# Social media data analysis & real world cases

Nikos Tsirakis

9/12/2017

# Agenda

1. Introduction to News & Social Media, Reputation Management Platforms

2. Data, Data, …. Big Data

3. Analyse, Transform, Enhance, Present

4. Challenges

5. Technologies

6. Case Study

7. Q & A

# Introduction

Brandwatch
Mention
Sysomos
Nuvi
Talkwalker
Agorapulse
Lithium
Synthesio
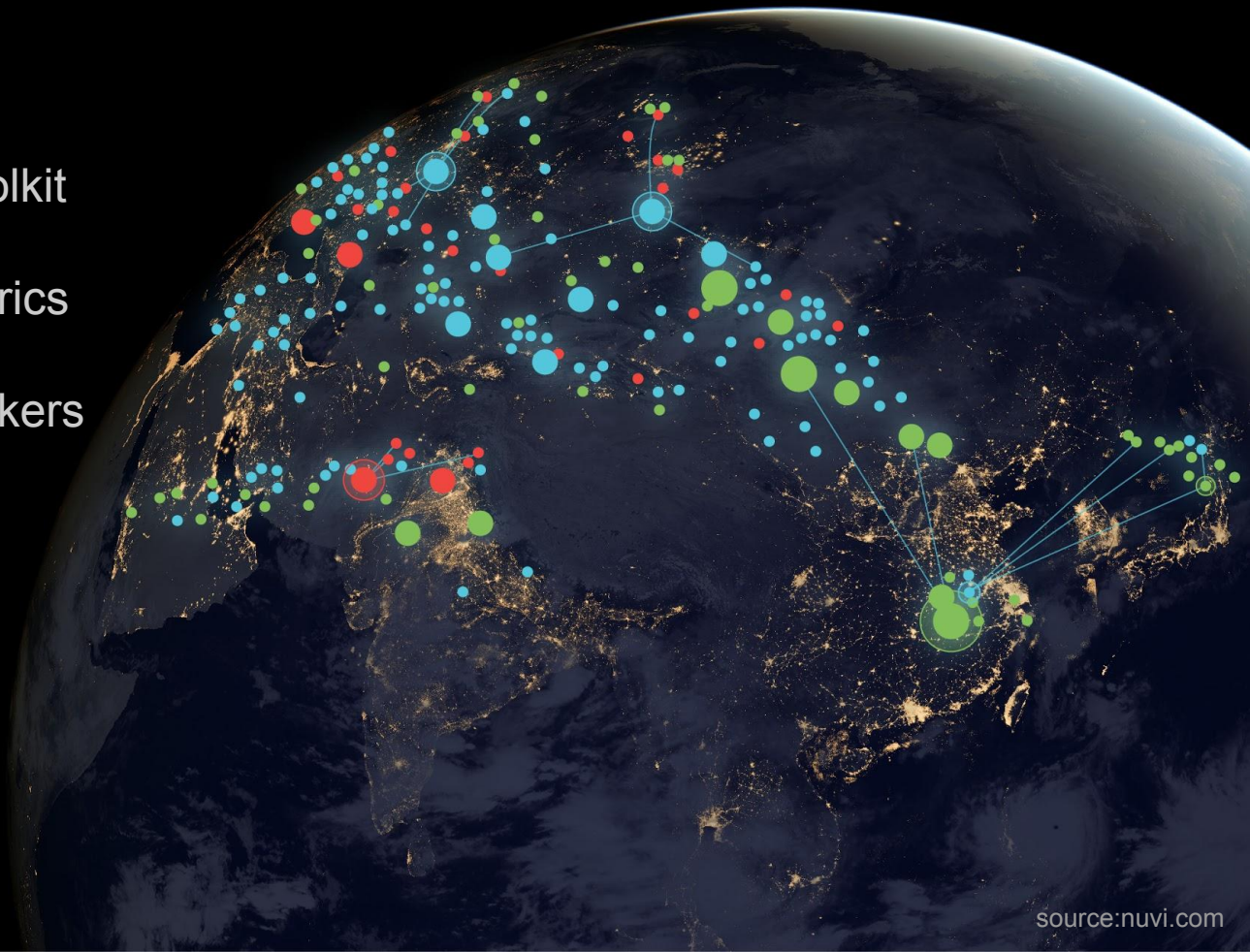Trackur
Keyhole
Buzzlogix
Mentionlytics

Mediatoolkit
Digimind
Ubermetrics
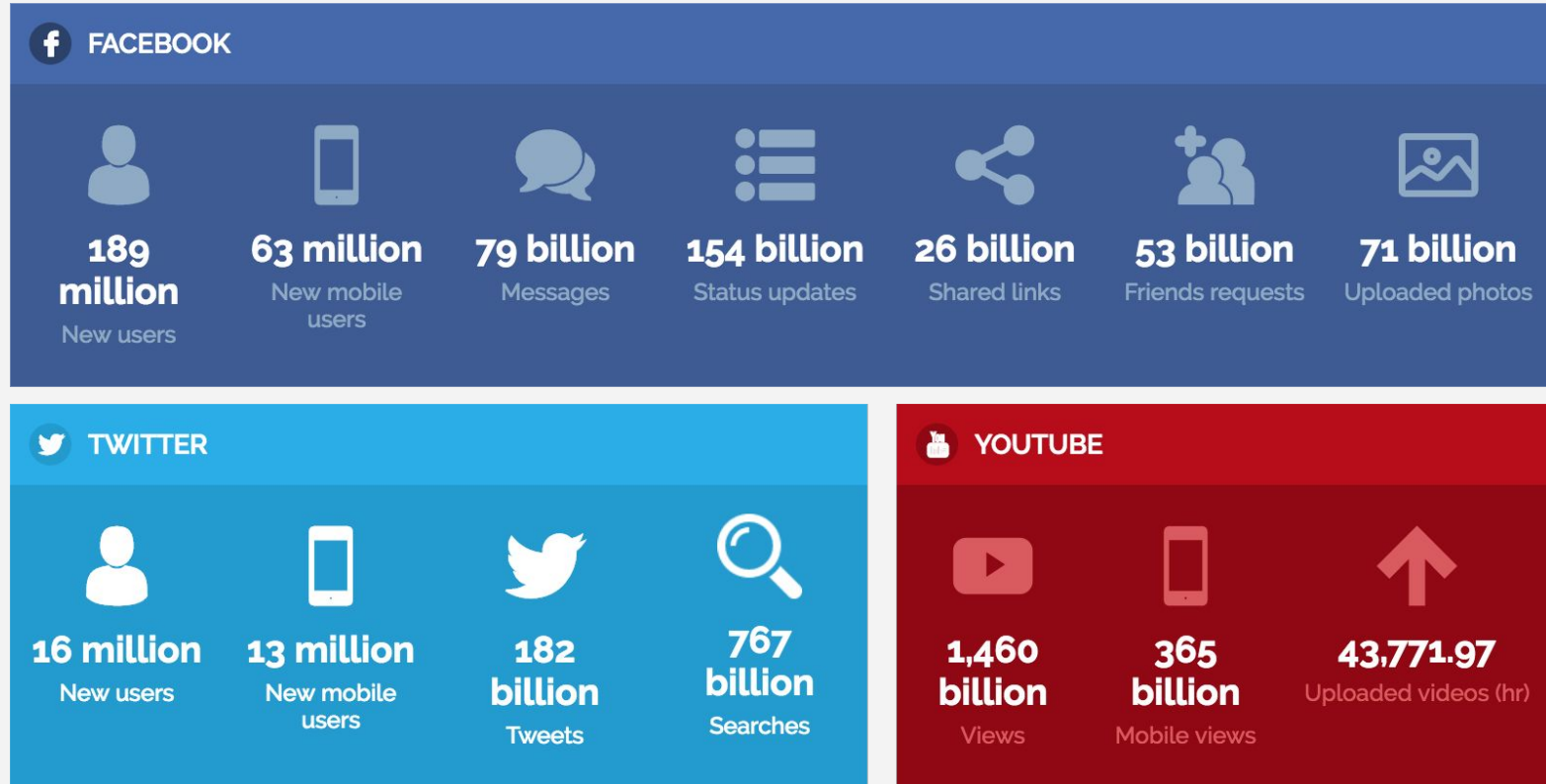Brand24
Socialbakers
Sprinklr
...

source:nuvi.com
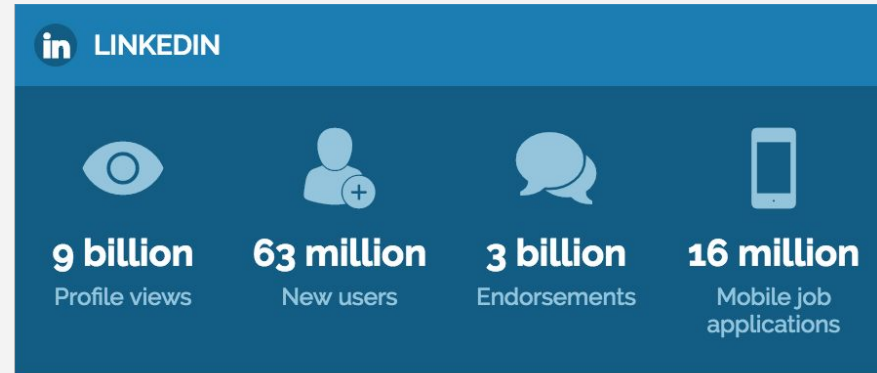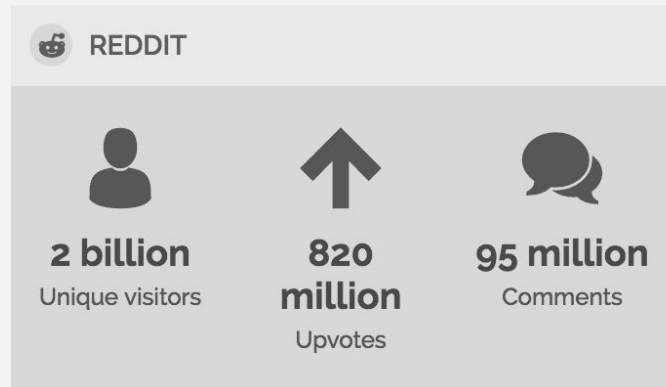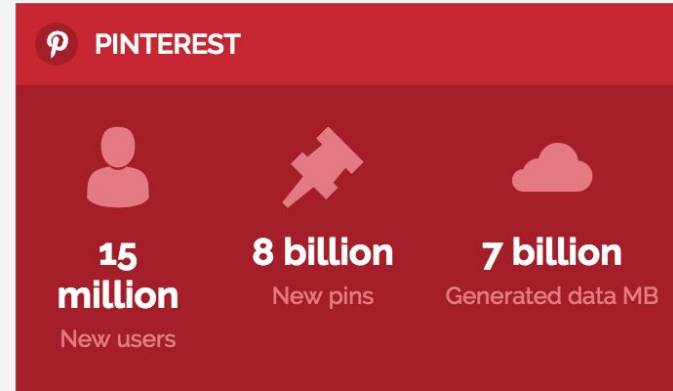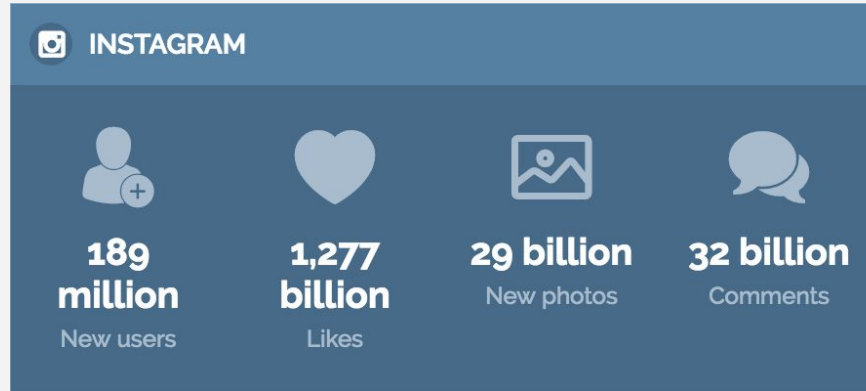
# **GR** Data & Precision in Numbers

- ~20k news sites & blogs

- +20 popular forums

- ~700k twitter users

- ~5m facebook users

- ~2.1m instagram users (mostly between 18-24)

- ~0.5m articles/d (news, blogs, twitter, facebook, youtube, forums)

- ~7.0m internet users

- ~1 million entities

- ~3 million entity terms

- Sentiment:
  - ~87% in sentence level
  - ~70% in entity level
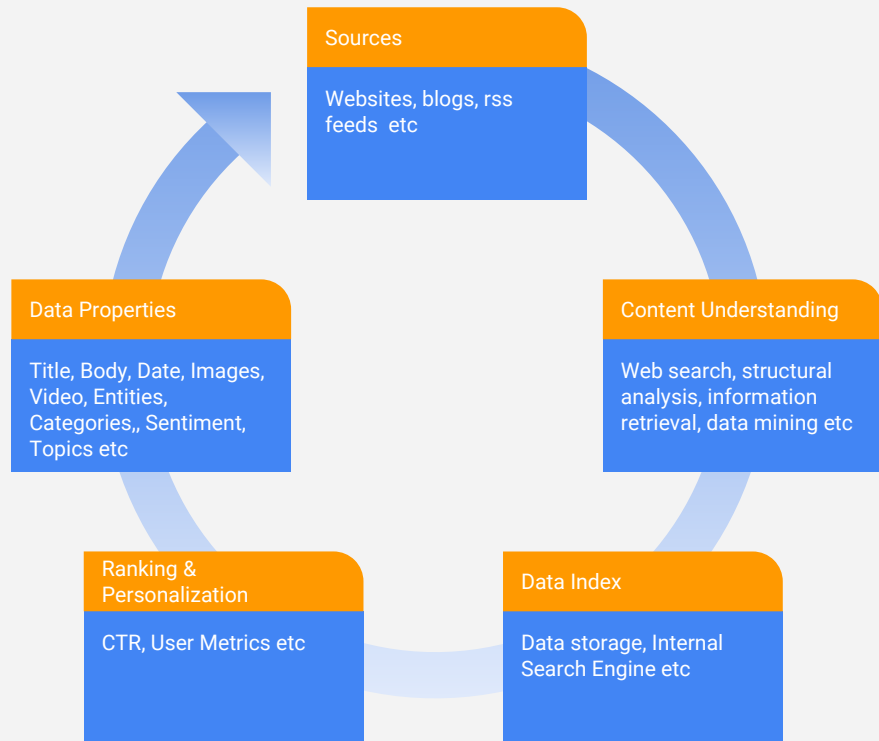
- NER
  - ~70% with 60% recall

# Social Media in numbers (1 year period)

## FACEBOOK

**189 million**
New users

**63 million**
New mobile users

**79 billion**
Messages

**154 billion**
Status updates

**26 billion**
Shared links

**53 billion**
Friends requests

**71 billion**
Uploaded photos

## TWITTER

**16 million**
New users

**13 million**
New mobile users

**182 billion**
Tweets

**767 billion**
Searches

## YOUTUBE

**1,460 billion**
Views

**365 billion**
Mobile views

**43.771.97**
Uploaded videos (hr)

# Social Media in numbers (1 year period)

## INSTAGRAM

**189 million**
New users

**1,277 billion**
Likes

**29 billion**
New photos

**32 billion**
Comments

## PINTEREST

**15 million**
New users

**8 billion**
New pins

**7 billion**
Generated data MB

## REDDIT

**2 billion**
Unique visitors

**820 million**
Upvotes

**95 million**
Comments

## LINKEDIN

**9 billion**
Profile views

**63 million**
New users

**3 billion**
Endorsements

**16 million**
Mobile job applications

source:coupofy.com

# Content Lifecycle



**Sources**

Websites, blogs, rss feeds etc

**Content Understanding**

Web search, structural analysis, information retrieval, data mining etc

**Data Properties**

Title, Body, Date, Images, Video, Entities, Categories,, Sentiment, Topics etc

**Data Index**

Data storage, Internal Search Engine etc
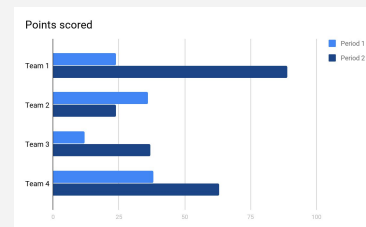
**Ranking & Personalization**

CTR, User Metrics etc

# Data Collection

- Web-page Crawlers (xpath queries, regex)
  - News Detection Content
  - ...
- Blogs APIs
  - Blogger
  - Wordpress
  - Medium
  - ....
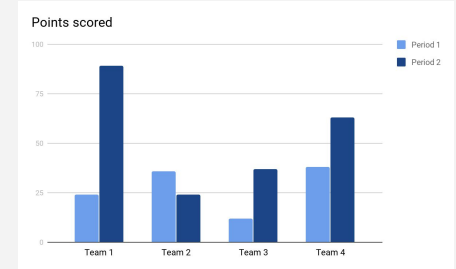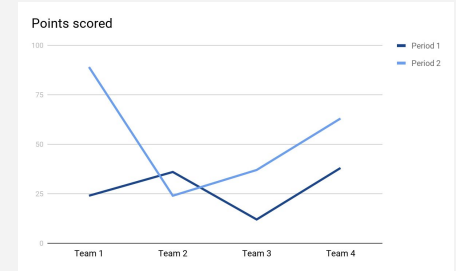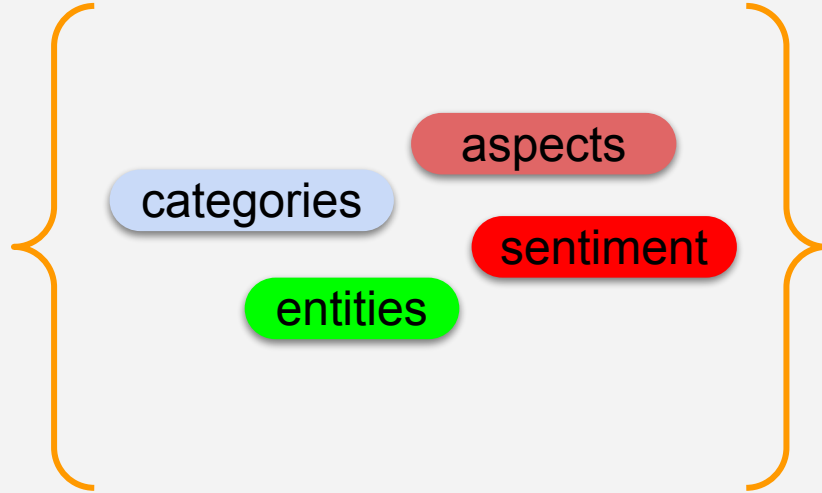- Forum Crawlers
  - PHPBB
  - vBulletin
  - ....

- Social Media Streams
  - Twitter, Facebook ...
- Social Media APIs
  - Twitter, Facebook, Youtube, Instagram ...
- Third Party Services
  - Amazon Alexa
  - Datasift
  - …
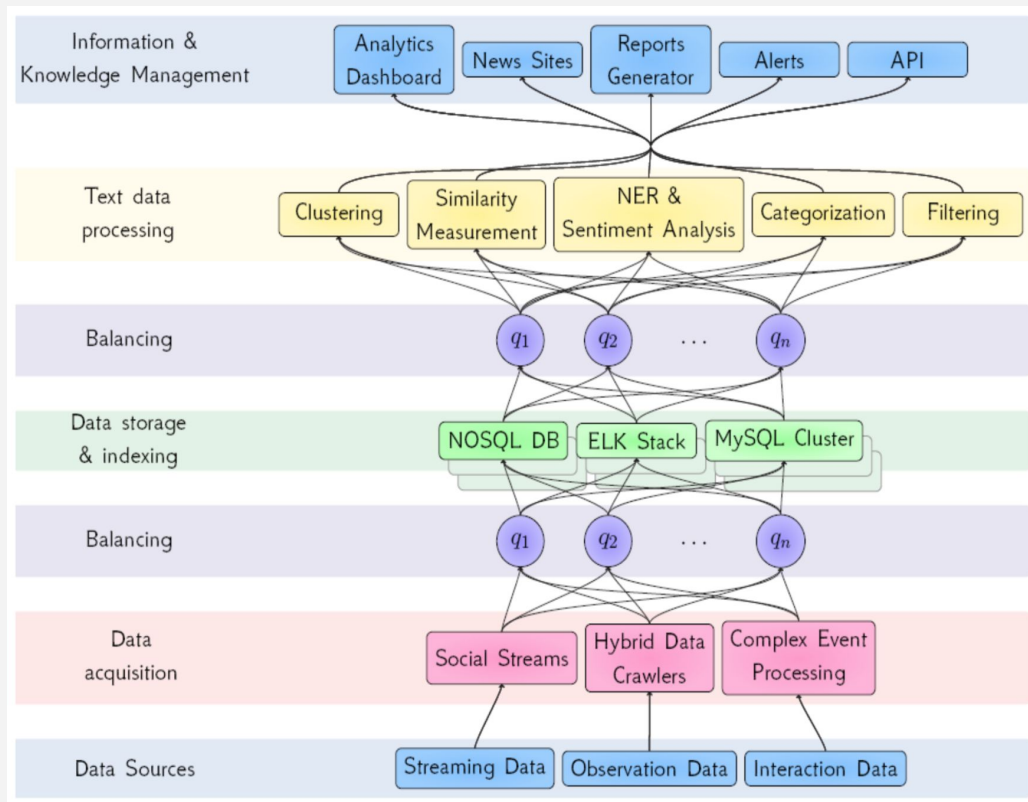- Linked Data
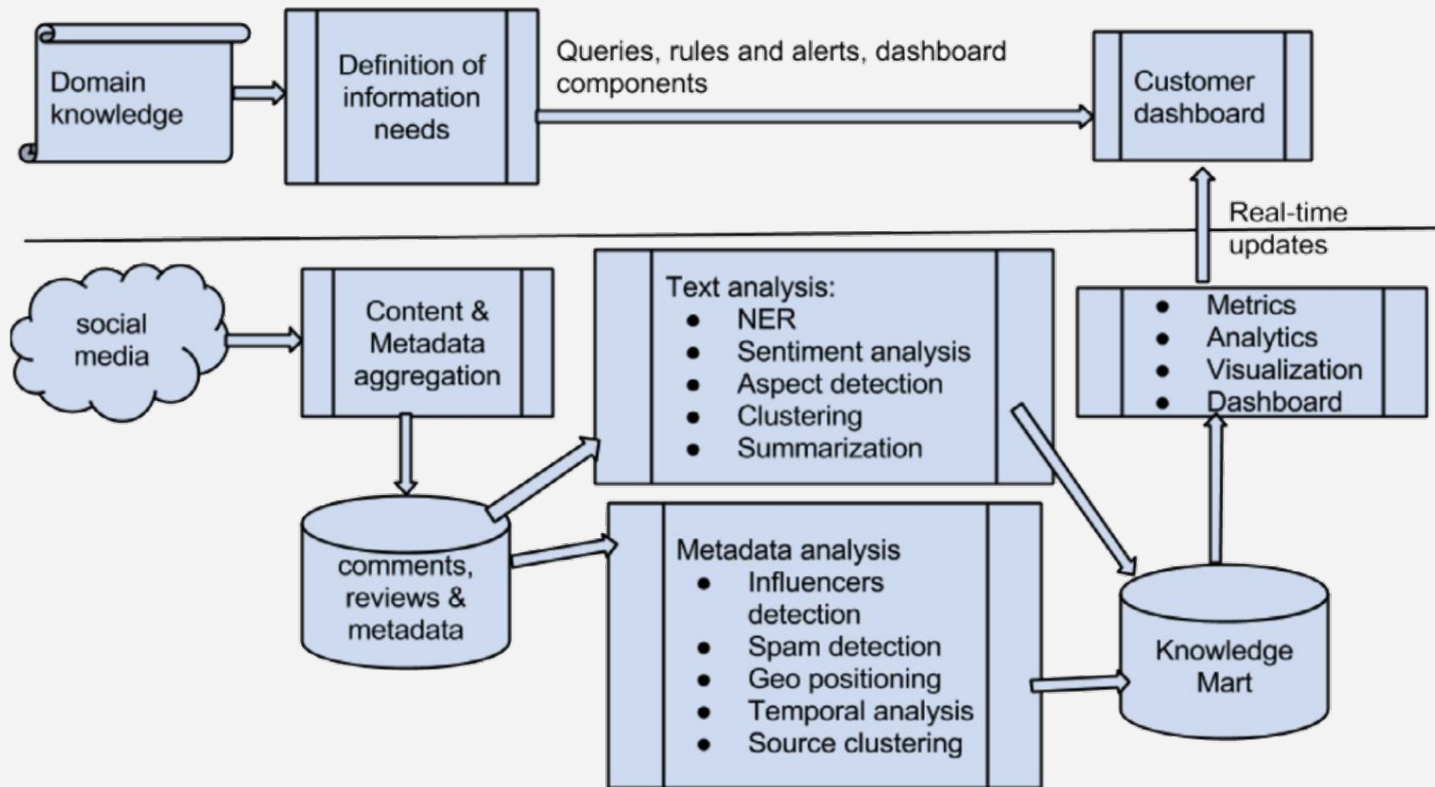  - DBpedia, Google Relation Extraction Corpus

# News Content

categories

hot topics

sentiment

entities

feeds

summaries

clusters

# Social Media Content

# Architecture

# Information Flow

# Challenges (I)



Hot topics

Top News / Breaking News

Clustering

Personalized Popular

# Challenges (II)

- Poor Content (outdated stories, missing media etc)

- Wrong Format (misspelling, wrong media etc)

- Legal Issues (copyright, plagiarism etc)

- Topic evolution (sport games, breaking news etc)

- Cluster Size (big news in short period etc)

# The palopro.io example



Version 1.0

Version 2.0

# Technologies

## Scale-out



**VS**

## Scale-up

# Technologies to scale up

Java & Python threaded services

RDBMS (MySQL)

Search (Solr)

Technologies to scale out

social media → Ingestion microservices

feedback → Kafka Cluster

Spark Cluster

Alerting microservices → Brand

Storage Clusters (ELK, Cassandra, Percona-MySQL…)

Case Study - Hands On

twitter    logstash    elasticsearch    kibana

# Ways to get data

- Logstash
- ~~Elasticsearch-twitter-river~~ (March 2015 - deprecated due to cluster stability)
- Tweepy

# The Logstash Pipeline

kibana

Discover
Visualize
Dashboard
Timelion
Dev Tools
Management

Full screen   Share   Clone   Edit   ❚❚   5 seconds   ❮   ⊙ Today   ❯

*

Uses lucene query syntax

Add a filter ✚

**Total Tweets**

Count

# 10,114

Count

**Geo Count**

Barcelona   Vatican City   Macedonia   Thessaloniki   Istanbul   Georgia   Ossetia
Naples   Greece   Ankara   Bursa   Turkey   Armenia   Azerbaijan   Nagorno-Karabakh Republic
Algiers   Athens   Izmir   Antalya   Tabriz
Malta   Cyprus   Aleppo   Mosul
Tunisia   Syria   Count
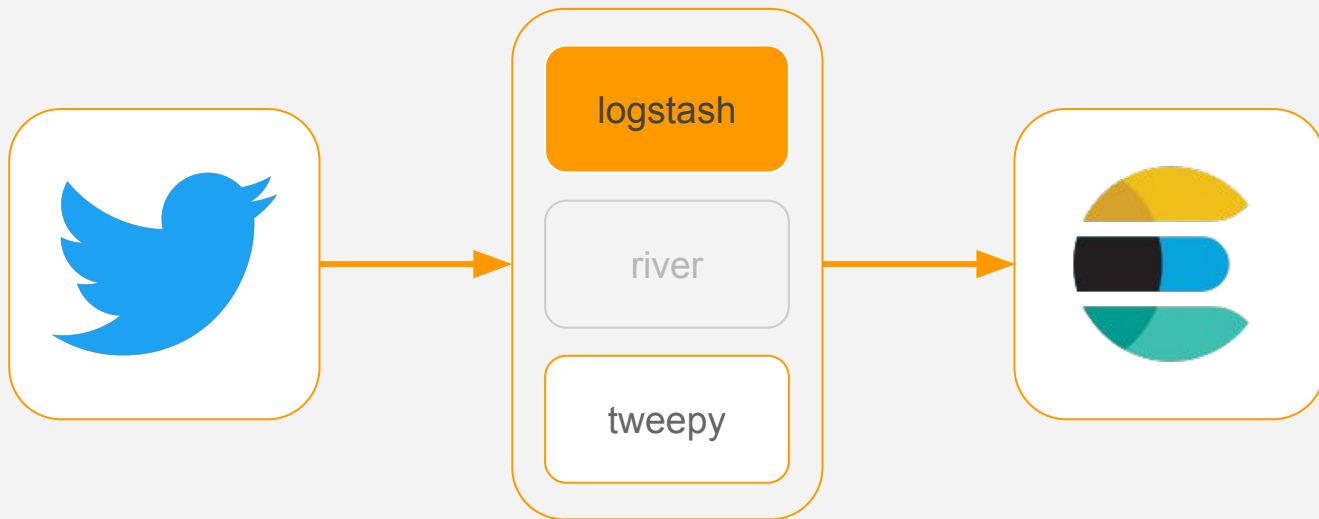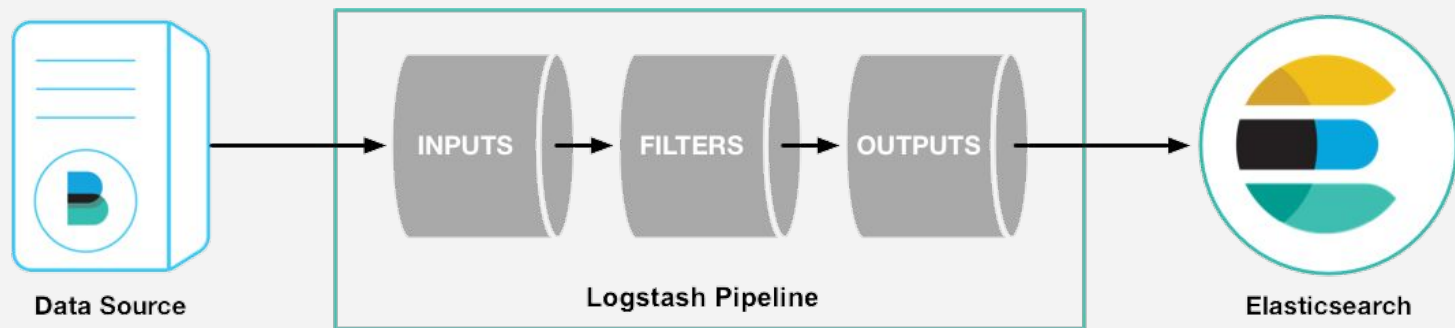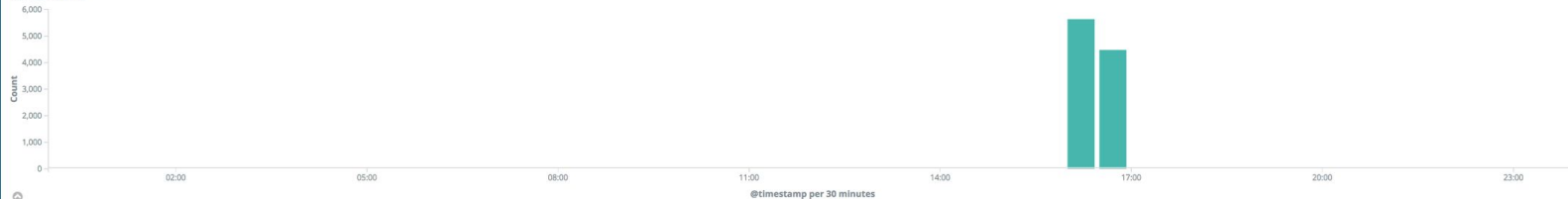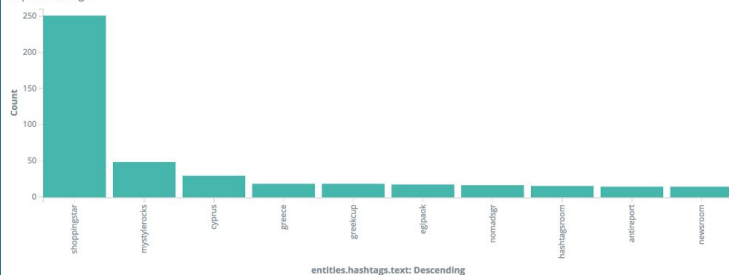Dama   OpenStreetMap contributors, Elastic Maps Service

**Tweets vs. time**

6,000
5,000
4,000
3,000
2,000
1,000
0

Count

02:00   05:00   08:00   11:00   14:00   17:00   20:00   23:00

@timestamp per 30 minutes

● Count

**Top 10 hashtags**

250
200
150
100
50
0

Count

shoppingstar   mystylerocks   cyprus   greece   greekcup   egipaok   nomadsgr   hashtagsroom   antireport   newsroom

entities.hashtags.text: Descending

● Count

**Top 10 influencers (by retweet volume)**

| retweeted_status.user.screen_name: Descending ⇕ | Max retweeted_status.favorite_count ⇕ |
| --- | --- |
| JohnStam13 | 5,382 |
| fayskorda | 4,680 |
| manoskrt | 4,196 |
| vasilis_tag | 3,603 |
| tsipouridhs | 3,198 |
| KaterinaGreece2 | 2,471 |
| vigeko | 2,437 |
| FKatsakis | 2,414 |
| stef_an | 2,284 |
| CheapEripo | 2,128 |

Export: Raw ⬇  Formatted ⬇

**Popular Hashtags**

500
400
300
200

Count

**Popular Hashtags per hour**

500
400
300
200

Count

● shoppingstar
● mystylerocks
● cyprus
● greece
● greekcup
● egipaok
● nomadsgr
● hashtagsroom
● antireport
● newsroom

Collapse

# Questions?

@tsirakis