

Mobility data storage and analysis

Antònia Tugores, Pere Colet



Index

- Motivation
- Data storage
- Insert/Query performance
- Preliminary results and on going work

Motivation

Study urban mobility by using their tweets (specially the geolocated ones)

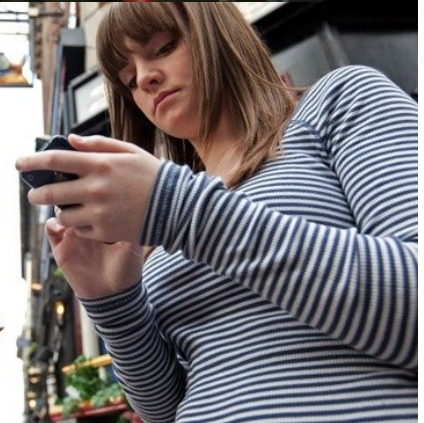
Compare with other on going studies

Barcelona

London

Zurich

Not only tweets (in the future)



Getty Images



JSON document

JavaScript Object Notation
text-based open standard
designed for human-readable data
interchange
language-independent
parsers available for many languages
alternative to XML
fields are in an arbitrary order

GeoJSON

open format for encoding collections of simple feature objects
points
line strings
polygons
and multi-part collections of points, lines and polygons
with their non-spatial attributes
using JavaScript Object Notation

```
1 {
2   "contributors": null,
3   "truncated": false,
4   "text": "Home, sweet home! I demà, a \u00passar la tarda\u00 a Valencia. ^^",
5   "in_reply_to_status_id": null,
6   "id": 277909591700942849,
7   "retweeted": false,
8   "coordinates":
9     {
10      "type": "Point",
11      "coordinates": [2.13029439, 41.3665202]
12    },
13   "source": "foursquare",
14   "in_reply_to_screen_name": null,
15   "id_str": "277909591700942849",
16   "retweet_count": 0,
17   "in_reply_to_user_id": null,
18   "favorited": false,
19   "source_url": "http://foursquare.com",
20   "user":
21     {
22      "geo":
23        {
24          "type": "Point",
25          "coordinates": [41.3665202, 2.13029439]
26        },
27      "in_reply_to_user_id_str": null,
28      "possibly_sensitive": false,
29      "created_at": "2012-12-09T22:56:24",
30      "in_reply_to_status_id_str": null,
31      "place":
32        {
33          "country_code": "ES",
34          "url": "http://api.twitter.com/1/geo/id/laca5ffbc553e95b.json",
35          "country": "Spain",
36          "place_type": "city",
37          "bounding_box":
38            {
39              "type": "Polygon",
40              "coordinates": [[[2.086323, 41.336062], [2.137433,
41              41.336062], [2.137433, 41.380548], [2.086323, 41.380548]]]]
42            },
43          "full_name": "Hospitalet de Llobregat, Barcelona",
44          "attributes": {},
45          "id": "laca5ffbc553e95b",
46          "name": "Hospitalet de Llobregat"
47        },
48      "id": "5107d50e66946222fce81631"
49    }
50 }
```



Tweet

(+ RT info)

tweet content (140 characters)

id

geolocation (if enabled)

source

user information

creation time

place (if defined)

fields are in an arbitrary order
fields scheme can change

```
1 {
2   "contributors": null,
3   "truncated": false,
4   "text": "Home, sweet home! I demà, a \u0022passar la tarda\u0022 a Valencia. ^^
5   (@ Granvia Centre w/ @rubenralc) [pic]: http://t.co/7iA2XK52",
6   "in_reply_to_status_id": null,
7   "id": 277909591700942849,
8   "retweeted": false,
9   "coordinates":
10  {
11    "type": "Point",
12    "coordinates": [2.13029439, 41.3665202]
13  },
14   "source": "foursquare",
15   "in_reply_to_screen_name": null,
16   "id_str": "277909591700942849",
17   "retweet_count": 0,
18   "in_reply_to_user_id": null,
19   "favorited": false,
20   "source_url": "http://foursquare.com",
21   "user":
22  {
23    "geo":
24    {
25      "type": "Point",
26      "coordinates": [41.3665202, 2.13029439]
27    },
28    "in_reply_to_user_id_str": null,
29    "possibly_sensitive": false,
30    "created_at": "2012-12-09T22:56:24",
31    "in_reply_to_status_id_str": null,
32    "place":
33    {
34      "country_code": "ES",
35      "url": "http://api.twitter.com/1/geo/id/1aca5ffbc553e95b.json",
36      "country": "Spain",
37      "place_type": "city",
38      "bounding_box":
39      {
40        "type": "Polygon",
41        "coordinates": [[[2.086323, 41.336062], [2.137433,
42        41.336062], [2.137433, 41.380548], [2.086323, 41.380548]]]
43      },
44      "full_name": "Hospitalet de Llobregat, Barcelona",
45      "attributes": {},
46      "id": "1aca5ffbc553e95b",
47      "name": "Hospitalet de Llobregat"
48    },
49    "_id": "5107d50e66946222fce81631"
50  }
51 }
```



description

language

fields are in an arbitrary order

id

on

description

language

fields are in an arbitrary order

```

21 "user":
22 {
23   "follow_request_sent": false,
24   "profile_use_background_image": true,
25   "contributors_enabled": false,
26   "id": 260287521,
27   "verified": false,
28   "profile_image_url_https":
29 "https://si0.twimg.com/profile_images/1533280952/tw_12403594_1315422640_normal.
30 jpg",
31   "profile_sidebar_fill_color": "252429",
32   "profile_text_color": "666666",
33   "followers_count": 128,
34   "protected": false,
35   "id_str": "260287521",
36   "default_profile_image": false,
37   "location": "L'Hospitalet de Llobregat",
38   "utc_offset": 3600,
39   "statuses_count": 1747,
40   "description": "Necesito meterme algo dentro, un poco de café o
41 algo... Y luego, de alguna manera, el mundo será un poco mejor. (Sam Vimes,
42 Hombres de Armas, Terry Pratchett)",
43   "friends_count": 570,
44   "profile_link_color": "AA0000",
45   "profile_image_url": "http://a0.twimg.com/profile_images/1533280952/
46 tw_12403594_1315422640_normal.jpg",
47   "notifications": null,
48   "geo_enabled": true,
49   "profile_background_color": "000000",
50   "profile_background_image_url": "http://a0.twimg.com/profile_background_images/
51 600683524/n9duka667b9ju3n9wgw4.jpeg",
52   "screen_name": "jX09A",
53   "lang": "ca",
54   "following": false,
55   "profile_background_tile": false,
56   "favourites_count": 7,
57   "name": "Jordi Cant\u00f3",
58   "url": null,
59   "created_at": "2011-03-03T15:45:46",
60   "profile_background_image_url_https": "https://si0.twimg.com/
61 profile_background_images/600683524/n9duka667b9ju3n9wgw4.jpeg",
62   "time_zone": "Madrid",
63   "profile_sidebar_border_color": "181A1E",
64   "default_profile": false,
65   "is_translator": false,
66   "listed_count": 1
67 }

```

Data size estimation



2KB/tweet

Tweet: aprox. 20 key/value

User: aprox 30 key/value

- id
- content (140chars)
- geolocation
- creation time
- place
- user information
 - id
 - screen name
 - description
 - creation time
 - location
 - friends count
 -

...

Data size estimation



2KB/tweet x 15million tweets/day
30 GB/day

Tweet: aprox. 20 key/value

User: aprox 30 key/value

- id
- content (140chars)
- geolocation
- creation time
- place
- user information
 - id
 - screen name
 - description
 - creation time
 - location
 - friends count

....

...

2500million tweets/year
6TB/year



<http://library.uoregon.edu/node/3475>

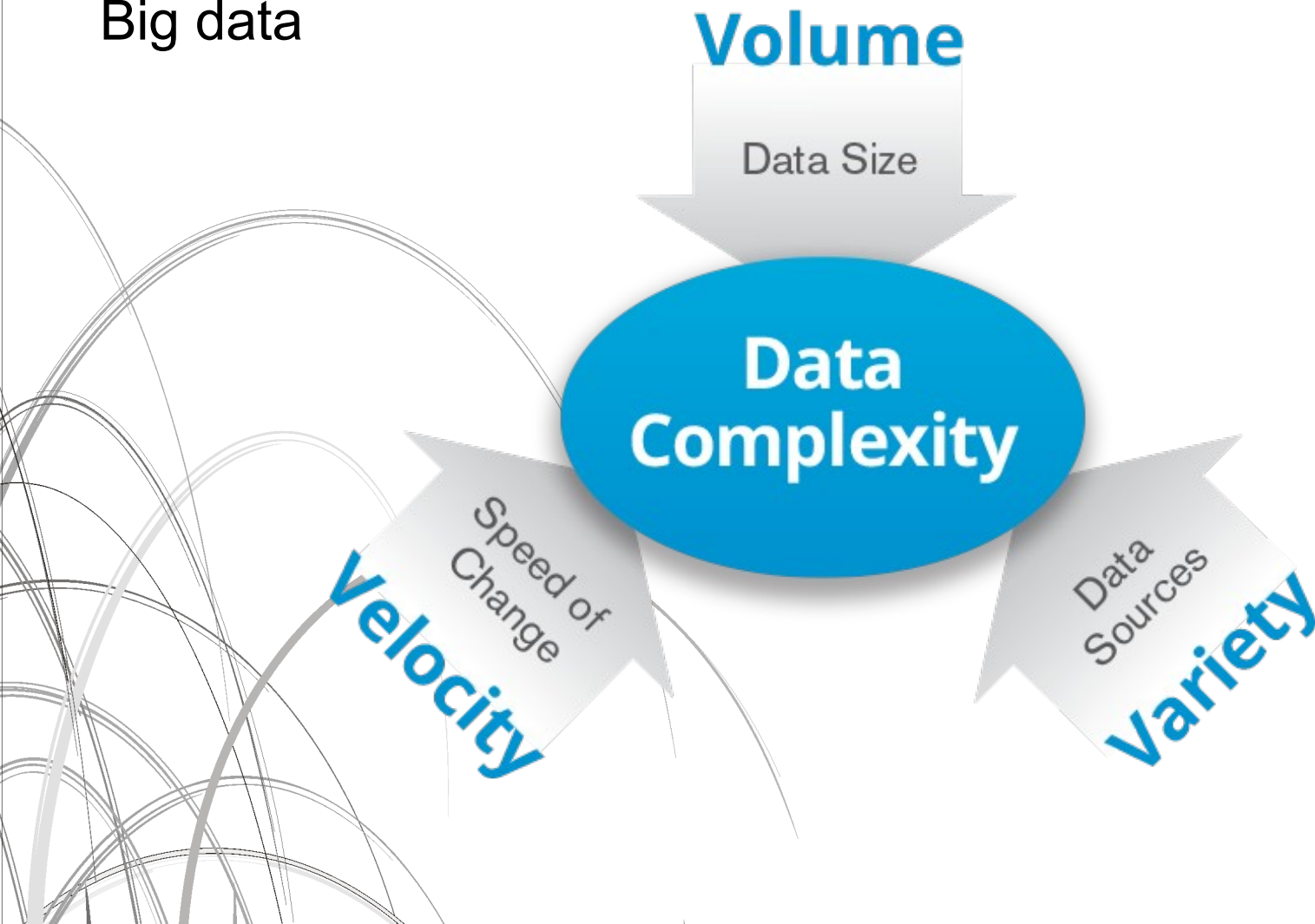
http://www.barclayandtodds.com/pages/24-7_Service.cfm

<http://ifisc.uib-csic.es>



Big data is like teenage sex:
everyone talks about it, nobody
really knows how to do it, everyone
thinks everyone else is doing it, so
everyone claims they are doing it ...

Big data



We need to...

Efficiently store data

Efficiently manage data (select subsets)

Efficiently analyse data



<http://bitzermobile.com/wp-content/uploads/2012/07/puzzle-piece.png>

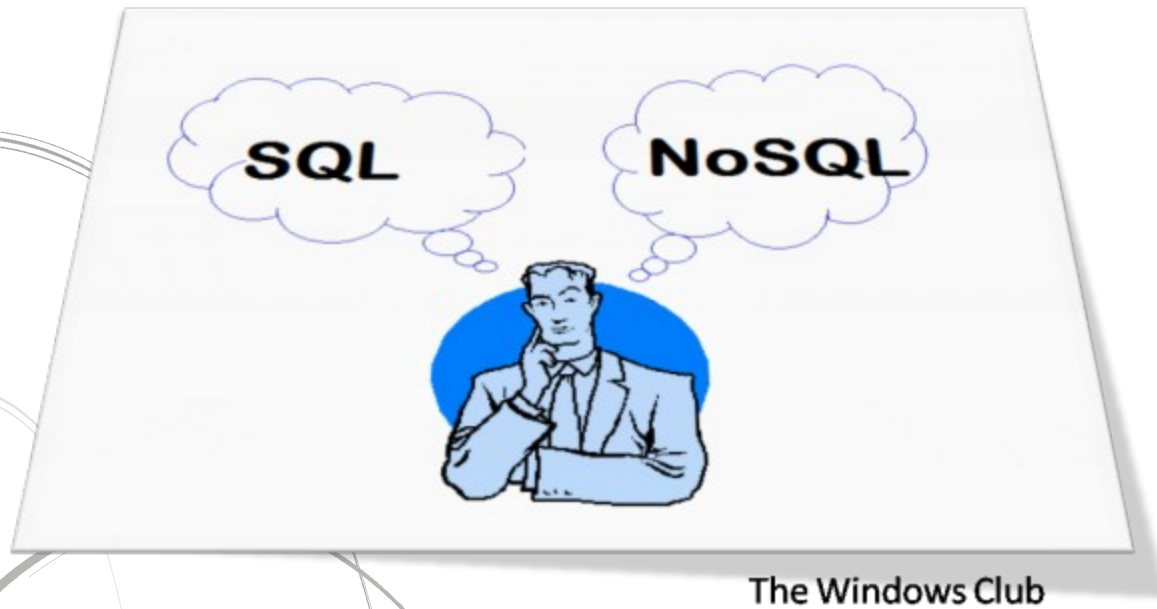
http://library.uoregon.edu/sites/default/files/node3393/data_management.jpg

<http://t2.gstatic.com/images?q=tbn:ANd9GcTIGtxf7mKEl6IU0ebDQhPHgq1zTJzQssuDQMEv5Kw7C7N4GGgw>

Requirements to store and manage the data

- * Capability to store **billions** of documents
- * Fast storage rate
- * Scalability
- * High availability
- * High search performance
- * Adaptative format (twitter can change the data format)

Databases



Databases

Relational (SQL) databases

has a collection of tables of data items all of which is formally described and organized according to the relational model

Examples

- Oracle
- PostgreSQL
- MySQL



PostgreSQL

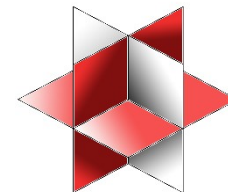


Non relational databases

A NoSQL database provides a mechanism for storage and retrieval of data that uses looser consistency models than traditional relational databases.

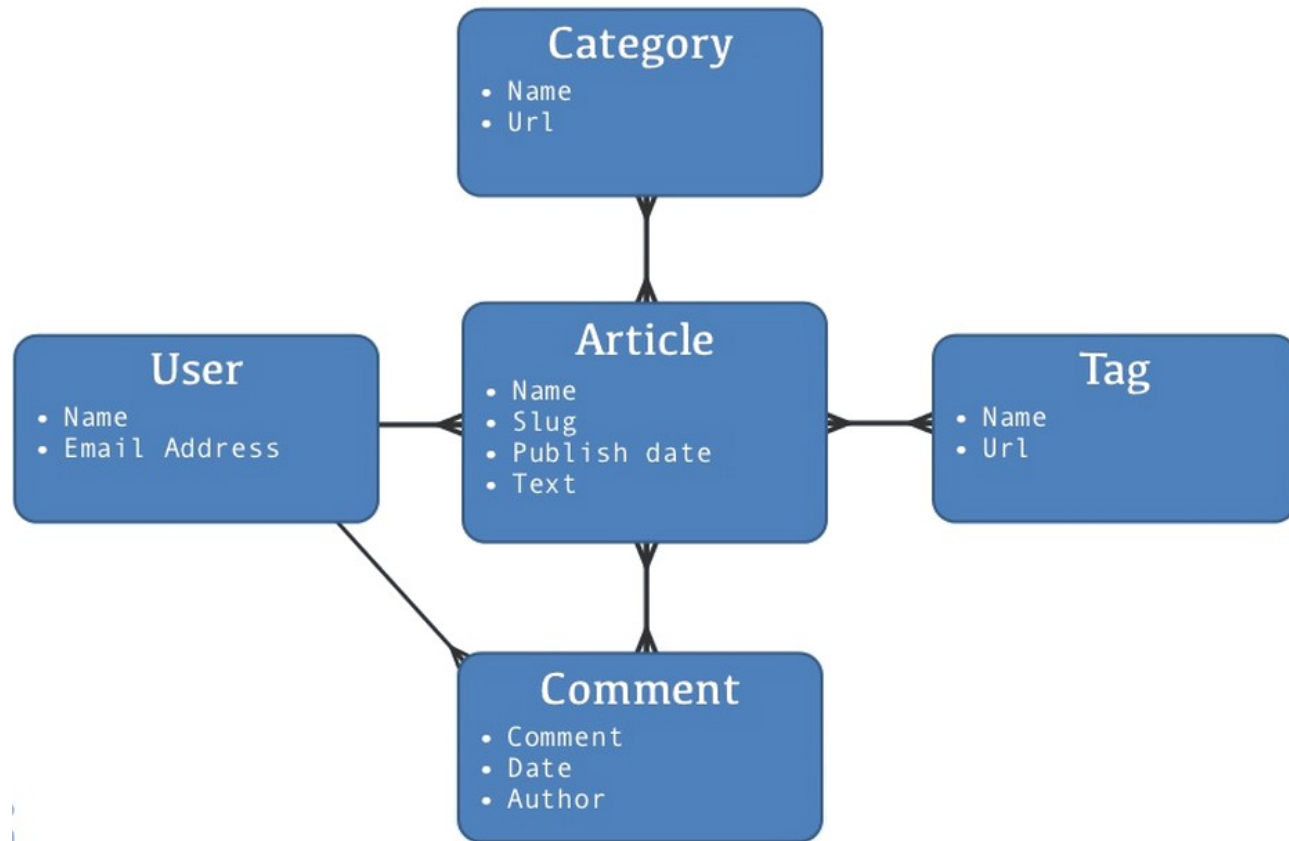
Examples

- MongoDB
- CouchDB
- HyperDex



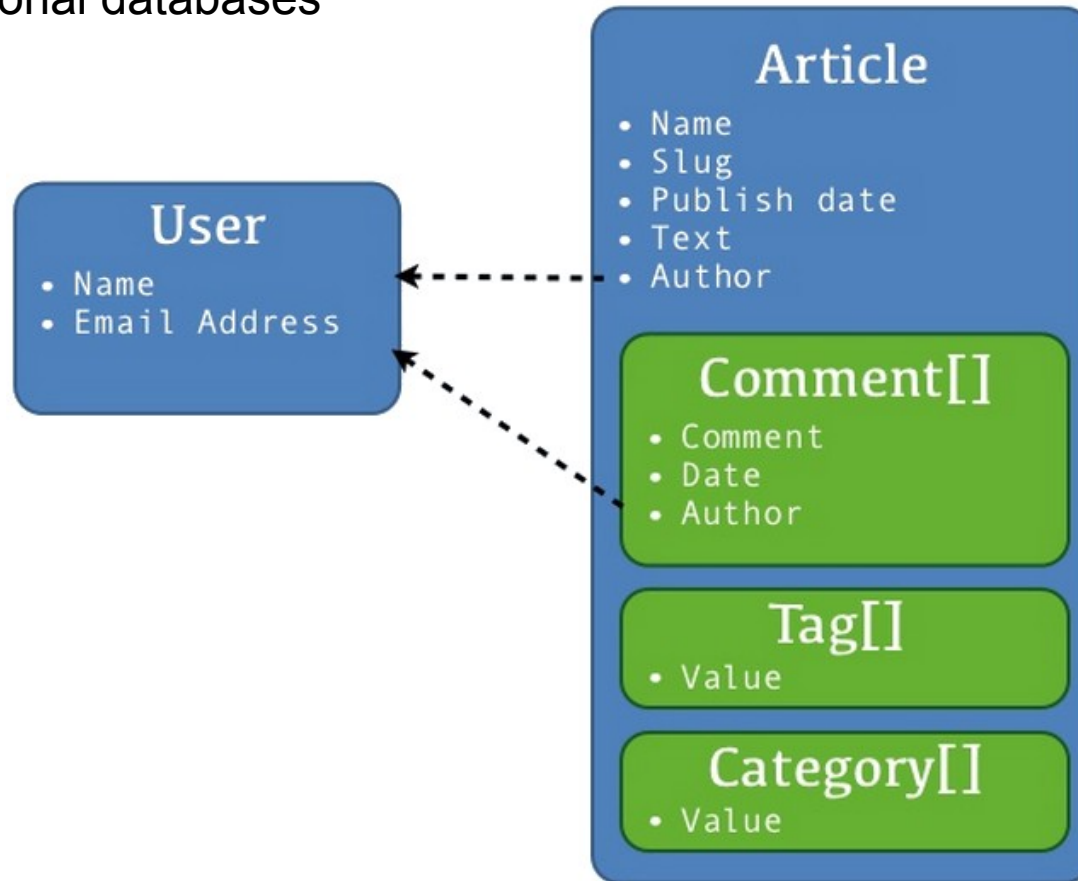
Databases

Relational (SQL) databases



Databases

Non Relational databases



DB comparison

SQL

noSQL

High performance

Indexes

Indexes

High availability

Automatic master failover
and recovery

Replica sets with automatic
master failover and recovery

Easy scalability

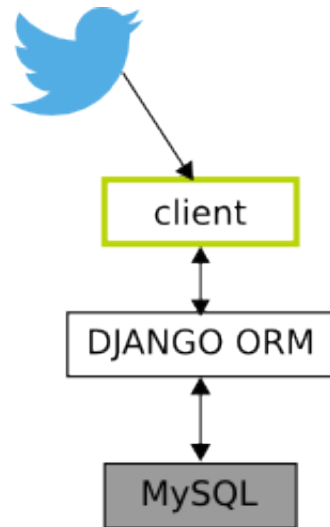
Hardware: \$\$\$\$\$\$\$
Soft: MySQL or \$\$\$\$\$\$
= \$\$\$\$\$\$\$\$\$\$

Hardware: Add ANY computer
Soft: MongoDB
= Automatic sharding

Empiric comparison: SQL vs noSQL



MySQL configuration



Physical computer:

16GB RAM

8 cores (2x Xeon L5520 @ 2.27GHz)

2TB, 7200 rpm

Indexing:

id

user.id

coordinates.coordinates

created_at

MongoDB minimal configuration

MongoDB “Data”

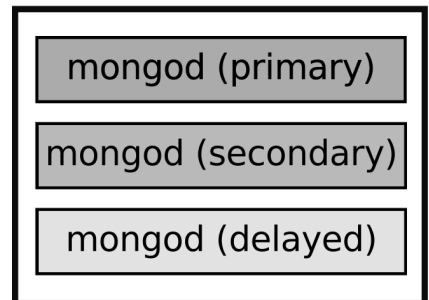
Holds the data

It can be

- A single server with the data
- A replica set with at least 3 servers with data replicated among them for security. One of the elements of the replica set acts as primary and the others are secondaries. Some of the secondary one can replicate the data with a delay time so that they can be used a back up.

Client / Application server

routes the reads and writes from applications to the data

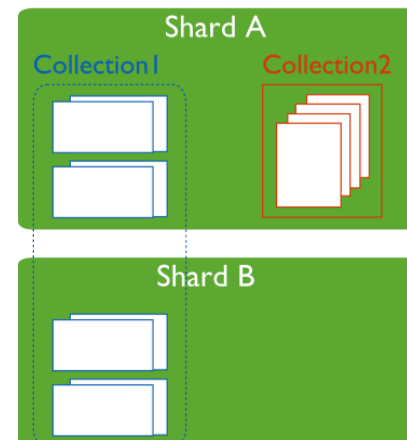


MongoDB sharding

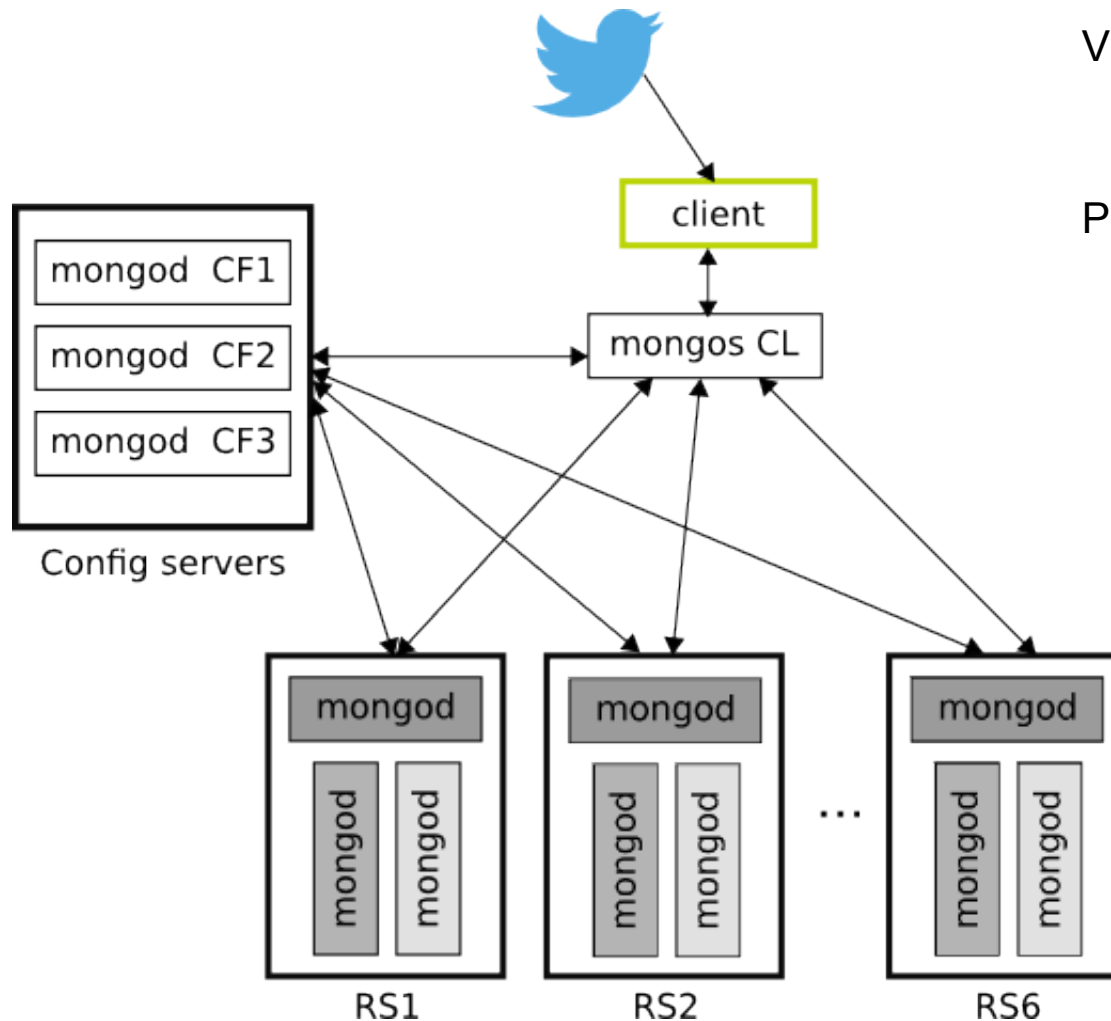
Several MongoDB shards (each formed by a replica set with at least three servers)
Data is distributed among the instances (sharding).

MongoDB configuration servers
distribute the data among the MongoDB instances (chunks)
hold the metadata of the cluster (relations between shard and chunks)

Client / Application server
routes the reads and writes from applications to the shards



MongoDB configuration



Virtual computers:

1-2GB RAM

1-2 cores

Physical computers:

16GB RAM

8 cores (2x Xeon L5520 @ 2.27GHz)

2TB, 7200 rpm

Indexing:

id

user.id

coordinates.coordinates (2d)

created_at

Data collection

Twitter APIs.

- StreamAPI

Client receives a small sample of all public statuses (aprox 1%)

Only 12% of them are geolocalized

and of those, only a small fraction are in the cities of interest.

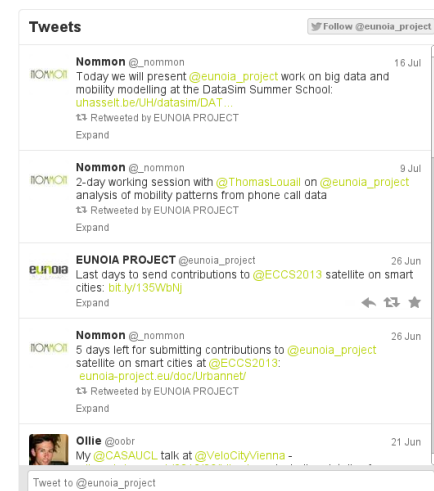
- REST API

Allows to get a users' timeline.

We use that for the most active users located in the cities of interest.

After one year of data collection, in London there are 240K users and 6M tweets

We use **tweepy** to interface with the Twitter APIs
simplejson to convert from/to text/JSON
pymongo to insert the data to MongoDB database
django-orm to insert the data to MySQL database



Data collection

Stream

```
from tweepy import Stream, OAuthHandler
from tweepy.streaming import StreamListener
```

```
class StdOutListener(StreamListener):
    def on_data(self, data):
        print data
        return True
```

```
auth = OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
auth.set_access_token(ACCESS_KEY, ACCESS_SECRET)
```

```
listen = StdOutListener()
stream = Stream(auth, listen, gzip=True)
```

```
stream.sample()
```

Data collection

Users timeline

continuously:

- get user ids with tweets geolocated in the cities of interest
(sorted by number of geolocated tweets)
- for the first 43200 uids (more active)
- get timeline (from last stored tweet until now)

Users network

continuously:

- get user ids with tweets geolocated in the cities of interest
(sorted by number of geolocated tweets)
- for all of them
- get current list of following and followers (uids)

Data collection

Write to DB (MySQL with Django-ORM)

```
class Tweet(Model):
    twid = BigIntegerField(primary_key=True,db_index=True)
    place = ForeignKey(Place, null=True)
    text = CharField(max_length=2048, blank=True)
    retweet_count = IntegerField(null=True)
    parent_id = BigIntegerField(null=True)
    source = CharField(max_length=2048)
    coordinates = ForeignKey(BoundingBox, null=True)
    contributors = CharField(max_length=2048, null=True)
    retweeted = BooleanField()
    truncated = BooleanField()
    created_at = DateTimeField(null=True)
    user = ForeignKey(User)
    entities = ForeignKey(Entities, null=True)
    in_reply_to_status_id = BigIntegerField(null=True)
    in_reply_to_user_id = BigIntegerField(null=True)
    in_reply_to_screen_id = BigIntegerField(null=True)
    deleted = BooleanField()
```

```
class Meta:
    app_label = 'twitter'
```

```
for line in tweets_file:
    tweet = fillTweet(line)
    tweet.save()
```

Idem for user, hashtag, coordinates, url, ...

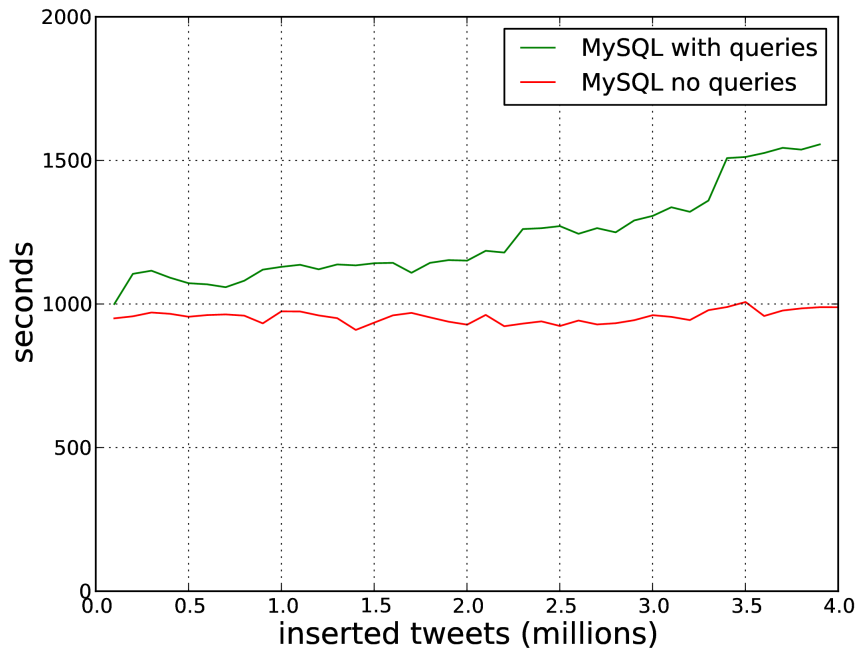
Data collection

Write to DB (MongoDB)

```
mongoserver_uri = "mongodb://" + user + ":" + pwd + "@" +  
                  host + ":" + port + "/" + dbname  
conection = MongoClient(host=mongoserver_uri)  
db = conection[dbname]  
collection = db[collname]  
  
for line in tweetsfile:  
    tweet = simplejson.loads(line.encode('utf8'))  
    collection.insert(tweet)
```

Insertion performance (time to insert 100K tweets)

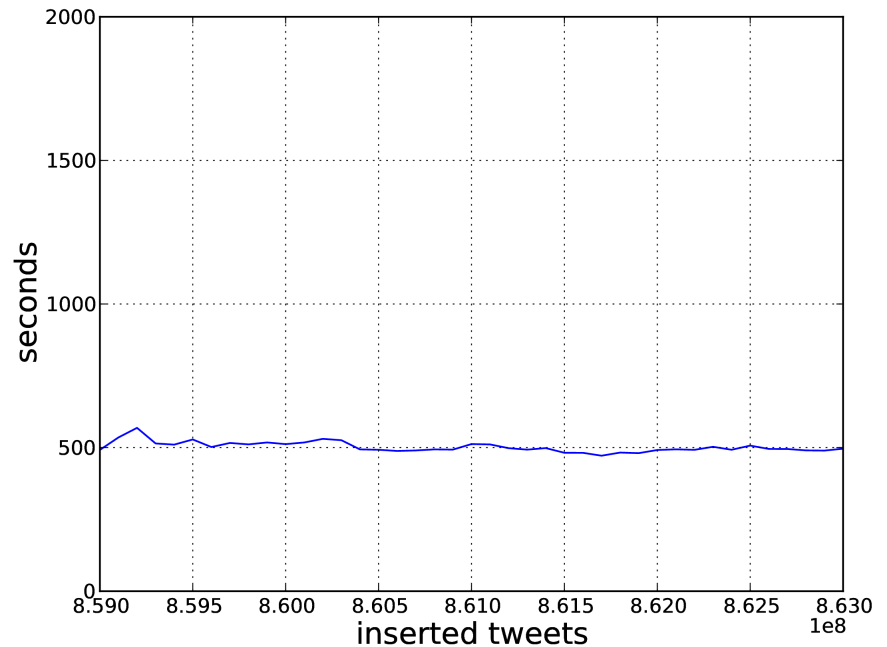
MySQL



empty DB

Django ORM
(+ relations)

MongoDB



850 million tweets
3 replicaset

JSON insertion + pymongo
(200 tweets/second)

Query performance (MongoDB)

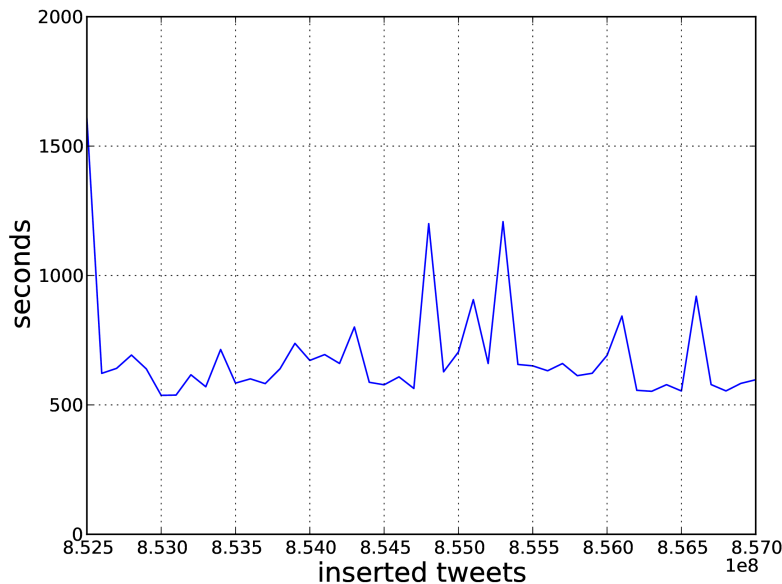
database with 1 thousand million tweets

geonear: get the closest documents to a given point with a maximum distance

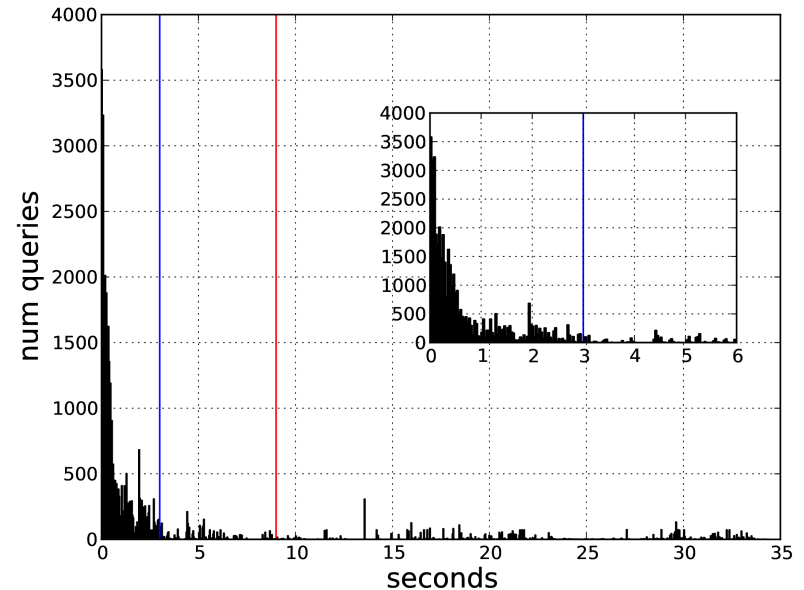
tessellate the city maps with cells of aprox 1km² (16MB size)

run a query per cell

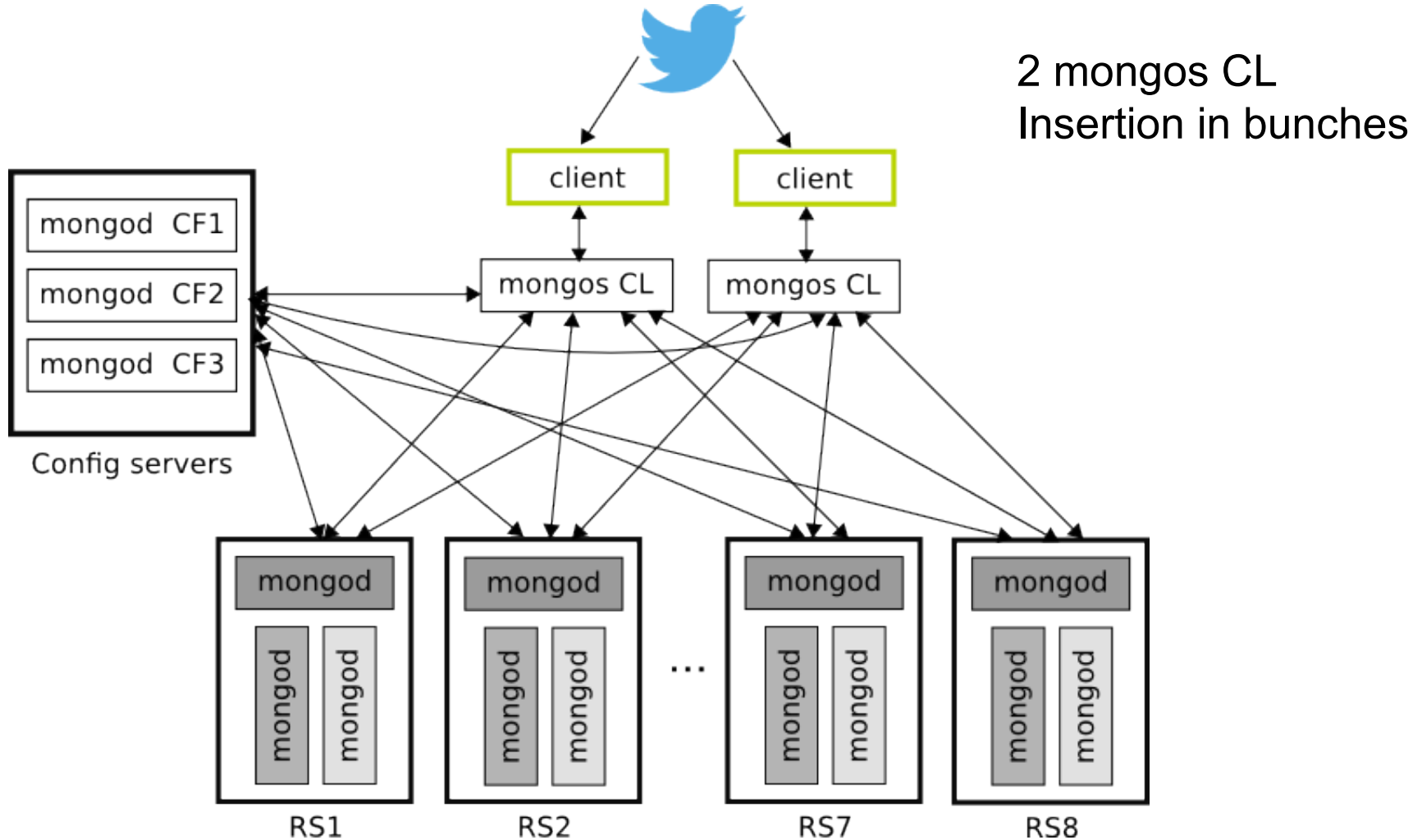
Insertion speed when querying the primary node. (not production)



Queries timing histogram for Barcelona metropolitan area. Blue line shows the median and red line the 70th percentile.



MongoDB configuration improvements



MongoDB issues

- Documents (JSON) and Javascript for queries
- Authorisation/authentication → 2.4
- Default write concern
- Synchronization between primary and secondaries
- Geoindexes: 2d and 2dsphere
- “Big” data backups are not easy → use delayed members in RS

Twitter (and other social network) data issues

- Bias: more tweets in cities
- Signal: people not represented
- Scale: sometimes, small data is better than big data
- Correlations between different subsets
- Big bad data (data quality analysis)
-

Preliminary results and on going work

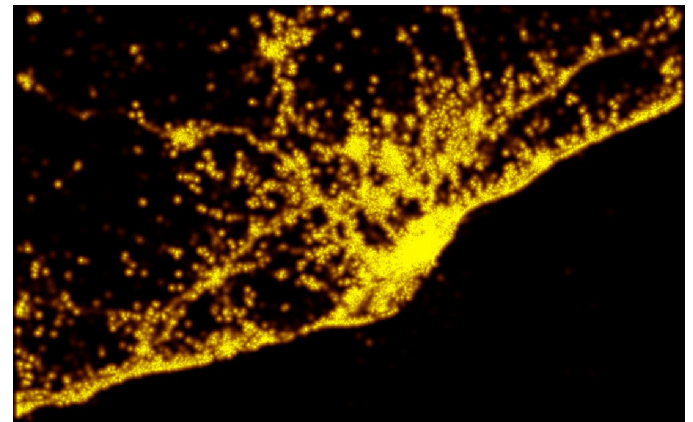
2 billion tweets stored in the database.
250 million geolocated (aprox 13%)

5.9 millions located in London
1.5 millions located in Barcelona
0.2 millions located in Zurich

Users with at least one tweet in the city
260 thousand users in London
25 thousand users in Barcelona
3 thousand users in Zurich

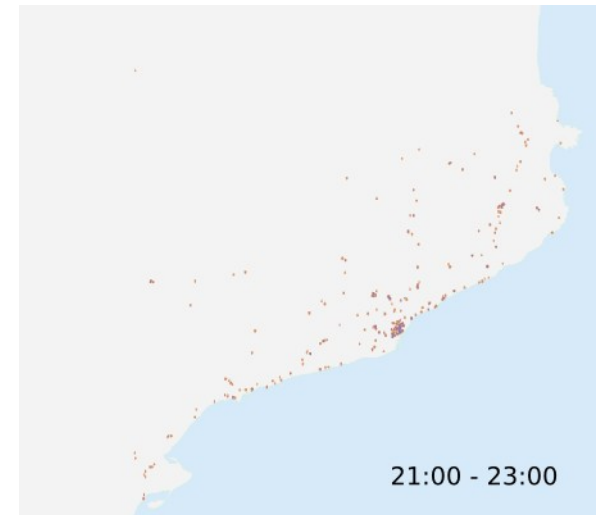
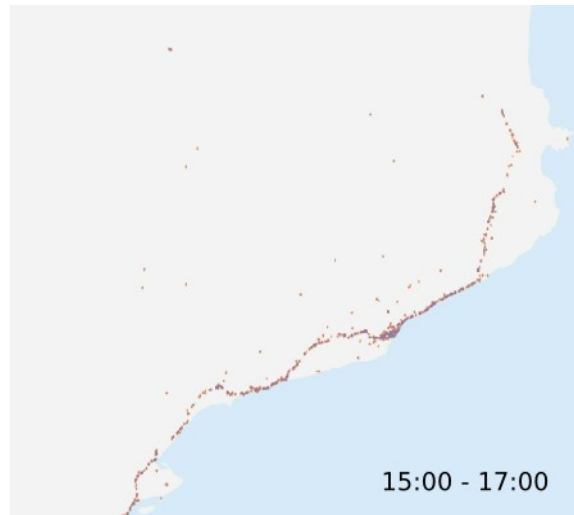
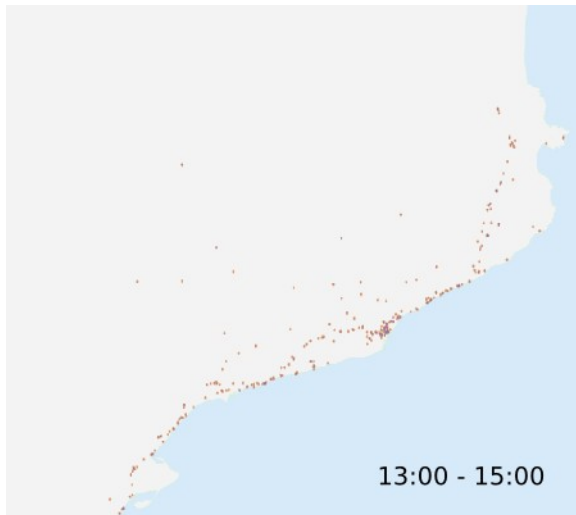
Preliminary data shows a good agreement with
population distribution and the transportation
network of the cities

On going work: more detailed comparison
with traffic data of the transportation networks
with traditional datasets within the European
project EUNOIA



#ViaCatalana

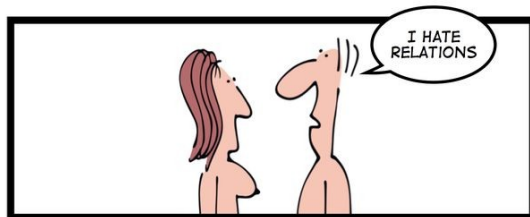
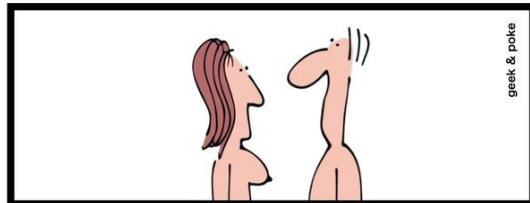
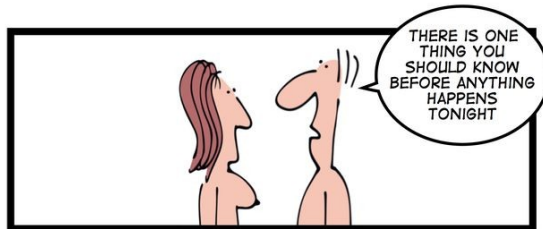
11/09/2013: demonstration in the center of Barcelona and a human chain crossing Catalonia from North to South along the Mediterranean coast.



<http://ifisc.uib-csic.es/humanmobility>

aprox 175000 tweets
4600 geolocated

The Hard Life of a NoSQL Coder



Part 1: The Outing

Thank you

Questions?