# Graph Visualization Tool for Twittersphere users based on a high-scalable Extract, Transform and Load System

Pablo Aragón, Íñigo García and Antonio García

May, 27th 2011

CIERZO DEVELOPMENT

smmart
social media marketing
analysis & reporting tool

WIMS'11

# INDEX

**INTRODUCTION**
DISTRIBUTED COMPUTATION
PIPELINE DESIGN
RESULTS

**CIERZO DEVELOPMENT AND SMMART**
STRUCTURE OF TWITTER
VOLUME OF TWITTER
DETECTION OF INFLUENCERS

# INTRODUCTION: CIERZO DEVELOPMENT AND SMMART



SMMART (Social Media Marketing Analysis and Reporting Tool) is the system developed by Cierzo Development for:

- ✓ Corporate social reputation

- ✓ Measuring effectiveness of marketing campaigns

- ✓ Detection of new trends

# INTRODUCTION: STRUCTURE OF TWITTER



Structure of a profile

## INTRODUCTION: STRUCTURE OF TWITTER

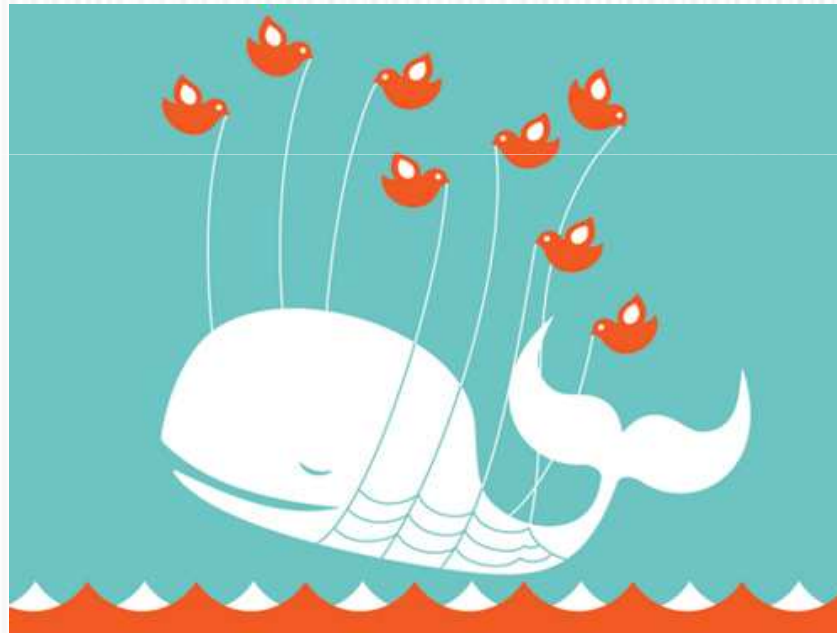A user can set a relationship with another user by:



❑ Reply:        Update that begins with @username

❑ Mention:     Update that contains @username in the body of the tweet

❑ Retweet:     Update that contains the body of another user tweet by specifying the original author

# INTRODUCTION: VOLUME OF THE TWITTER



More than 200M users publishing millions of tweets per day

# INTRODUCTION: DETECTION OF INFLUENCERS



| | | | | | |
|---|---|---|---|---|---|
| #1 | | | Lady Gaga | http://www.ladygaga.com | 9.814.456 |
| #2 | | | Justin Bieber | http://www.youtube.com/justinb | 9.467.962 |
| #3 | | | Barack Obama | http://www.barackobama.com | 7.862.926 |
| #4 | | | Britney Spears | http://www.britneyspears.com | 7.717.271 |
| #5 | | | Kim Kardashian | http://kimkardashian.celebuzz. | 7.403.808 |
| #6 | | | Katy Perry | http://www.katyperry.com | 7.149.480 |
| #7 | | | ashton kutcher | http://www.facebook.com/Ashton | 6.633.714 |
| #8 | | | Ellen DeGeneres | http://www.ellentv.com | 6.577.865 |
| #9 | | | taylorswift13 | http://twitter.com/taylorswift | 6.264.929 |
| #10 | | | Shakira | http://www.shakira.com | 5.858.085 |

Top 10   Top 100   Top 200   Top 300   Top 400   Top 500   Top 600   Top 700   Top 800   Top 900   Top 1000

Old metrics based on data as:

❑ Absolute info:  Number of followers

❑ Relative info:  Quotient of following users and followers

**INTRODUCTION**
DISTRIBUTED COMPUTATION
PIPELINE DESIGN
RESULTS

CIERZO DEVELOPMENT AND SMMART
STRUCTURE OF TWITTER
VOLUME OF TWITTER
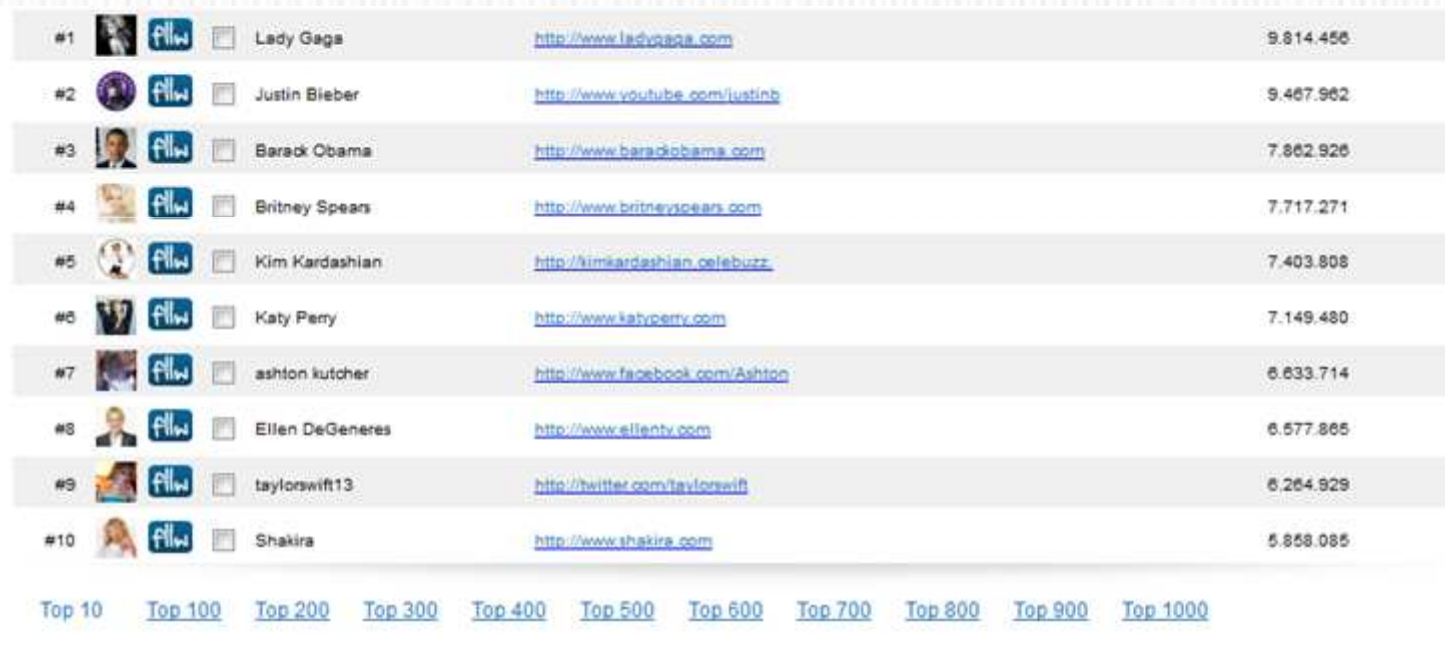**DETECTION OF INFLUENCERS**

# INTRODUCTION: DETECTION OF INFLUENCERS
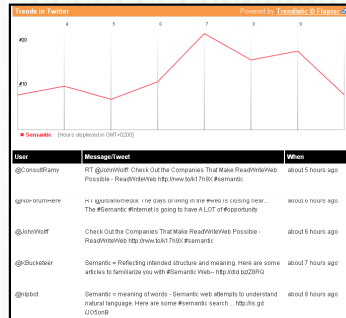


Available search engines track Twitter and list results,

but they do not set a value to the users from the response.

# Spanish voters head to the polls, as city square protests continue

Millions turn out to elect municipal councils and regional governments, despite protests over politics-as-usual

**#spanishrevolution**
**#yeswecamp**
**#15m**

**Giles Tremlett** in Madrid
guardian.co.uk, Sunday 22 May 2011 19.46 BST
Article history

## Opinion

### 'Yes we camp' activists hit Spanish streets

When Spanish officials banned a protest camp in central Madrid ahead of elections, thousands defied the ruling.

Leila Nachawati  Last Modified: 22 May 2011 10:34

Email    Print    Share    Send Feedback        Tweet  258    Like  2K

Listen

Protesters react a
vote in regional el

...ement that protests against the ongoing financial crisis, politicians and bankers sit on the roof a
...bway station as they camp at Madrid's Puerta del Sol May 17, 2011 [Reuters]
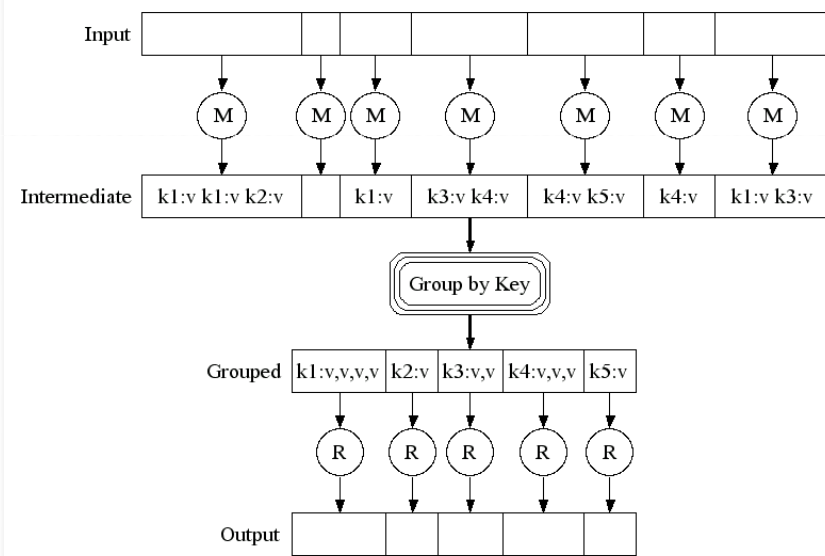
# DISTRIBUTED COMPUTATION

- ❑ Management of large volumes at the lowest cost

- ❑ Automatic adjustment to the daily growth of users and the oscillations in the frequency of publication

INTRODUCTION
**DISTRIBUTED COMPUTATION**
PIPELINE DESIGN
RESULTS

**HADOOP**
AMAZON EC2

# DISTRIBUTED COMPUTATION: HADOOP



Map Reduce



Distributed File System

INTRODUCTION
**DISTRIBUTED COMPUTATION**
PIPELINE DESIGN
RESULTS

HADOOP
**AMAZON EC2**

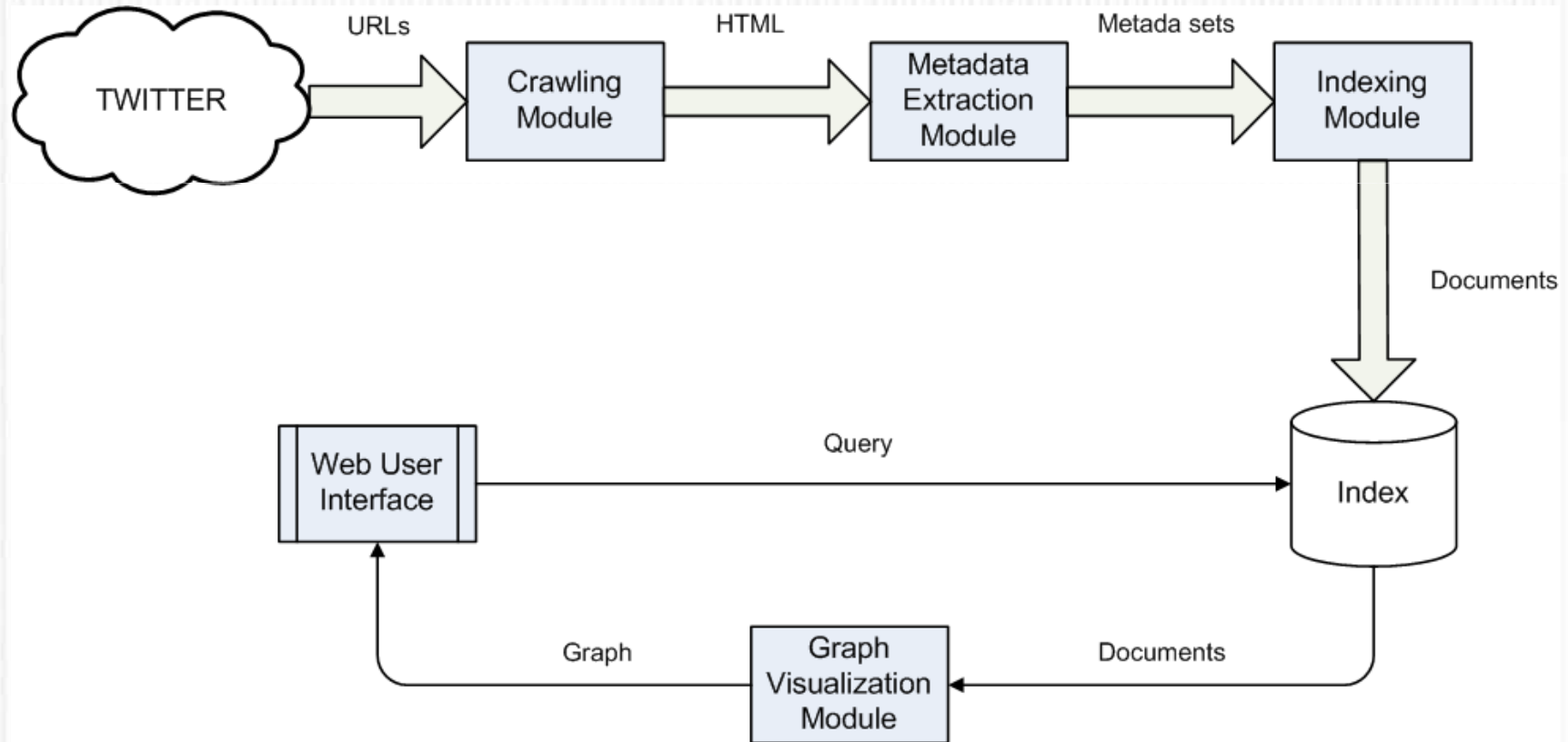# DISTRIBUTED COMPUTATION: AMAZON EC2

Definition of a Hadoop node as a machine image in Amazon Elastic Compute Cloud.

The system balancing mechanism adds and removes Hadoop nodes in real time on demand.

# PIPELINE DESIGN

INTRODUCTION
DISTRIBUTED COMPUTATION
**PIPELINE DESIGN**
RESULTS

**CRAWLING MODULE**
METADATA EXTRACTION MODULE
INDEXING MODULE
GRAPH VISUALIZATION MODULE

# PIPELINE DESIGN: CRAWLING MODULE

## Based on Nutch

1. Crawl the Twitter profiles stored in a DB

2. Extract outlinks to new profiles

INTRODUCTION
DISTRIBUTED COMPUTATION
**PIPELINE DESIGN**
RESULTS

CRAWLING MODULE
**METADATA EXTRACTION MODULE**
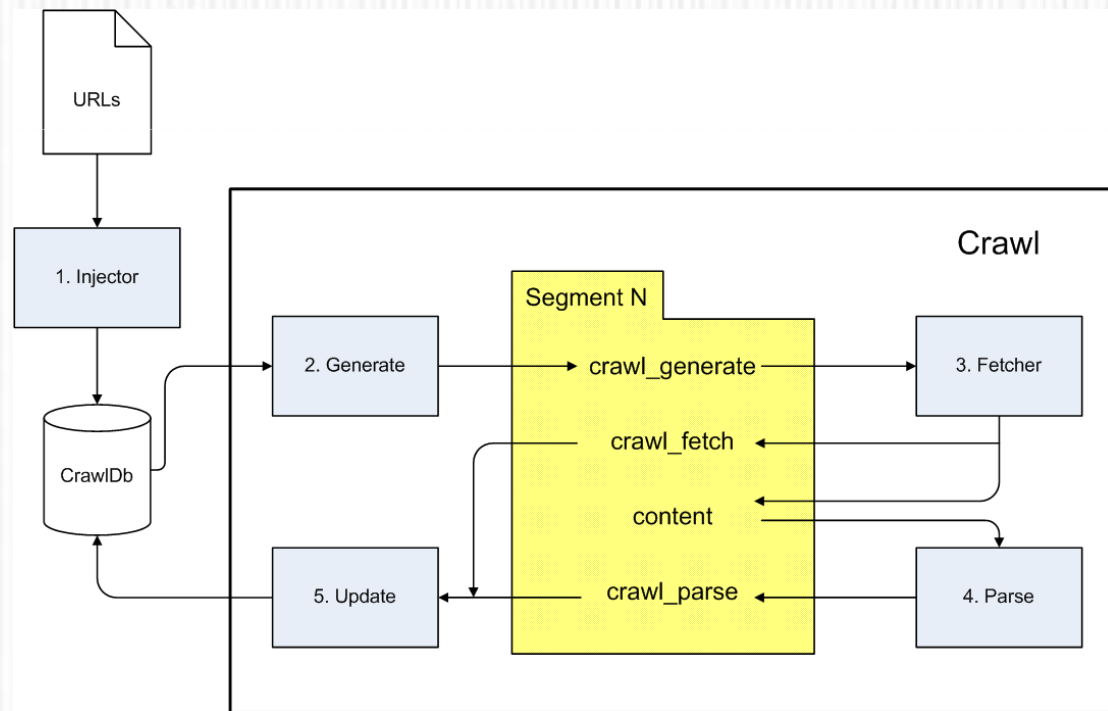INDEXING MODULE
GRAPH VISUALIZATION MODULE

# PIPELINE DESIGN: METADATA EXTRACTION MODULE

The portion of HTML of a tweet contains a set of metadata:

- ❑ Textual content

- ❑ Publication date

- ❑ Author

- ❑ Mention to other users

INTRODUCTION
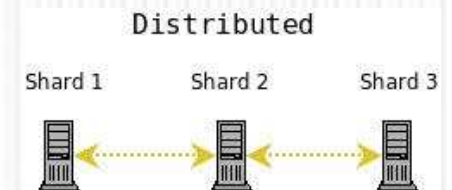DISTRIBUTED COMPUTATION
**PIPELINE DESIGN**
RESULTS

CRAWLING MODULE
METADATA EXTRACTION MODULE
**INDEXING MODULE**
GRAPH VISUALIZATION MODULE

# PIPELINE DESIGN: INDEXING MODULE

Apache Solr (enterprise search server based on Lucene)

✓ Sorting algorithms

✓ Stemming

✓ Stopwords filters

✓ Faceted searchs

Multicore architecture sharding by publication date.

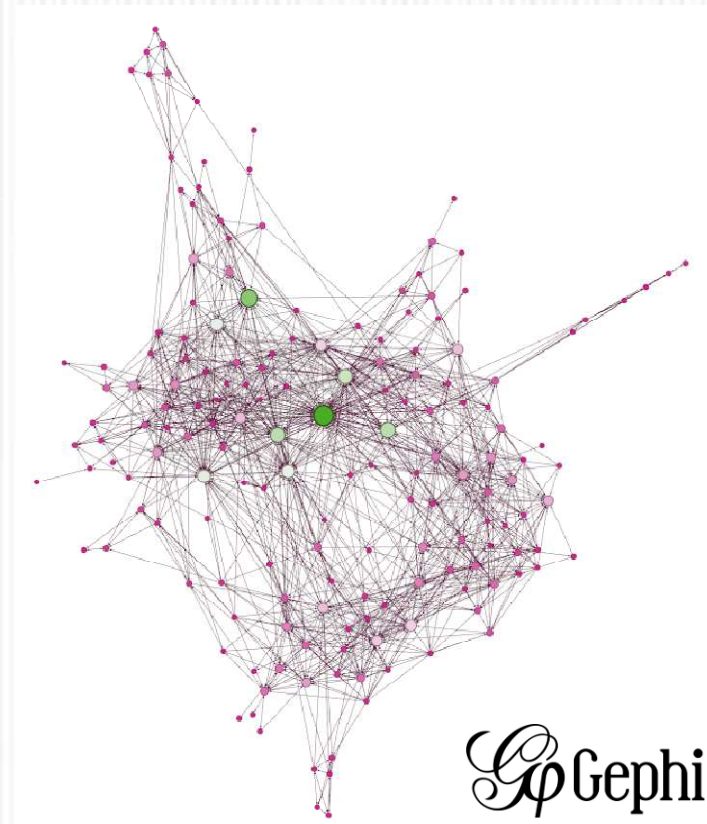INTRODUCTION
DISTRIBUTED COMPUTATION
**PIPELINE DESIGN**
RESULTS

CRAWLING MODULE
METADATA EXTRACTION MODULE
INDEXING MODULE
**GRAPH VISUALIZATION MODULE**

# PIPELINE DESIGN: GRAPH VISUALIZATION MODULE

The Graph Visualization module transforms the responses from the index into a graph by the force-based multilevel layout Yifan Hu's algorithm provided in Gephi Toolkit.
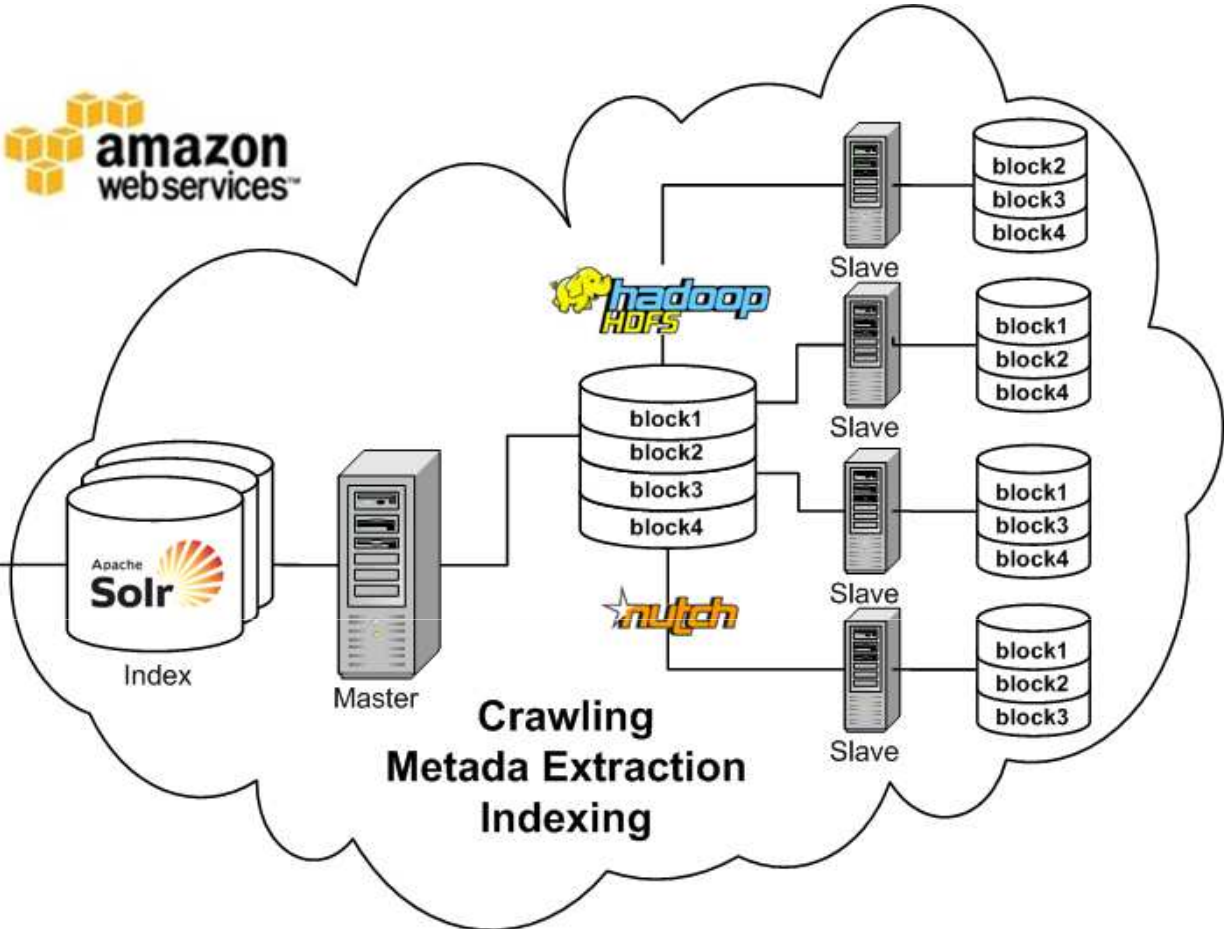
INTRODUCTION     **WESTERN SAHARA CONFLICT**
DISTRIBUTED COMPUTATION     PATXI LÓPEZ
PIPELINE DESIGN     CONCLUSIONS
**RESULTS**     FUTURE WORK

# RESULTS: WESTERN SAHARA CONFLICT
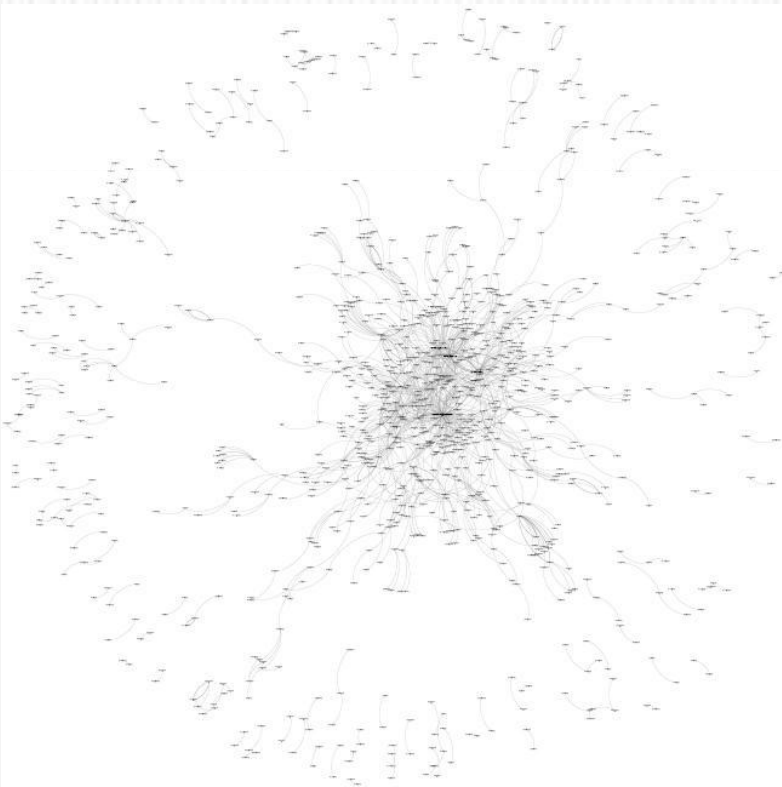


In November 2010, Moroccan security forces involved in a camp in Western Sahara. This action was criticized by part of the Spanish society.

# RESULTS: WESTERN SAHARA CONFLICT



## Search

- ❑ content:'sahara'
- ❑ language:'es'
- ❑ date:[2010-11-10 TO 2010-11-18]

## Results

- ✓ 1721 users
- ✓ 3925 tweets
- ✓ 707 mentions

INTRODUCTION
DISTRIBUTED COMPUTATION
PIPELINE DESIGN
**RESULTS**

**WESTERN SAHARA CONFLICT**
PATXI LÓPEZ
CONCLUSIONS
FUTURE WORK

# RESULTS: WESTERN SAHARA CONFLICT

INTRODUCTION
DISTRIBUTED COMPUTATION
PIPELINE DESIGN
**RESULTS**

WESTERN SAHARA CONFLICT
**PATXI LÓPEZ**
CONCLUSIONS
FUTURE WORK

# RESULTS: PATXI LÓPEZ
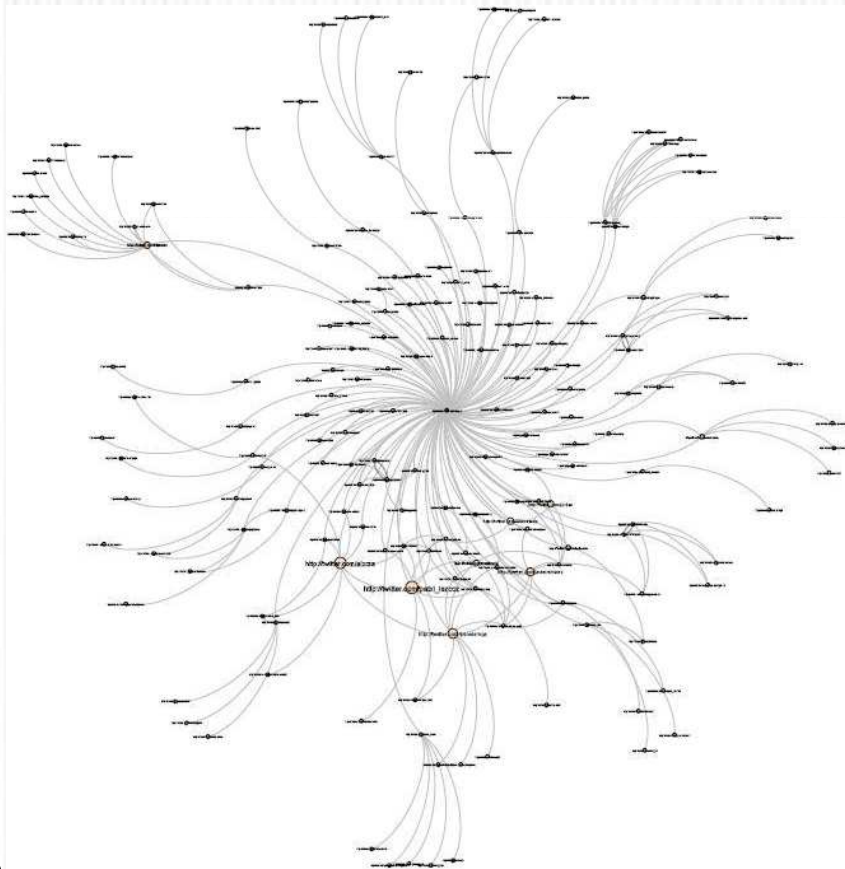


Patxi López holds the position of the President of the Basque Country Government. His campaign included strategies in social networks.

# RESULTS: PATXI LÓPEZ



## Search

❑ mention:'patxi_lopez'

❑ language:'es'
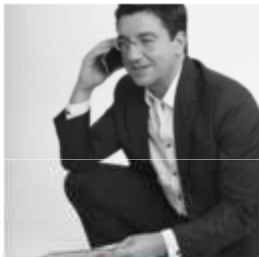
❑ date:[2010-11-10 TO 2010-11-18]

## Results

✓ 186 users

✓ 196 tweets

✓ 366 mentions

INTRODUCTION
DISTRIBUTED COMPUTATION
PIPELINE DESIGN
**RESULTS**

WESTERN SAHARA CONFLICT
**PATXI LÓPEZ**
CONCLUSIONS
FUTURE WORK

# RESULTS: PATXI LÓPEZ

## RESULTS: CONCLUSIONS

- ✓ The implemented tool identifies main influencers in a specific topic or around a concrete user

- ✓ The high-scalable design adapts to a large social network as Twitter

- ✓ Enterprises can deploy social media monitoring systems using exclusively open source technologies

- ✓ The tool provides information for performing crisis management

| Tag | Evolución |
|---|---|
| SYKES | 13.7 |
| estabilidad | 10 |
| GSS | 4.9 |
| UNISONO | 4.7 |
| COHEREN* | 2.1 |
| ERROR | 2 |
| APLAUD* | 0.4 |
| SERTEL | 0.3 |

Actualizar 22-05-201

# RESULTS: FUTURE WORK

- ✓ New versions for more social media sources

- ✓ Real-time results

- ✓ New data mining applications

  - ❑ Predictive models

*Thanks for your attention*