# Monitoring the Socio-Economic Climate at a Global Scale

## Jorge Louçã

### Lisbon University Institute

Seminar "Modeling Complex Socio-Economic Systems and Crises", organized by the ETH Zurich CCSS - Center of Competence for Coping with Crises in Socio Economic Systems

8 March 2011

## I - The Observatorium

- main goal
- collecting vs. analyzing data

## II - On-going research

- topic detection in on-line newspapers
- dynamics of scientific trends
- political opinion dynamics
- socio-economic "climate" at a global scale

## III - Open source library of tools and free data for computational social science
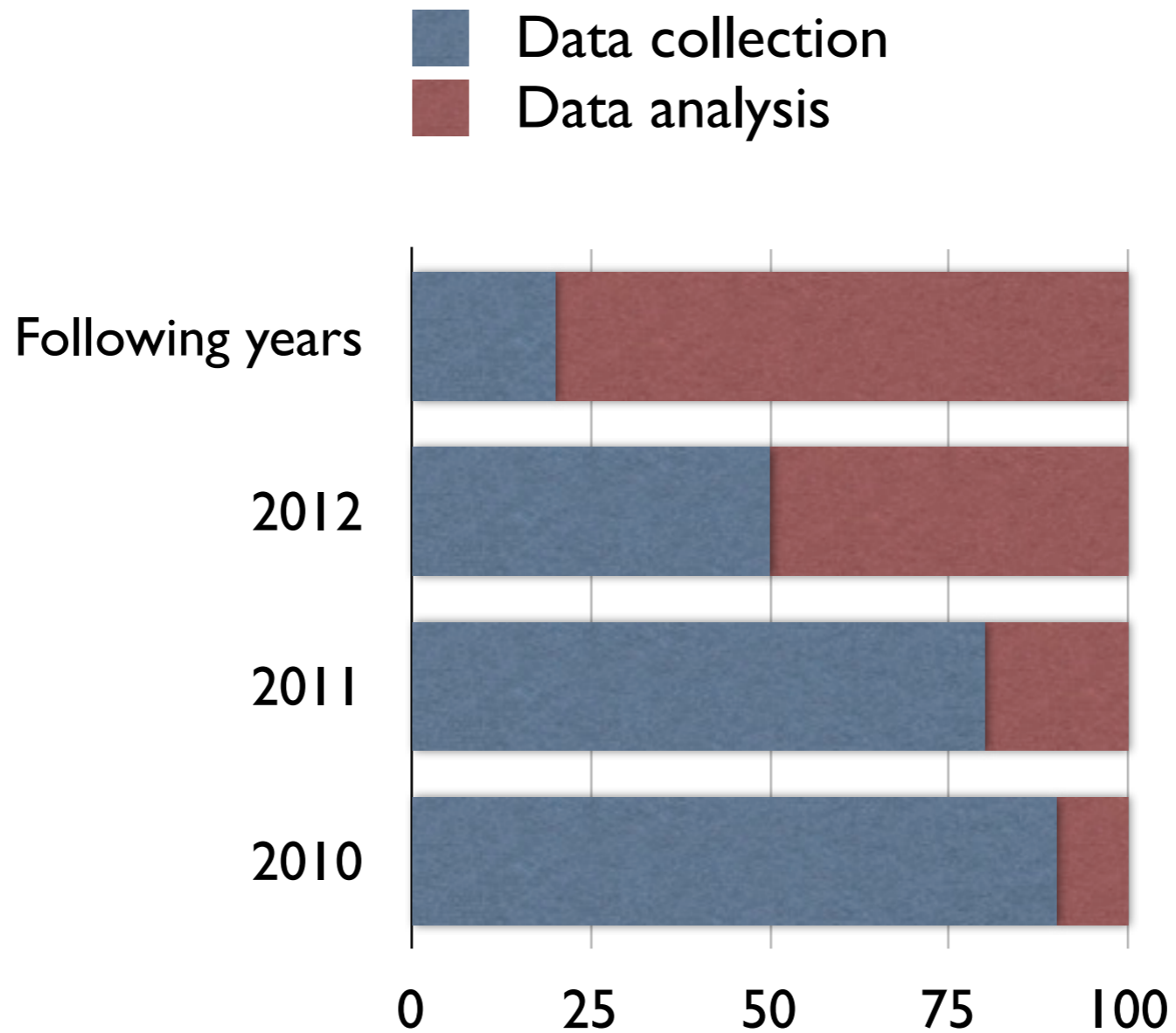
"Intelligence is ten million rules."

Douglas Lenat

Goal:

>  understand more about how and what
>  are people discussing - what are their
>  interests, topics, opinions, arguments,
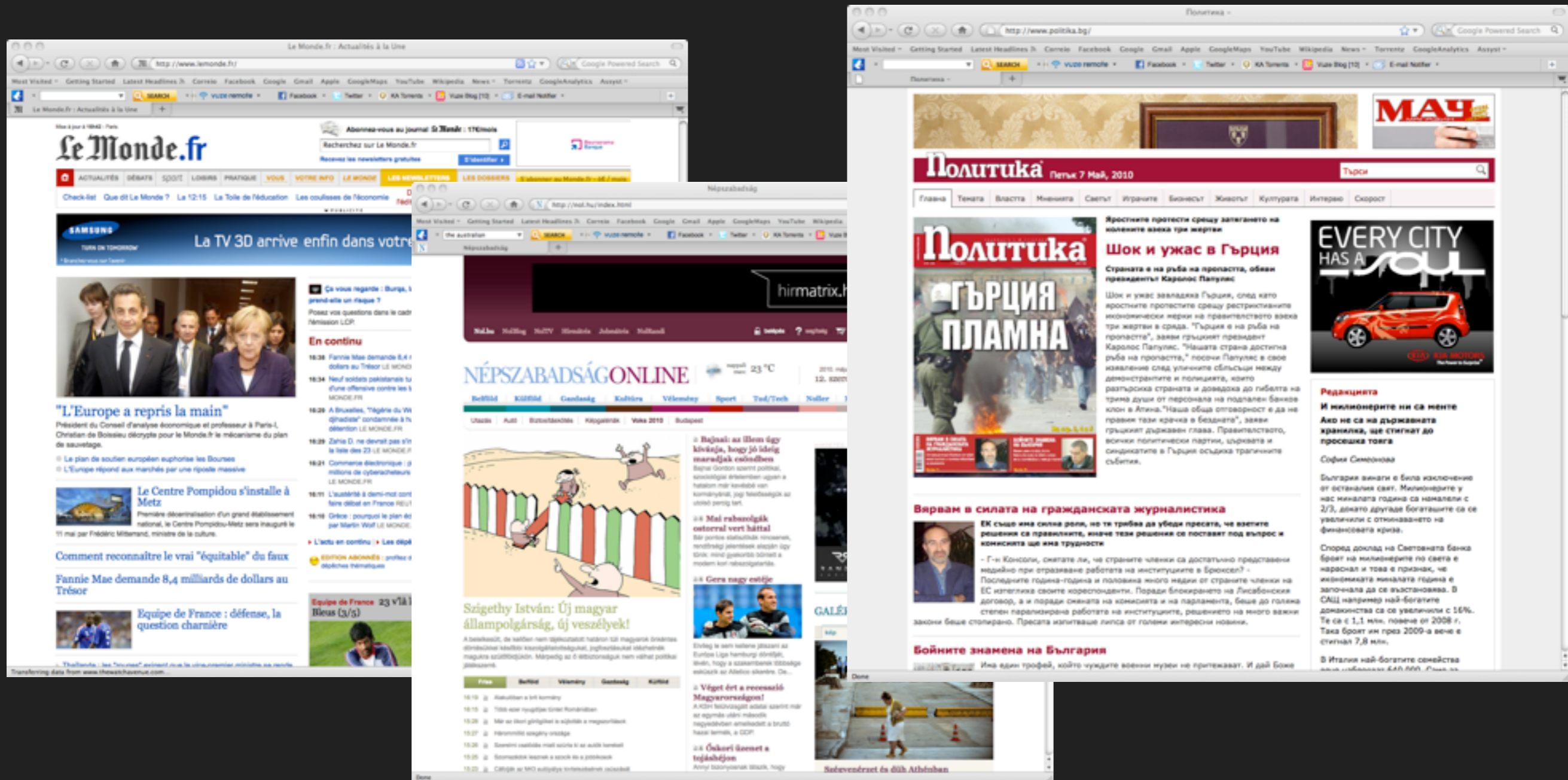>  structure of communication

>  at a global level...

Approach:

>  data-driven

On Feb 17th: 27,9GB (of text) corresponding to 424.850 news

Frequent approaches: automatic categorisation of text from linguistics, natural language, statistics and information retrieval.  Strategies:

- regression models

- Bayesian approaches

- nearest neighbor classification

- neural networks

- hierarchical clustering          (Miao and Qiu, 2010) (Sole et al., 2010)

Frequent approaches: automatic categorisation of text from linguistics, natural language, statistics and information retrieval.  Strategies:

- regression models

- Bayesian approaches

- nearest neighbor classification

- neural networks

- hierarchical clustering          (Miao and Qiu, 2010) (Sole et al., 2010)

**Supervision is needed !**

A non-semantic methodology for extraction and real-time analysis of documents from the Internet

A non-semantic methodology for extraction and real-time analysis of documents from the Internet

Extraction of the text from each HTML file by employing the Text to Tag ratio (Weninger, 2008)

http://www.politika.bg/article?id=17018

## Extraction of the text from each HTML file by employing the Text to Tag ratio (Weninger, 2008)



http://www.politika.bg/article?id=17018

Extraction of the text from each HTML file by employing the Text to Tag ratio (Weninger, 2008)

## Extraction of the text from each HTML file by employing the Text to Tag ratio (Weninger, 2008)

Extraction of the text from each HTML file by employing the Text to Tag ratio (Weninger, 2008)
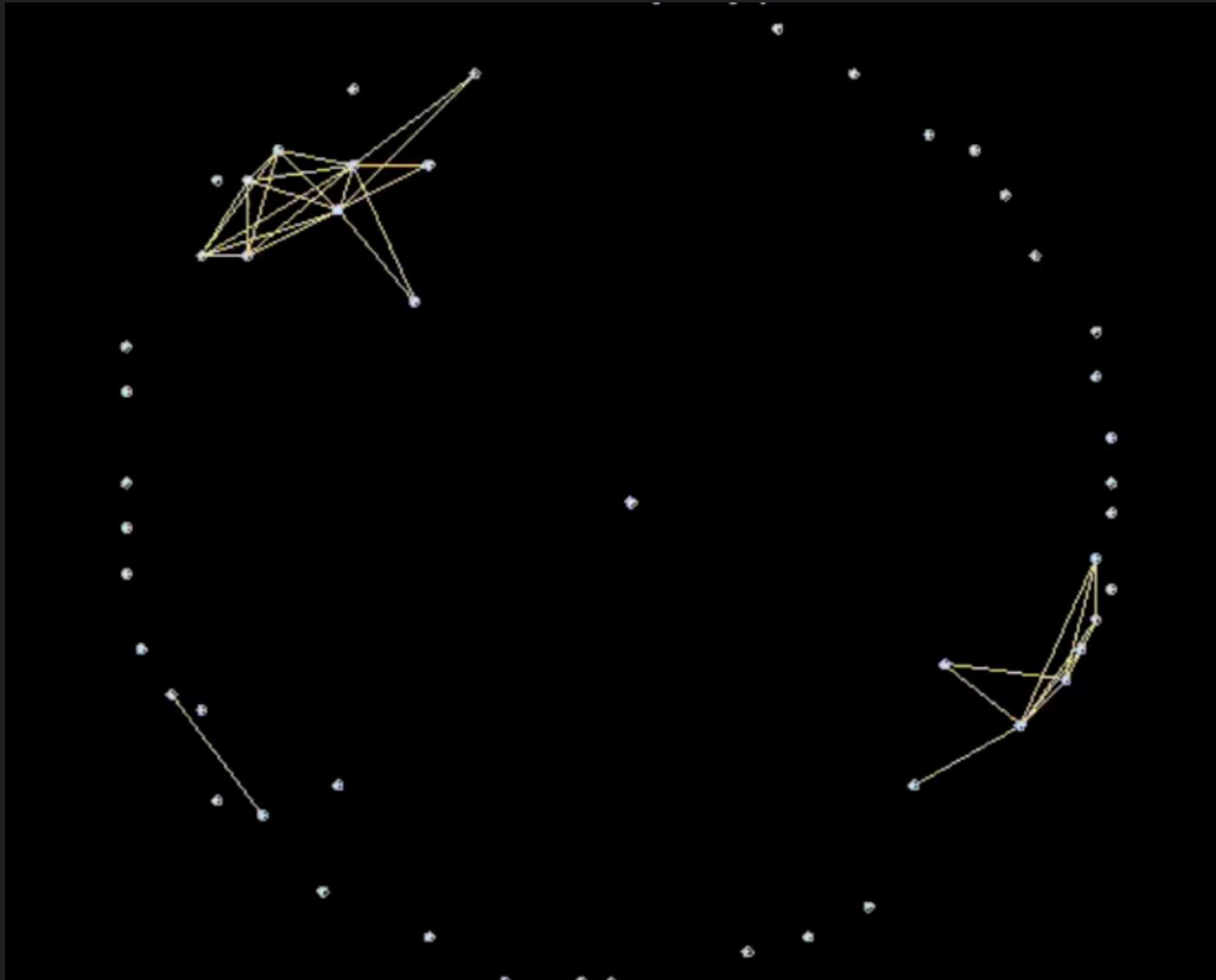
Graph composition:

- Each news has a set of words, represented by a node

- Each node is added to a graph G

- A node is linked to other nodes according to a distance

- The Jaccard similarity coefficient is used to calculate distances

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

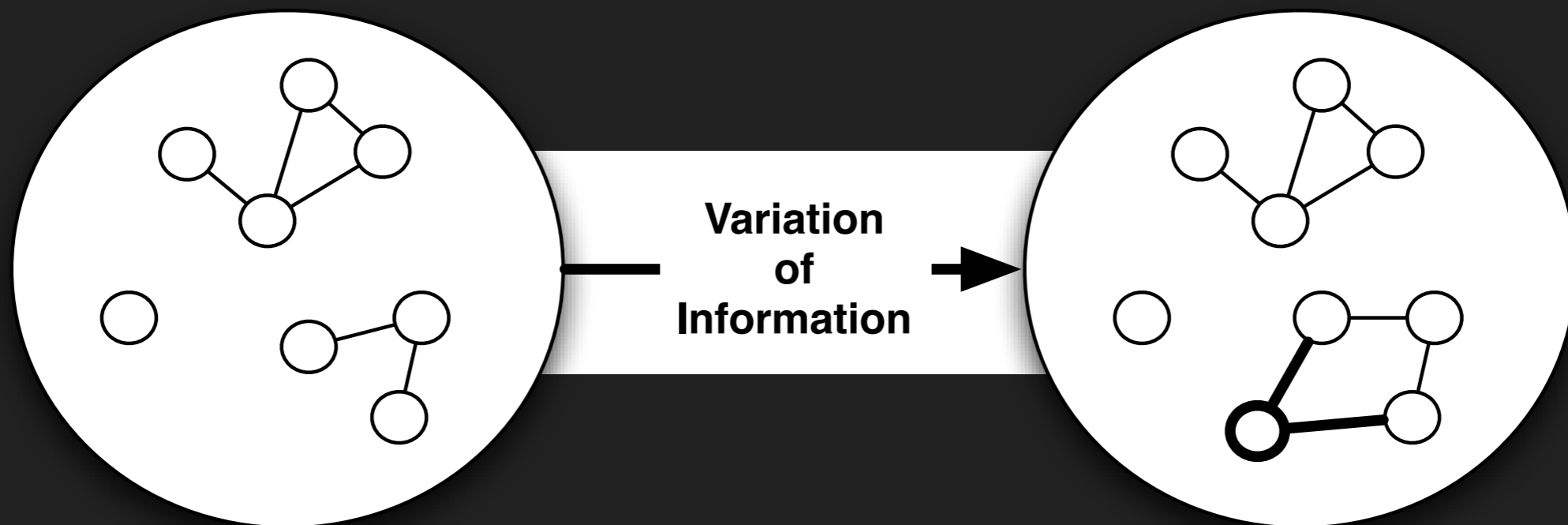Problem: real-time monitoring, amount of data

The life expectancy of a node to graph G is given by the *time to live* of this node. Each time we added a node to graph G, the TTL of the nodes that have established connections to it is incremented, and for others this value is decremented.

The variation of information evaluates the distance between different clusters. Variation in information measures the amount of information lost and gained when changing from one version of the graph to a new one.

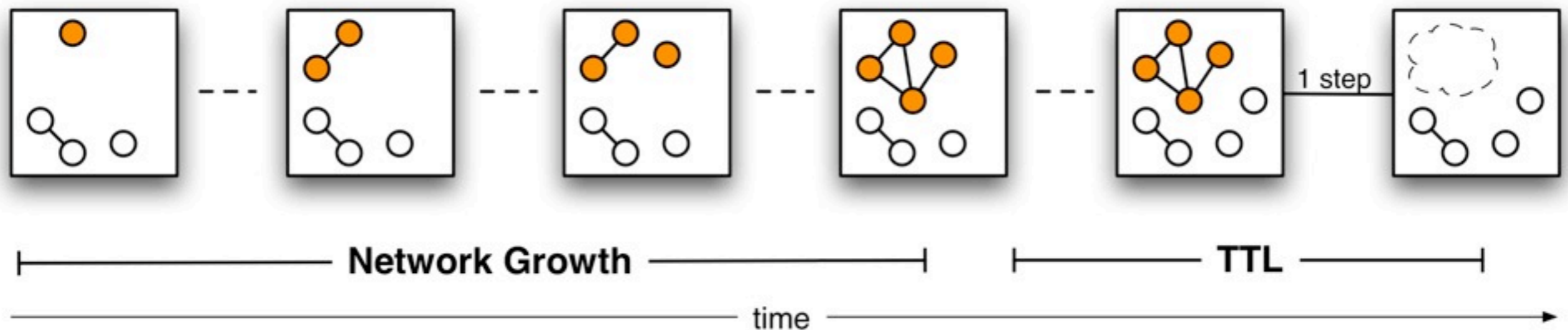VI = entropy of the 1st graph + entropy of the 2nd – mutual information of both

(Meilã, 2007)



Variation
of
Information

Sunday, March 27, 2011

## TTL + VI



Schematics of topic growth and detection via VI

High VI

Network Growth — TTL

1 step

time

Variation of Information

Sunday, March 27, 2011

now supervision is required...

Sunday, March 27, 2011

now supervision is required...

Sunday, March 27, 2011

now supervision is required...

Sunday, March 27, 2011

# On-going research - dynamics of scientific trends



http://work.theobservatorium.eu/science/

A case study of presidential elections in January 2011

Data collected from the 1st of November 2010 to the 22nd January 2011

Collection of:

- News from a representative set of portuguese media - TV stations; radio stations; newspapers and news agencies.

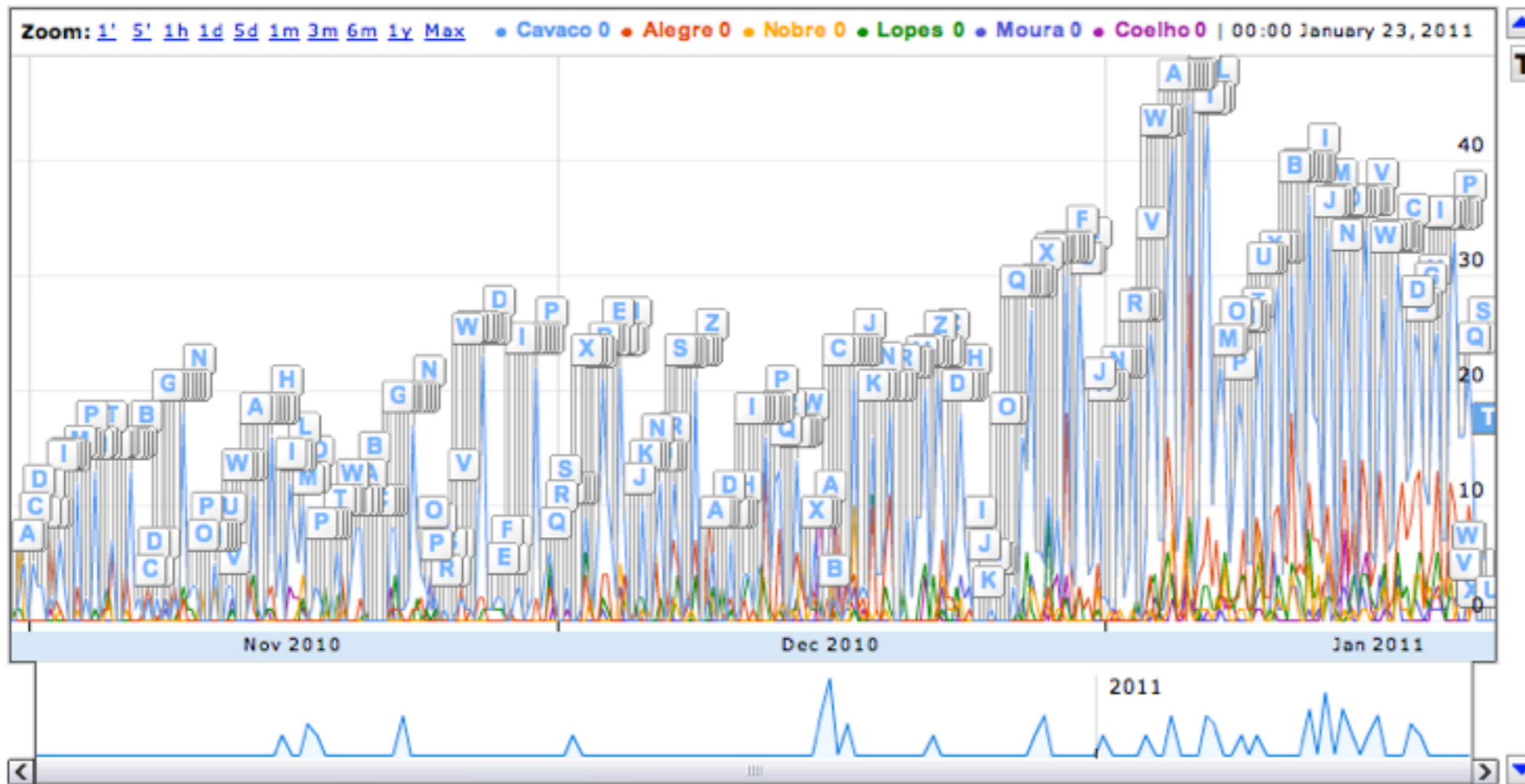- Twits

- Facebook

**Final results**

After the election pool the final results were the following:

- **Cavaco** Silva - 52.95%
- Manuel **Alegre** - 19.76%
- Fernando **Nobre** - 14.1%
- Francisco **Lopes** - 7.14%
- Manuel **Coelho** - 4.49%
- Defensor **Moura** - 1.57%
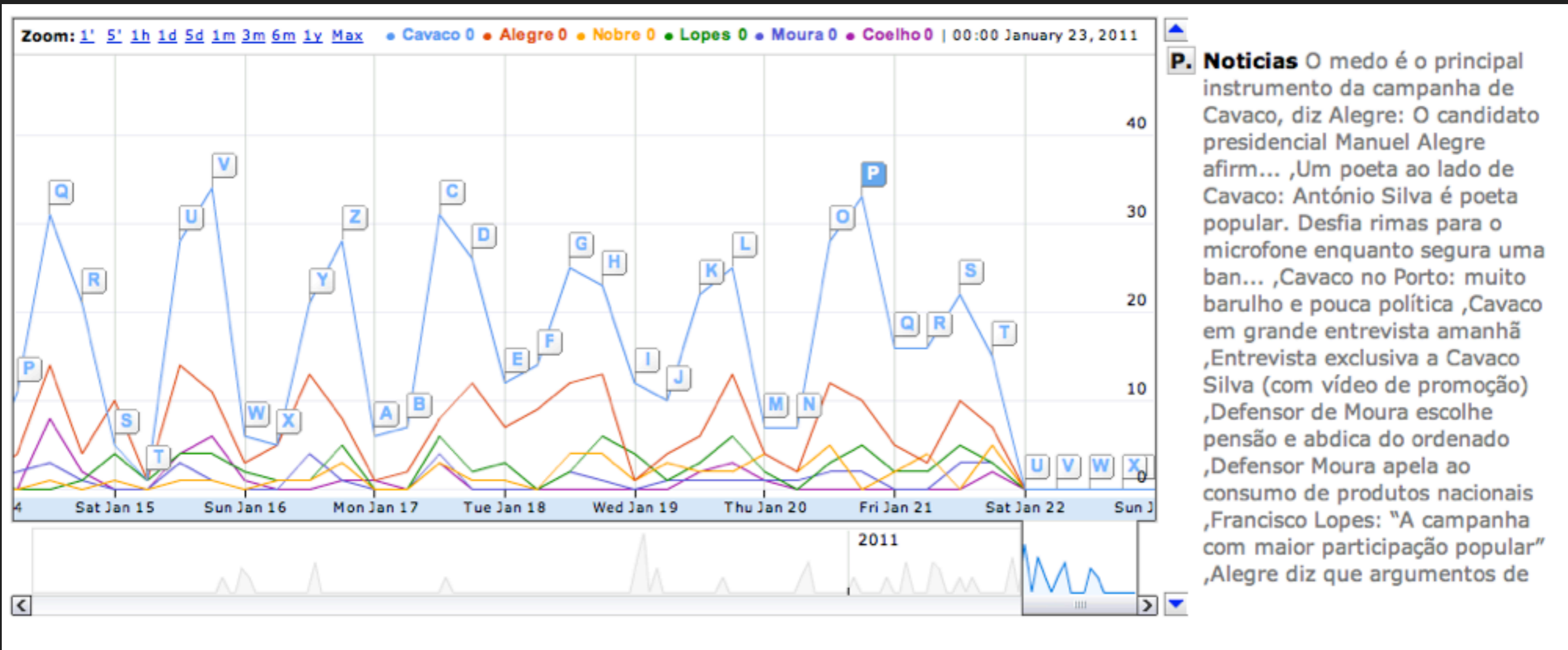
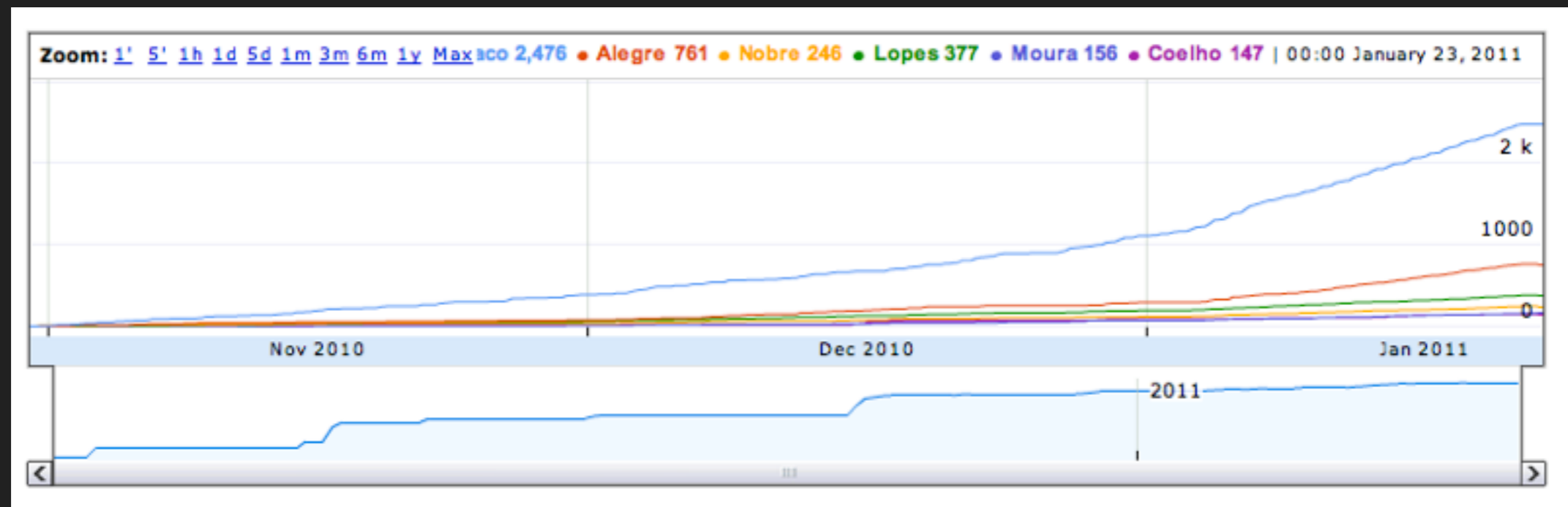## News

# On-going research - political opinion dynamics

## News

## News

## News



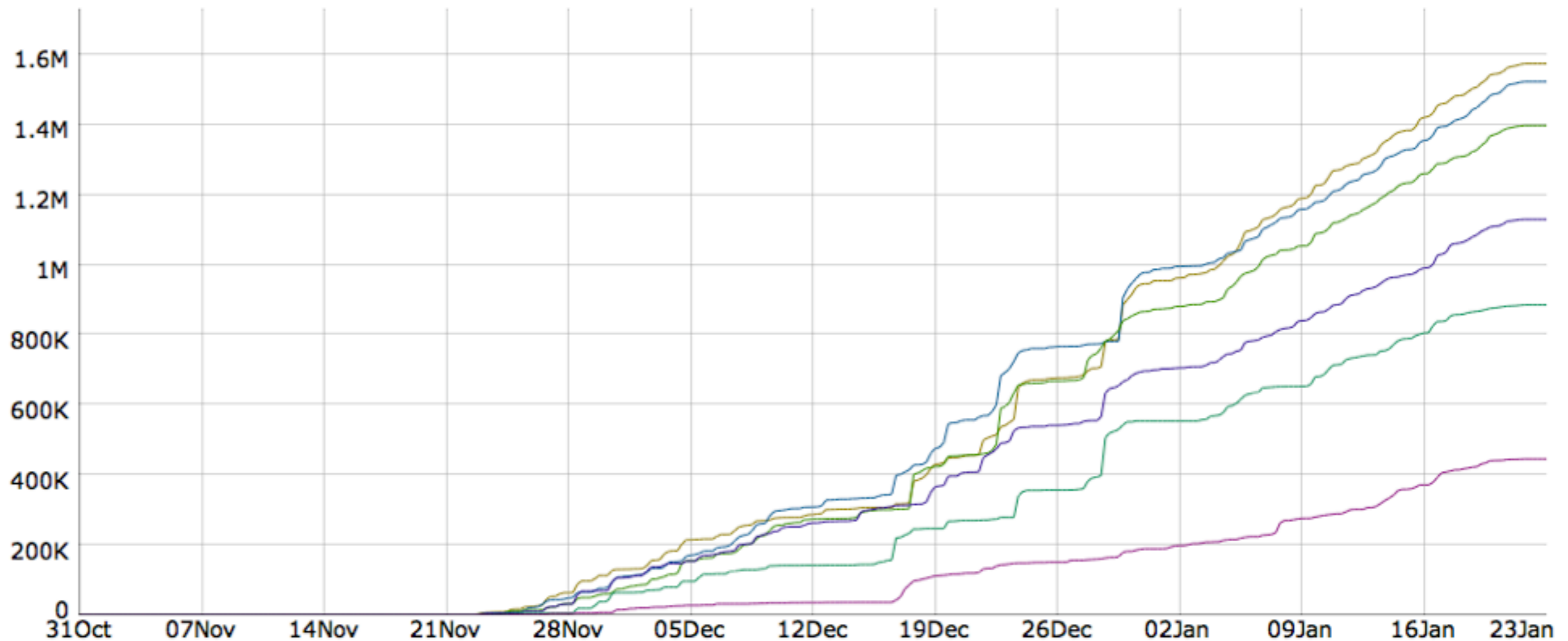Figure 1. Daily number of tweets per candidate posted on Twitter.

Sunday, March 27, 2011

## News



Figure 2. Cumulative summing of the previous chart.

Sunday, March 27, 2011

## Social networks

Goal:

monitoring financial and economic data

+

monitoring opinion dynamics concerning economy, markets, etc

How to decide what data to collect ?

## GDP – Gross Domestic Product

"economy strenght"; data from the International Monetary Fund, relative to 2008

## DMC – Domestic Market Capitalization

"market size"; data from the World Federation of Exchanges, relative to 2008

## AI - Financial Market Activity

stocks, bonds and derivatives traded

$$TI = \frac{1}{4} AI + \frac{2}{4} DMCI + \frac{1}{4} GDPI$$

| Rank | Country | Market | Index |
|---|---|---|---|
| 1 | United States | NYSE Euronext - New York | 1.0000 |
| | | NASDAQ OMX | |
| | | Chicago Board Options Exchange(CBOE) | |
| | | CME Group | |
| | | International Securities Exchange | |
| | | Intercontinental Exchange (NYSE ICE) | |
| 2 | China | Shangai Stock Exchange | 0.9183 |
| | | Shenzhen Stock Exchange | |
| 3 | Japan | JASDAQ Securities Exchange | 0.9154 |
| | | Osaka Securities Exchange | |
| | | Tokyo Stock Exchange Group | |
| 4 | United Kingdom | London Stock Exchange | 0.8906 |
| | | ICE Futures Europe | |
| | | London Metal Exchange | |
| 5 | Brazil | BM&FBOVESPA | 0.8768 |
| 6 | India | National Stock Exchange of India (NSE) | 0.8746 |
| | | Bombay Stock Exchange | |
| 7 | Korea, Rep. of | Korea Exchange | 0.8734 |
| 8 | Spain | BME Spanish Exchanges | 0.8541 |
| 9 | Italy | Borsa Italiana | 0.8506 |
| 10 | Finland | NASDAQ OMX Nordic Exchange - Helsinki | 0.8474 |
| | Denmark | NASDAQ OMX Nordic Exchange - Copenhagen | |
| | Sweden | NASDAQ OMX Nordic Exchange - Stockholm | |
| | Iceland | NASDAQ OMX Nordic Exchange - Iceland | |
| | Estonia | NASDAQ OMX Nordic Exchange - Tallinn | |
| | Latvia | NASDAQ OMX Nordic Exchange - Riga | |
| | Lithuania | NASDAQ OMX Nordic Exchange - Vilnius | |
| 11 | Australia | Australian Securities Exchange | 0.8462 |
| 12 | Canada | TMX Group | 0.8443 |
| 13 | Mexico | Bolsa Mexicana de Valores | 0.8275 |
| 14 | Belgium | NYSE Euronext - Brussels | 0.8172 |
| | France | NYSE Euronext - Paris | |
| | Netherlands | NYSE Euronext - Amesterdam | |
| | Portugal | NYSE Euronext - Lisbon | |
| 15 | South Africa | Johannesburg Stock Exchange | 0.8133 |
| 16 | Taiwan | Taiwan Stock Exchange | 0.8072 |
| 17 | Turkey | Istambul Stock Exchange | 0.8057 |

| Rank | Country | Media 1 | RSS | Media 2 | RSS |
|---|---|---|---|---|---|
| 1 | United States | REUTERS | 🟧 🟧 🟧 | THE WALL STREET JOURNAL. ONLINE | 🟧 |
| 1 | United States | Bloomberg Businessweek | 🟧 🟧 🟧 | CNN Money.com | 🟧 |
| 2 | China | 经济观察网 eeo.com.cn | 🟧 🟧 | English.news.cn NEWS | 🟧 |
| 3 | Japan | asahi.com | 🟧 | The Japan Times ONLINE | 🟧 |
| 4 | United Kingdom | FT | 🟧 | Economist.com | 🟧 |
| 5 | Brazil | Brasil Econômico | 🟧 | Valor ECONÔMICO | 🟧 |
| 6 | India | THE ECONOMIC TIMES | 🟧 | THE FINANCIAL EXPRESS | 🟧 |
| 7 | Korea | The Chosunilbo english.chosun.com | 🟧 | mk mbn | 🟧 |
| 8 | Spain | CincoDias.com | 🟧 | elEconomista.es | 🟧 |
| 9 | Italy | Il Sole 24 ORE.com | 🟧 | MF MILANO FINANZA | 🟧 |
| 10 | Sweden | di.se | 🟧 | DN.se | 🟧 |
| 10 | Finland | TALOUS SANOMAT | 🟧 | HS.fi | 🟧 |

| | | | | |
|---|---|---|---|---|
| 10 | Iceland | mbl.is | 🔲 | IceNews |
| 10 | Denmark | Børsen | 🔲 | Business.dk |
| 10 | Estonia | DELFI MAJANDUS | 🔲 | aripaev.ee |
| 11 | Australia | TRADING ECONOMICS | 🔲 | THE AUSTRALIAN |
| 12 | Canada | les affaires.com | 🔲 | Benefits CANADA |
| 13 | Mexico | EL ECONOMISTA | 🔲 | ElFinanciero EN LINEA |
| 14 | Belgium | DE TIJD | 🔲 | L'Echo |
| 14 | France | LesEchos.fr | 🔲 | LA TRIBUNE.fr |
| 14 | Netherlands | fd.nl | 🔲 | vkml |
| 14 | Portugal | Diário Económico | 🔲 | AF |
| 15 | South Africa | MONEYWEB | 🔲 | FM |
| 16 | Taiwan | Taiwan News | 🔲 | CommonWealth magazine |
| 17 | Turkey | DAILY NEWS | 🔲 | Hürriyet hurriyet.com.tr |
| 18 | Germany | FINANCIAL TIMES DEUTSCHLAND | 🔲 | FAZ FINANCE.NET |
| 19 | Hong Kong | FT中文网 | 🔲 | THE WALL STREET JOURNAL ASIA |
| 20 | Norway | DinSide | 🔲 | DN.no |
| 21 | Argentina | ámbito.com | 🔲 | CRONISTA |

## Table Of Contents

## Next topic

What is Theseus?

## This Page

Show Source

## Quick search

[Go]

Enter search terms or a module, class or function name.

# Theseus's documentation

Latest release: 0.7.1

Latest update: February 24, 2011

## Contents

- What is Theseus?
- What will Theseus be at version 1.0?
- Theseus now!
    - theseus.processor
    - theseus.crawler
    - theseus.utils
    - theseus.examples
- Download Theseus
- How to
    - Process 11 TXT files inside a "TXT" folder
- Frequently Asked Questions (FAQ)
- Contact The Observatorium
- Roadmap
    - 0.8
    - 0.7.1 **Present Version**
    - 0.7
    - 0.6
    - 0.5.1

## Indices and tables

- *Index*
- *Module Index*
- *Search Page*

http://www.theobservatorium.eu/html/

Sunday, March 27, 2011

# The team

António Fonseca

David Rodrigues

José António Silva

Manuel Tânger

Pedro Loureiro

Jorge Louçã

software is available under open source licence

data is available under request

# Thank you !