Sassy: Facebook



(c) UPB: RECAST



(d) UPB: Facebook

Figure 12: The % of users who were reached over time grouped by RECAST classes and Facebook (FB) friendship.

Next...

- Having this data, exhibit the correlations between content and context
 - Do users have regular habits in data usage?
 - If yes, is it possible to model these networks with the content plane in mind?
- Using network models, deriving data pre-fetching strategies to adjust the load off available networks

....