# Information Search in Web Archives

Miguel Costa

Portuguese Foundation for National Scientific Computing
LaSIGE @ Faculty of Sciences, University of Lisbon
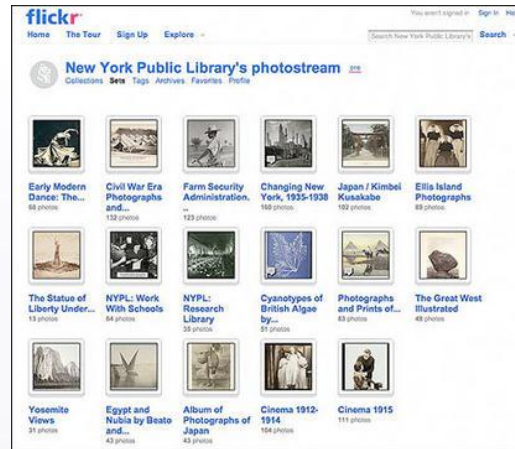
1st International Workshop on Archiving Community Memories
September 6, 2013, Lisbon, Portugal
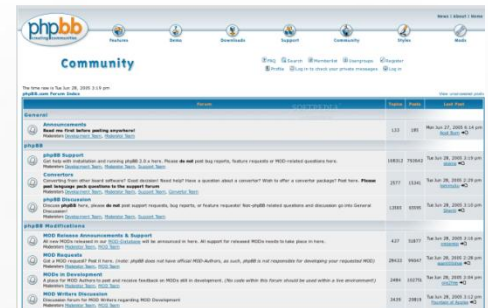
# Our Memory is in Digital Form

## E-books



## Web photo galleries



## Forums



## Blogs



## Online newspapers



## Social networks

# The Web is Ephemeral

- 50 days - 50% of documents are changed

  (Cho and Garcia-Molina. 2000)


- 1 year - 80% of documents become inaccessible

  (Ntoulas, Cho and Olson. 2004)


- 27 months - 13% of web references disappear

  ([http://webcitation.org/](http://webcitation.org/). 2007)

# Will we face a Digital Dark Age?



The page cannot be found

The page you are looking for might have been removed, had its name changed, or is temporarily unavailable.
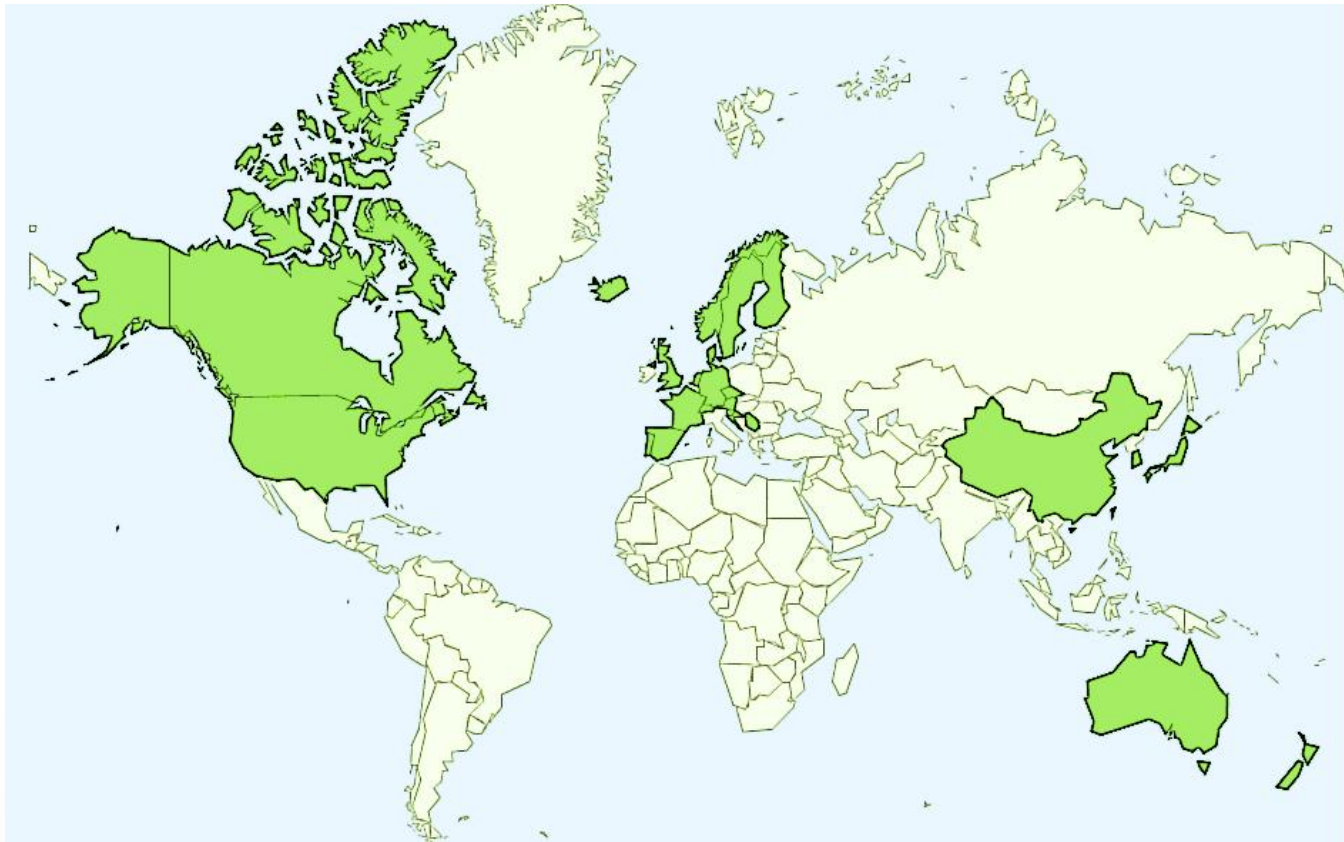
Please try the following:

- If you typed the page address in the Address bar, make sure that it is spelled correctly.
- Open the httpd.apache.org home page, and then look for links to the information you want.
- Click the ← Back button to try another link.
- Click 🔍 Search to look for information on the Internet.

HTTP 404 - File not found
Internet Explorer

# 2010: Worldwide Web Archiving Initiatives



- +42 initiatives in 26 countries

- +180 billions of web contents since 1996 (6.6 PB)

# Wikipedia Page with Updated List

Article | Discussion

Read | Edit | View history | Search

## List of Web archiving initiatives

From Wikipedia, the free encyclopedia
(Redirected from List of Web Archiving Initiatives)

This page contains a list of Web archiving initiatives worldwide. For easier reading, the information is divided in three tables: web archiving initiatives, archived data and access methods.

**Contents** [hide]
1 Web archiving initiatives
2 Archived data
3 Access methods
4 References

Map of Web archiving initiatives worldwide in June, 2011.

### Web archiving initiatives                                    [edit]

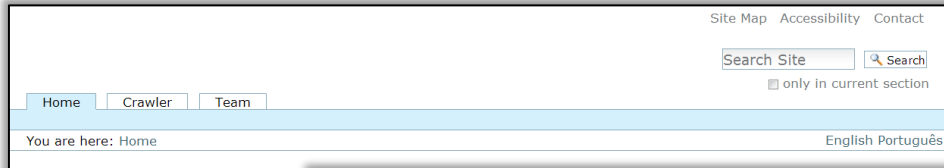| Name ⋈ | Country ⋈ | Creation Year ⋈ | Technologies ⋈ | Number of Employees ⋈ | | Comments ⋈ |
|---|---|---|---|---|---|---|
| | | | | Full-time | Part-time | |
| Australia's Web Archive[1] | Australia | 1996 | PANDORA Digital Archiving System (PANDAS), NLA Trove, HTTrack. | 4 | >4.25 | It is a collaborative program of 11 agencies that provide an estimate average monthly staffing equivalent to 4 FTE. IT outsourced support: 0.25 person-month. Whole Domain Harvests are conducted by the Internet Archive using Heritrix, Wayback Machine. |
| Our digital island, a Tasmanian Web Archive[2] | Australia | 1996 | HTTrack, Experimentally: Web Curator, Heritrix and Wayback Machine | | 1 | |
| | | | Archive-access tools and | | | |

http://en.wikipedia.org/wiki/List_of_Web_archiving_initiatives

6

# 2013: Worldwide Web Archiving Initiatives

- +77 initiatives in 39 countries
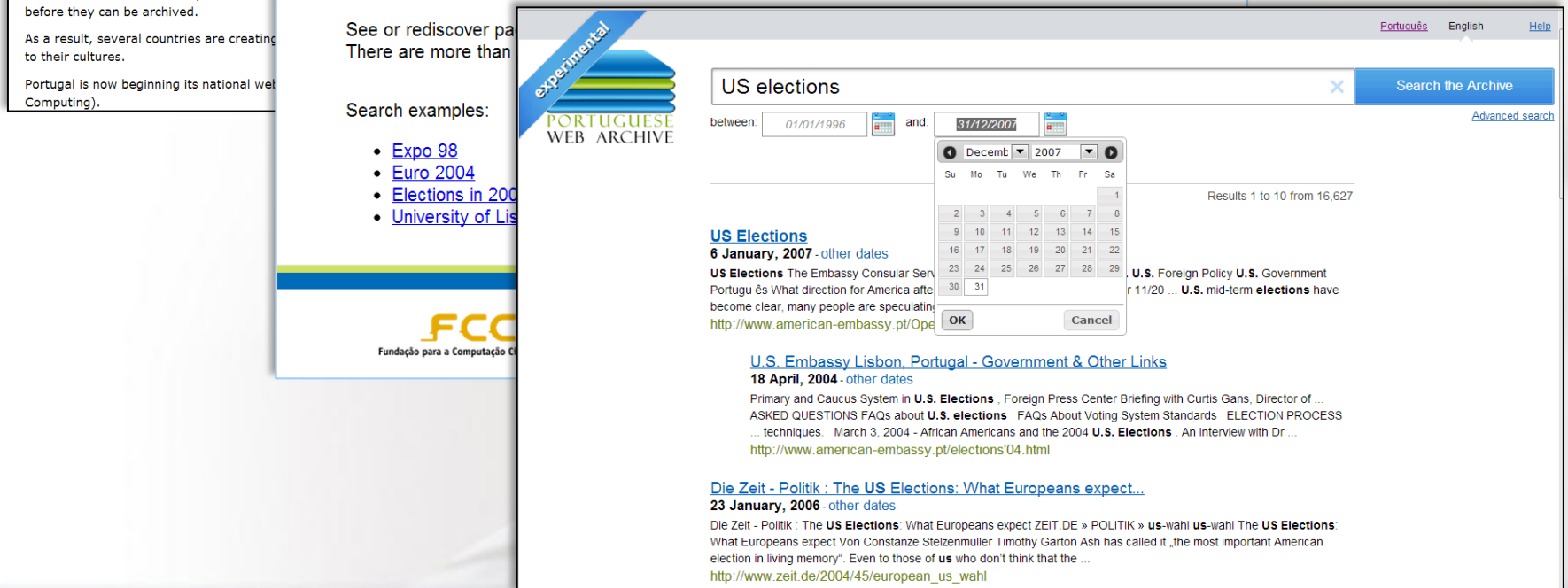- +294 billions of web contents since 1996 (8.5 PB)
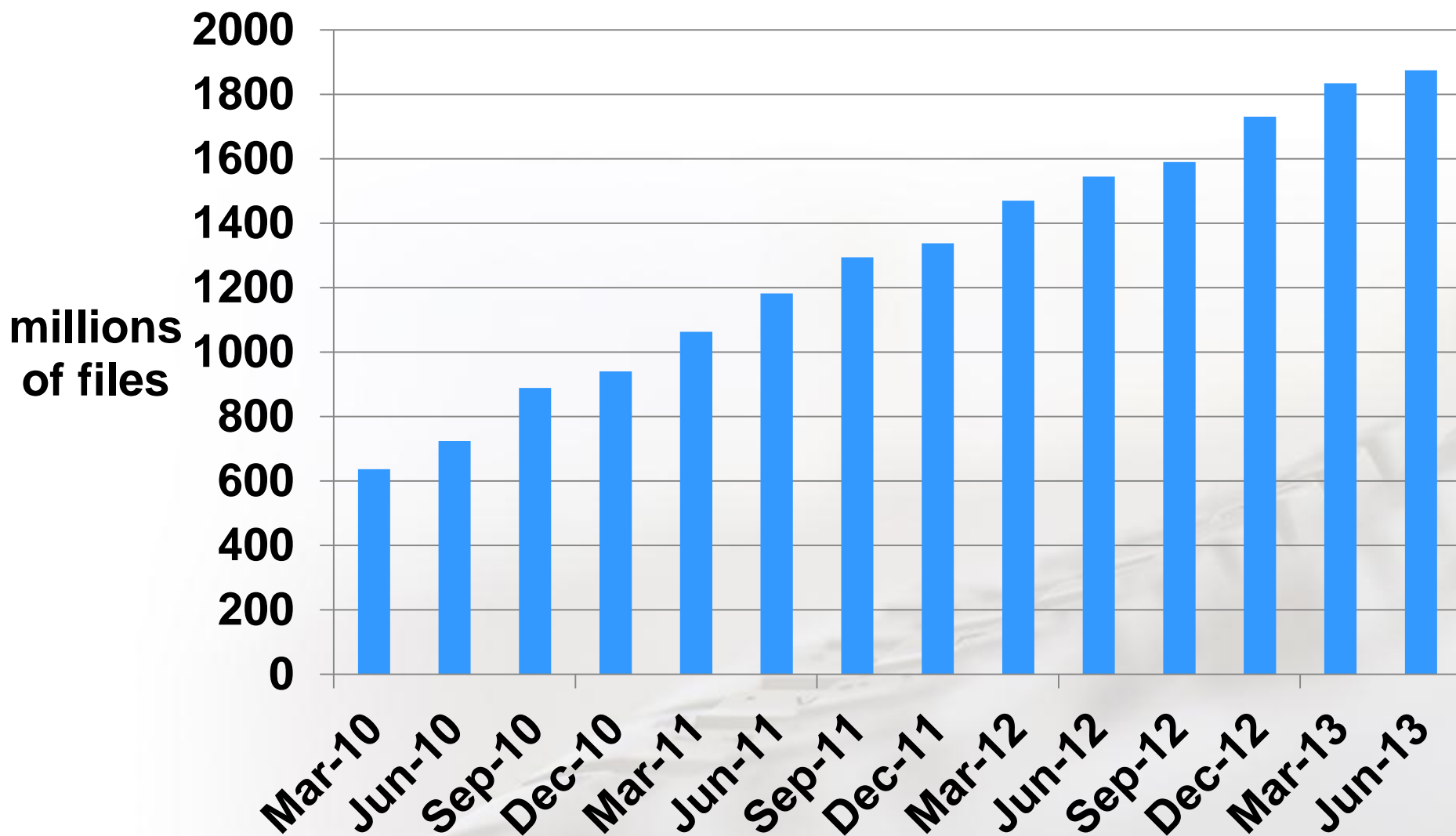
7

# Portuguese Web Archive

# 1.8 Billion Archived Files (52 TB)

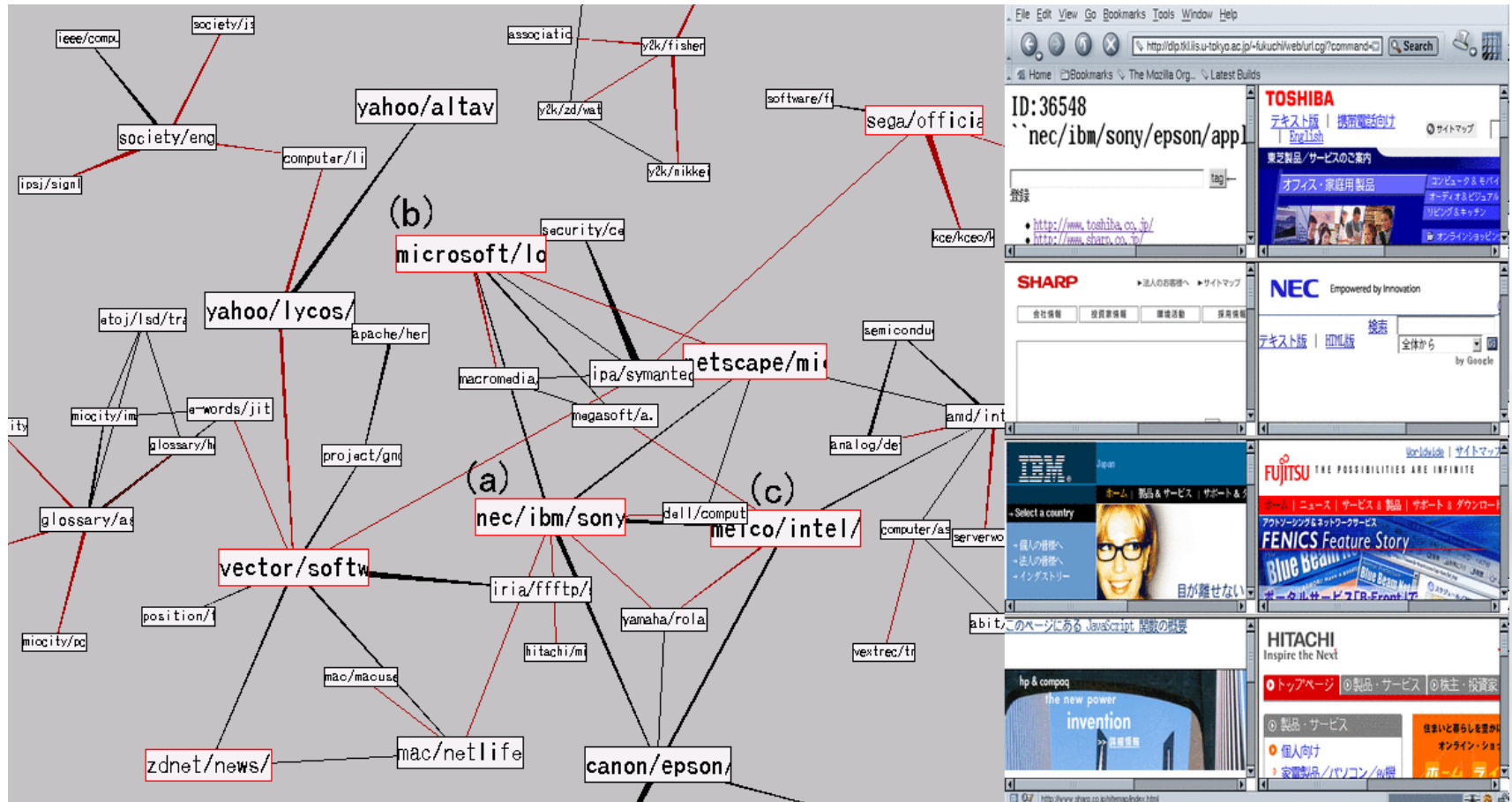# What can we do with all this data?

# Portuguese WA: Web Characterization

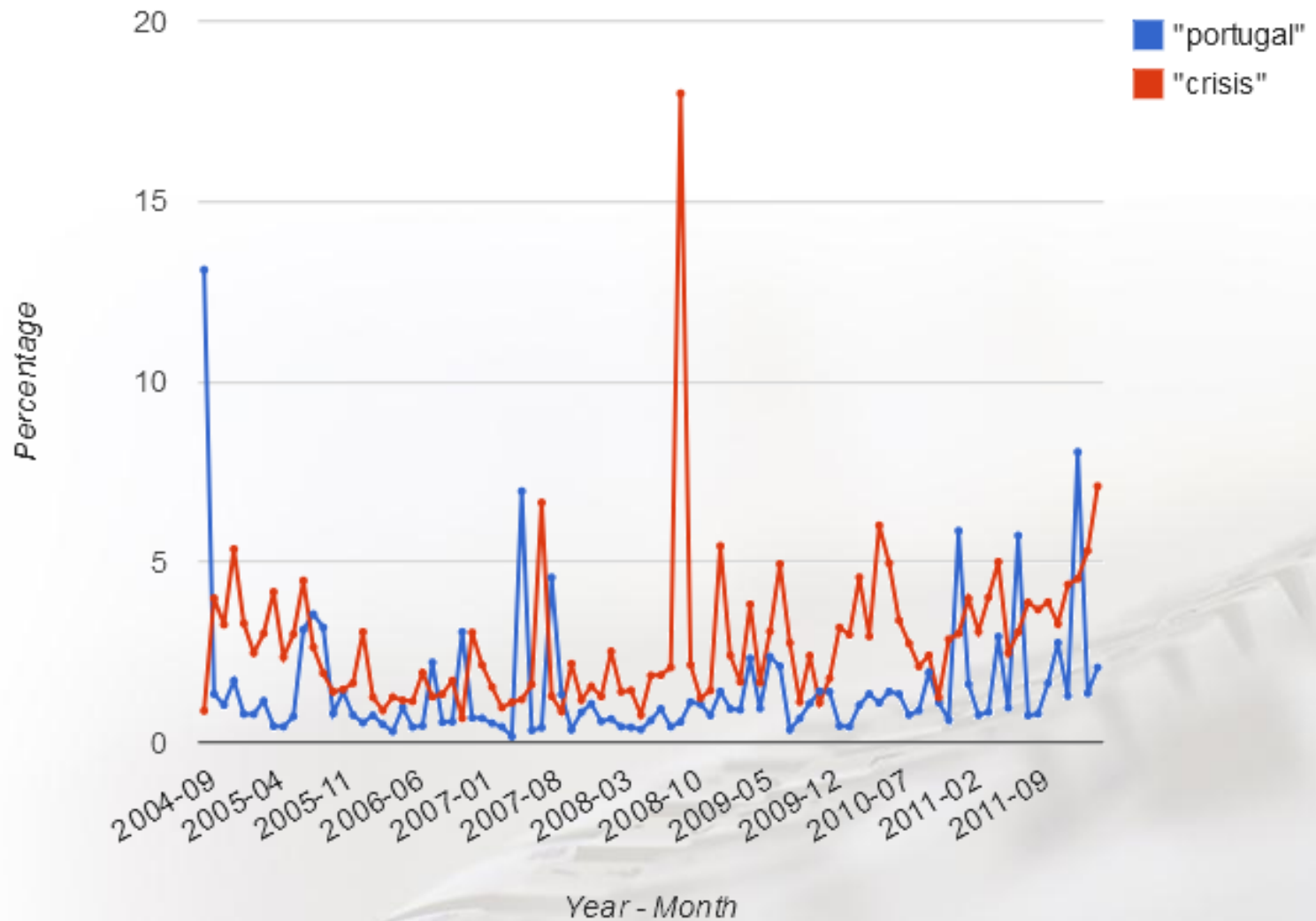| Media type | % contents 2005 | % contents 2008 | Trend |
|---|---|---|---|
| Text/html | 61.2% | 57.8% | -5.5% |
| Image/jpeg | 22.6% | 22.8% | +1.2% |
| Image/gif | 11.4% | 9.4% | -17.4% |
| Text/pdf | 1.6% | 1.9% | +18.5% |
| Other | 3.2% | 8.1% | - |

Trends in Web Characteristics, *7th Latin American Web Congress* 2009.

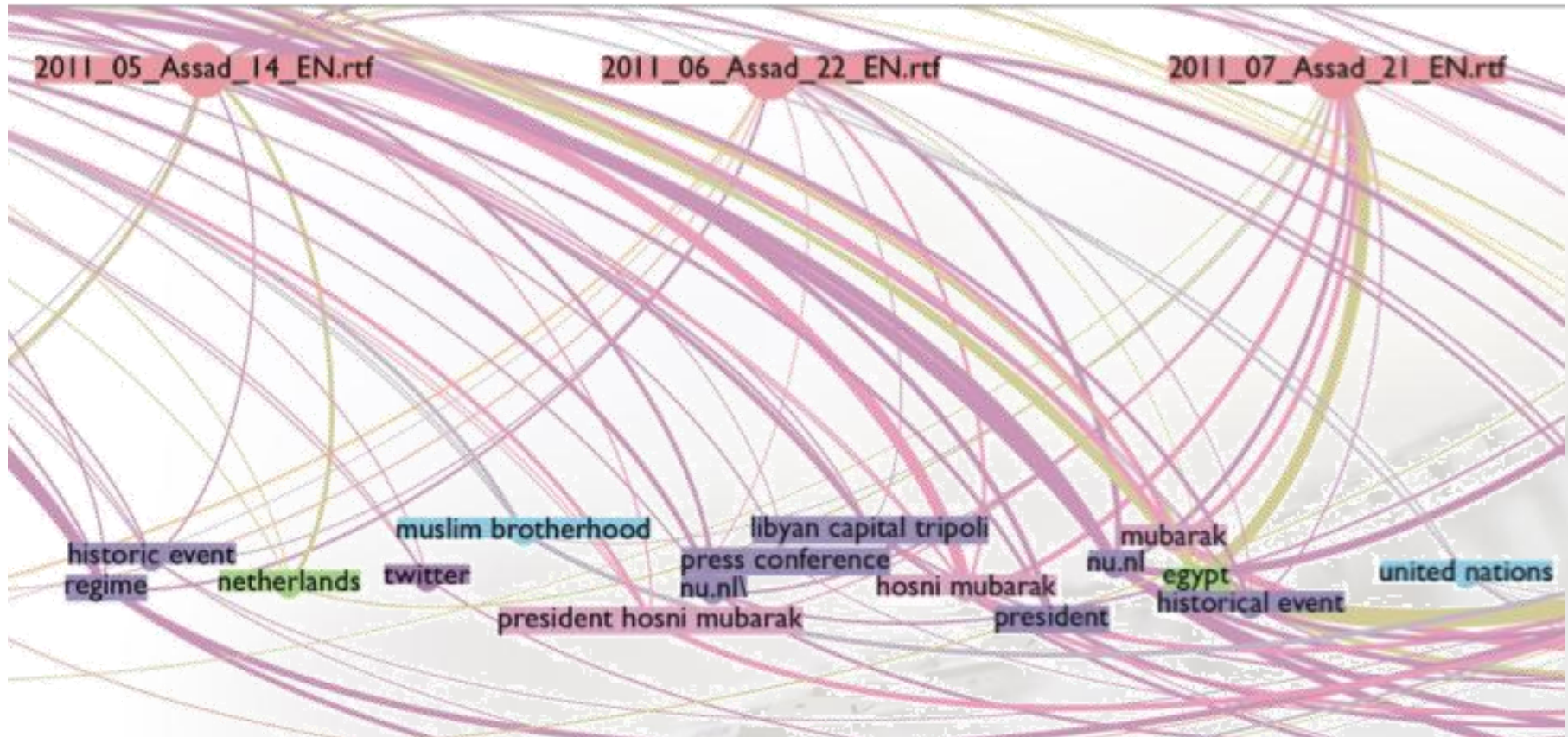# Japanese WA: Evolution of Web Communities



Extracting Evolution of Web Communities from a Series of Web Archives, *14th ACM Conference on Hypertext and Hypermedia* 2003.

# UK WA: Word Frequency Analysis



http://www.webarchive.org.uk/ukwa/ngram/

# WebART: Co-word Analysis
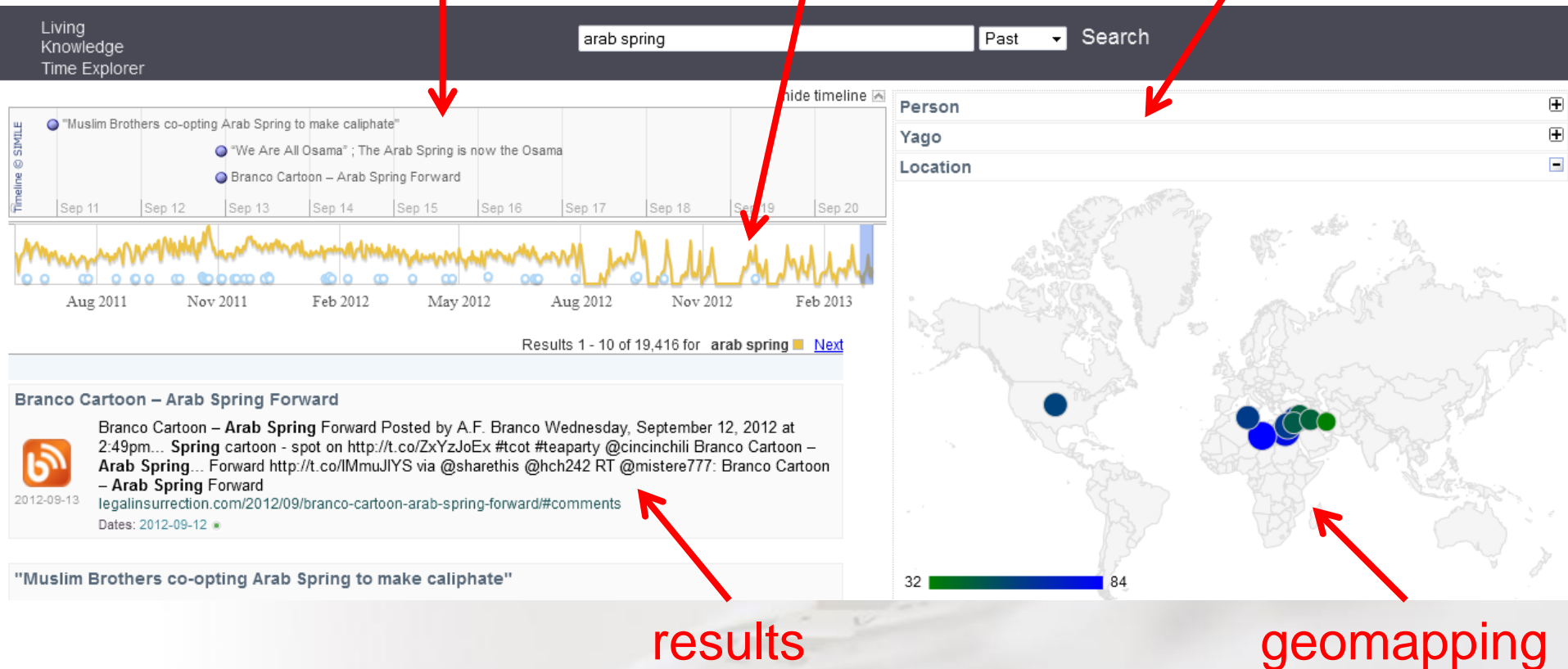


http://www.webarchiving.nl/

# Living Knowledge (Yahoo!): News Analysis

timeline

frequency graph

entity selection



results

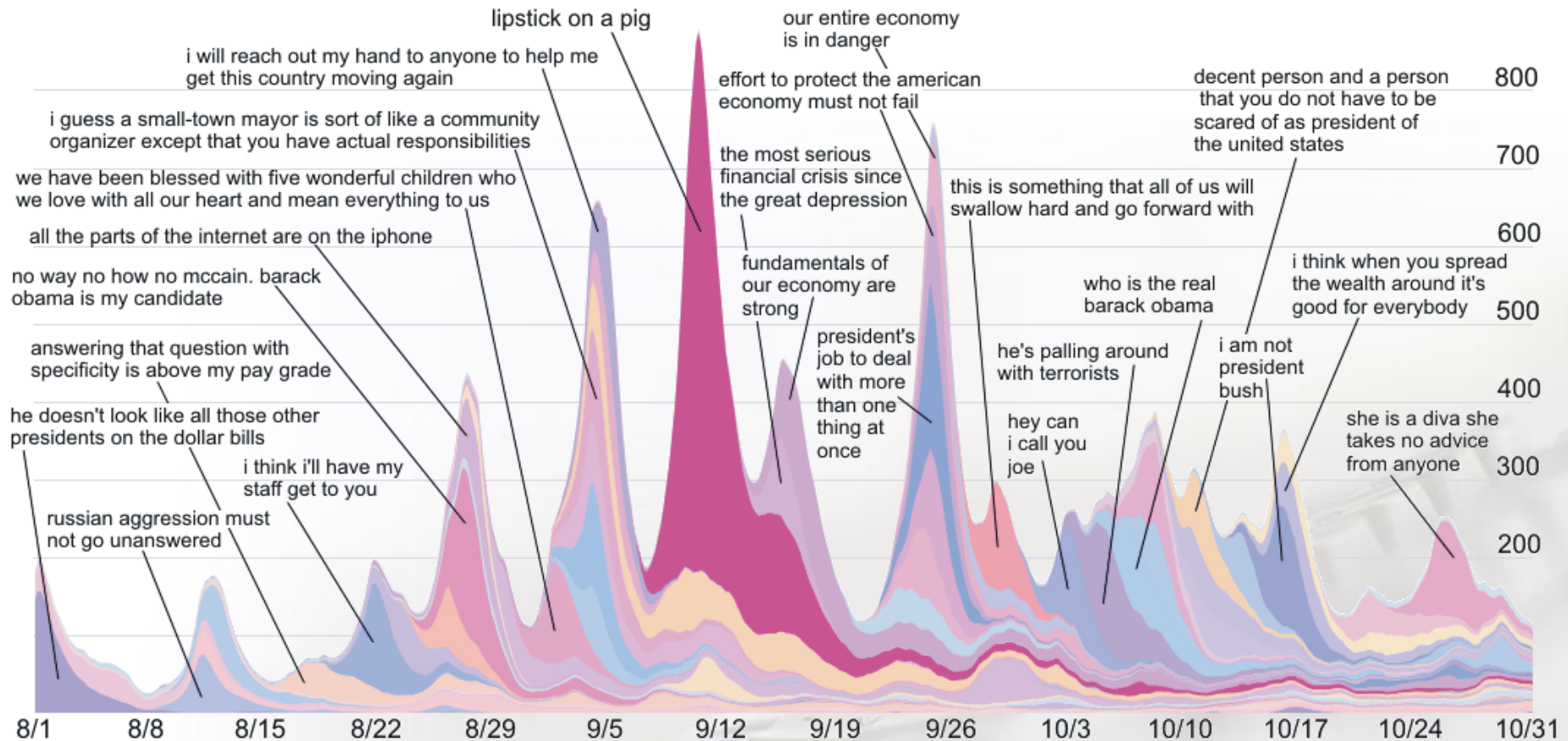geomapping

Searching through time in the New York Times, *Human Computer and Information Retrieval* 2010.

# MemeTracker: News Analysis



Meme-tracking and the Dynamics of the News Cycle, *Knowledge Discovery and Data Mining* 2009*.*

# Facebook: Social Analysis



Paul Butler created this friend relationship visualization map using Facebook data.

# Twitómetro: Sentiment Analysis



http://dmir.inesc-id.pt/project/Reaction

# NYT Archive: Forecast Events

In deep **drought**, at 104 degrees, dozens of Africans are **dying**.
New York Times 02/17/2006

Angola **cholera** cases
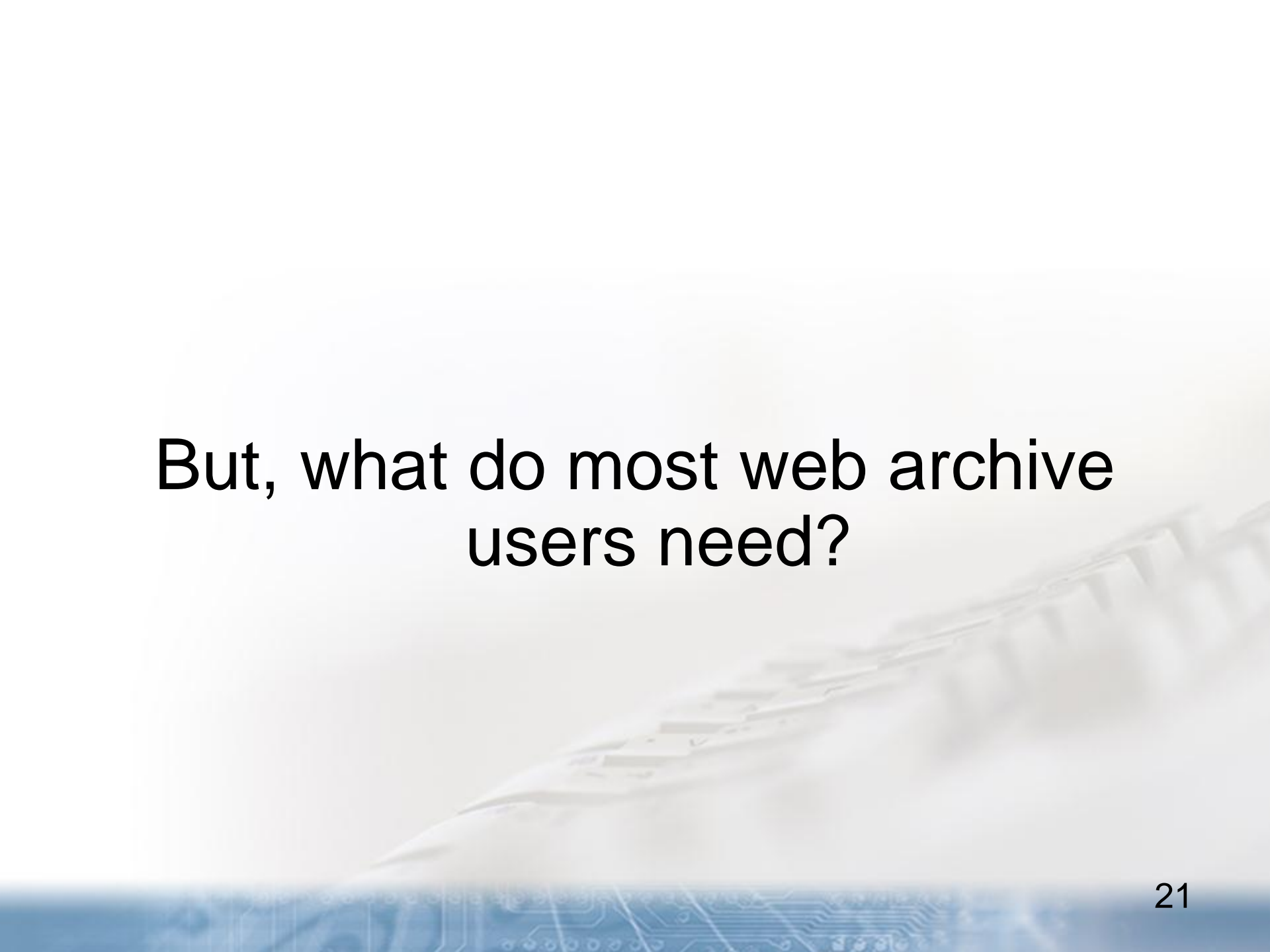rise sharply after **floods**.
New York Times 01/30/2007

Mining the Web to Predict Future Events, *Web Search and Data Mining* 2013.

# What can we do with all this data?

## all kinds of machine learning over time
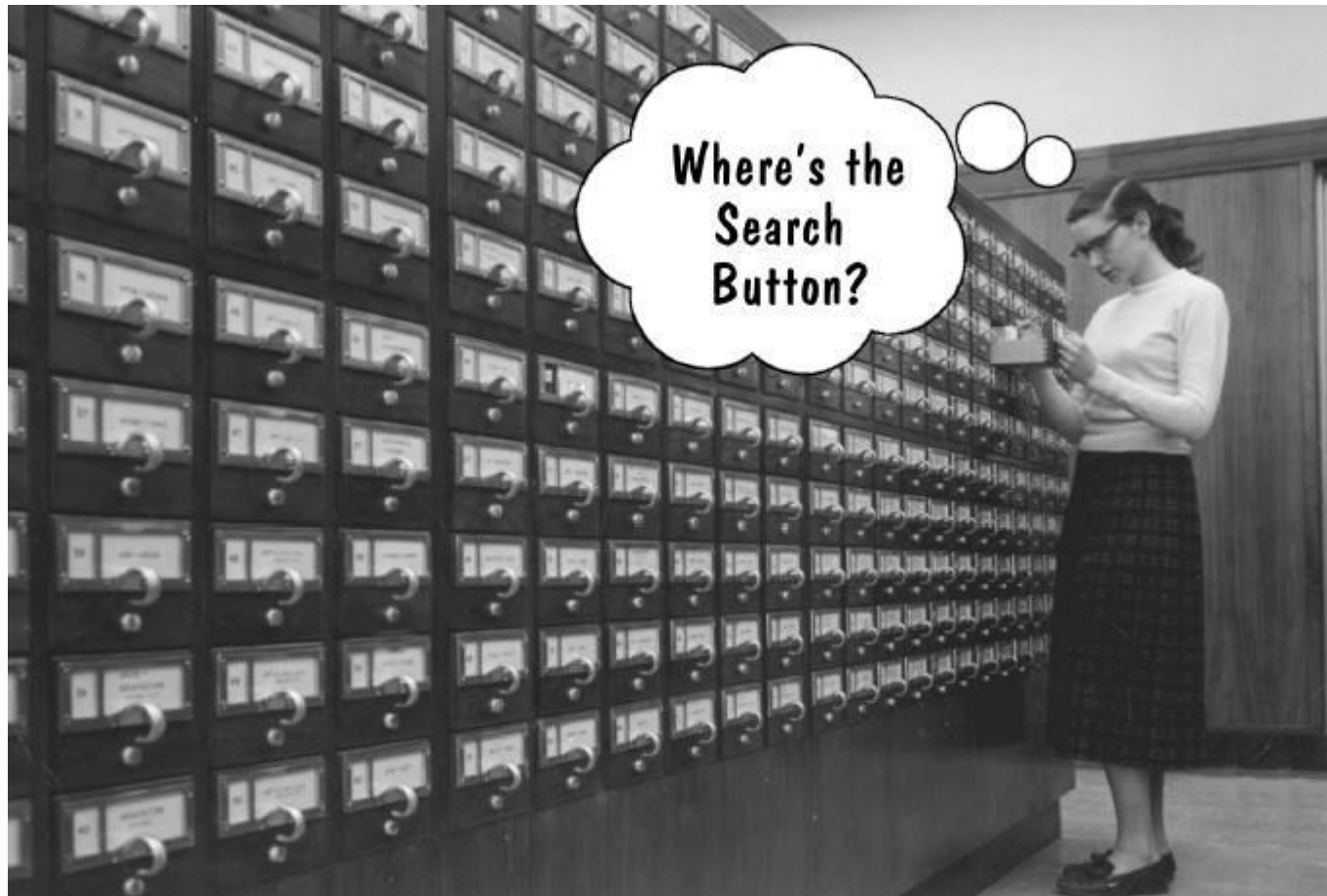
model the past and predict the future

# But, what do most web archive users need?

# Use Cases

- **User** visits a missing bookmark
- **Journalist** investigates past information
- **Webmaster** recovers the lost site
- **Historian** searches for digital documents
- **Web designer** creates portfolio of sites
- **Professor** downloads missing slides
- **Lawyer** looks for evidences

# Users don't understand Web Archives

## What do you want? I don't know!

# PWA Search System



- Available since 2010: http://archive.pt
- 1.2 billion documents
  - searchable by full-text and URL
  - range between 1996 and 2012

# URL Search

| sapo.pt | × | **Search the Archive** |

between: 01/01/1996 📅 and: 31/12/2012 📅

Advanced search

Did you want to see webpages with the text: http://sapo.pt?

## Versions of the archived the Web pages

We archived 1,832 versions of the Web page http://sapo.pt from 1 January, 1996 and 26 August, 2013.

| 1997 2 | 1998 4 | 1999 23 | 2000 87 | 2001 58 | 2002 20 | 2003 29 | 2004 199 | 2005 444 | 2006 119 | 2007 120 | 2008 5 | 2009 6 | 2010 255 | 2011 368 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 8 Oct | 10 Jan | 25 Jan | 29 Feb | 5 Jan | 24 Jan | 5 Feb | 16 Feb | 1 Jan | 1 Jan | 2 Jan | 1 Jan | 20 May | 26 Mar | 1 Jan |
| 10 Dec | 29 Jan | 25 Jan | 29 Feb | 6 Jan | 6 Feb | 10 Feb | 19 Mar | 2 Jan | 1 Jan | 5 Jan | 14 Mar | 24 Jun | 1 Apr | 2 Jan |
| | 7 Feb | 8 Feb | 29 Feb | 7 Jan | 30 Mar | 19 Feb | 5 Apr | 3 Jan | 2 Jan | 7 Jan | 14 Mar | 26 Sep | 5 Apr | 3 Jan |
| | 7 Feb | 8 Feb | 29 Feb | 8 Jan | 1 Apr | 20 Feb | 20 May | 4 Jan | 2 Jan | 7 Jan | 22 Oct | 26 Sep | 8 Apr | 4 Jan |
| | | 9 Feb | 1 Mar | 19 Jan | 29 May | 24 Mar | 3 Jun | 4 Jan | 5 Jan | 9 Jan | 22 Oct | 18 Dec | 9 Apr | 5 Jan |
| | | 20 Feb | 3 Mar | 24 Jan | 30 May | 12 Apr | 9 Jun | 5 Jan | 6 Jan | 11 Jan | | 18 Dec | 12 Apr | 6 Jan |
| | | 20 Feb | 3 Mar | 30 Jan | 4 Jun | 19 Apr | 9 Jun | 5 Jan | 10 Jan | 12 Jan | | | 13 Apr | 7 Jan |
| | | 21 Apr | 4 Mar | 4 Feb | 6 Jun | 22 Apr | 11 Jun | 6 Jan | 10 Jan | 14 Jan | | | 16 Apr | 8 Jan |
| | | 23 Apr | 4 Mar | 10 Feb | 7 Jun | 24 Apr | 12 Jun | 7 Jan | 11 Jan | 16 Jan | | | 19 Apr | 9 Jan |

# SAPO.PT 1997

# Full-text Search

# Understanding the Users' Information Needs



data richness

Laboratory
Studies

Interactive
Questionnaires

Search
Logs

generalization

[03/02/2012  21:16:11] QUERY fcul
[03/02/2012  21:16:19] CLICK RANK=1

Understanding the Information Needs of Web Archive Users, *10th International Web Archiving Workshop*  2010.

# What are the Users' Information Needs?

- **Navigational** – 53% to 81%
  - seeing a web page in the **past** or how it evolved

- **Informational** – 14% to 38%
  - collecting information about a topic written in the **past**

# What is the best tool to support **navigational** (and **informational**) information needs?

## Searching vs Analytical tools

- **URL Search** – Internet Archive's Wayback Machine
  - difficult to remember or unknown

- **Full-text Search** – Lucene extensions (NutchWAX & Solr)
  - does not scale for large collections
  - slow searches
  - poor quality results

~~How to improve?~~

How to evaluate?

New Technology

Evaluation

Is it better than State-of-the-Art?

- Test Collection (Cranfield Paradigm):
  - Corpus
    - **What** are the typical web collections?
  - Topics
    - **Why**, **what** and **how** do users search?
  - Relevance Judgments
    - **What** is relevant for users?
  - Measures
    - **What** and **how** many documents do users see?

- Wrong assumptions lead to wrong conclusions

- Test Collection (Cranfield Paradigm):
  - **Corpus**: 6 web collections, 255M contents, 8.9TB
    - broad crawls, selective crawls, integrated collections
  - **Topics**: 50 navigational
    - I need the page of Público newspaper between 1996 and 2000.
  - **Relevance Judgments**: 3 judges, 3-level scale of relevance, 267 822 versions assessed
  - **Measures**: (S@k, NDCG@k, P@k | k=1,5,10)
    - only 14% see the second page (> top 10)

Evaluating Web Archive Search Systems, *13th International Conference on Web Information System Engineering* 2012.

# Evaluation Metric: Success@k

– 1 if a relevant version has been found on the top-k

– 0 otherwise

– Example:
- avg. Success@5 = 3/4

# Search Effectiveness of the State-of-the-Art

# State-of-the-Art (SoA) Effectiveness



Success@1

0.28
SoA

0.65
Best in TREC

0.84
Google

0
1

Success@5

0.6
SoA

0.84
Best in TREC

0.9
Google

0
1

Success@10

0.8
SoA

0.88
Best in TREC

0.92
Google

0
1

# How to improve?

# Goal: Maximize Relevance

How?

Relevant document

Non-relevant document

Query: sapo (toad in English)

http://www.sapo.pt



1 de Janeiro — Faça do SAPO a sua homepage

Acesso: Internet SAPO | Meo Fibra | Meo | Telepac | TMN | TVT

Área de Cliente SAPO

Web | Imagens | Vídeos Novo! | Notícias | Blogs | Produtos | PAi | PBi

Insira o texto a pesquisar

Pesquisar

Mail Blogs Carros Casas Fotos Mapas Vídeos Notícias Messenger

Todo o SAPO

Qual o clube que José Mourinho passou a comandar em 2009? Entre no Jogo do Ano 2010!

PUB

Destaques Desporto | Economia | Vida | Tecnologia | Local Vídeos

Fontes: DD | DZ | Lusa | RTP | SIC | Sol

Mail Tempo Horóscopo Emprego

personalizar página

Mulheres no poder
192 países, apenas nove mulheres presidentes (SAPO)

Transportes
Utentes criticam aumento dos preços dos transportes (Sol)

Argentina
Sismo de 6,9 na escala de Richter abalou Argentina (SIC)

Ano Novo
Balanço de atentado em Alexandria ascende a 21 mortos (SAPO)

Mulheres no poder
192 países, apenas nove mulheres presidentes (SAPO)

São 192 os Estados com assento nas Nações Unidas e a maioria são repúblicas, mas o ano 2011 começa

HTTP Status 404 -
/wayback
/20110101160218
/about:blank

type Status report

message /wayback/20110101160218/about:blank

description The requested resource (/wayback
/20110101160218/about:blank) is not available.

42

# What other Features can we use?

# Document Change Over Time

## Change vs. Relevance

documents with higher relevance are more likely to change

documents with higher relevance tend to change to a greater degree

## Change Amount vs. Relevance

Leveraging Temporal Dynamics of Document Content in Relevance Ranking, *Web Search and Data Mining* 2010.

# Vocabulary Change Over Time



**nytimes.com**

chilies
saviano

neediest
wheelspin
sebnem

gotbaum
n.r.a.
archibold
grynbaum
greentech

m.t.a.

mazzetti
u.a.w.
lede
tavernise
gamete
s.u.v.
stolberg
nagourney
nytimes

0  100  200  300  400  500  600

time since t0

**allrecipes.com**

sugared
yams
merrymaking

imparts
wontons
soups

tiera
latkes
pureed
pregn
hila
simmered
challah
frightfully
marinades
stews
cook's
mouthwatering
weeknight

cooks

0  100  200  300  400  500  600

time since t0

the most persistent terms are descriptive of the main topic

The Web Changes Everything: Understanding the Dynamics of Web Content, *Web Search and Data Mining* 2009.

# # Versions vs Relevance



documents with higher relevance tend to have more versions

# Modeling Temporal Information

$$f_{Versions}(d) = log_y(x)$$

Parameters:
x = number of versions of document d
y = maximum number of versions of
    a document in the collection

Assumption: persistent documents are more relevant

# Search Effectiveness of the State-of-the-Art + #Versions

# New Ranking Model: $f_{SoA} + f_{Versions}$



50

# If one is good, more is better

**Term-weighting:**
BM25
Lucene
NutchWAX

...

**Term-distance:**
MinPair
MinSpanOrdered
MinSpanUnordered

...

**Web-graph based:**
Inlinks
Outlinks
PageRank

...

**Temporal:**
NumberVersions
BoostOlder
Age

...

**URL based:**
UrlLength
UrlDepth
UrlSlashes

...

# Supervised *Learning-to-Rank* Framework

Loss function for minimization

$$= \sum_{i=1}^{m} L(y^{(i)}, f(\vec{x^{(i)}}))$$

Training data → Learning System

model h()

Test data → Ranking System → Prediction

$$f(\vec{x}) = \sum_{n=1}^{p} \lambda_n {}_* f_n(\vec{x})$$

# Dataset for L2R in Web Archives

- 39 608 quadruples <query, version, relevance grade, features>
  - 50 **queries** (navigational topics)
  - 843 **versions** assessed on average per query
  - **3-level scale** of relevance
  - 68 ranking **features** extracted

- File Format:     **used in training**

| Rel. | Query | Features | Doc. Version |
|------|-------|----------|--------------|
| 2 | qid:21 | 1:0.70  2:0.34  3:0.27 ... 68:0.86 | # id114746079 |
| 0 | qid:22 | 1:0.05  2:0.18  3:0.14 ... 68:0.43 | # id172346033 |
| 1 | qid:22 | 1:0.75  2:0.33  3:0.84 ... 68:0.54 | # id456334535 |

- Normalized Discounted Cumulative Gain at cut-off k
  – total gain accumulated at a particular rank $p$

$$DCG_p = \sum_{i=1}^{p} \frac{2^{rel_i} - 1}{log(1+i)}$$

NDCG@5 = 1     NDCG@5 ≈ 0.4

55

# Search Effectiveness of the 68 features

# Results of L2R Algorithms

| | State-of-the-Art | | L2R algorithms (68 features) | | | |
|---|---|---|---|---|---|---|
| Metric | Lucene | NutchWAX | RankBoost | AdaRank | ListNet | Random Forests |
| NDCG@1 | 0.220 | 0.250 | 0.530 | 0.400 | 0.450 | **0.650** |
| NDCG@5 | 0.157 | 0.215 | 0.535 | 0.426 | 0.432 | **0.665** |
| NDCG@10 | 0.133 | 0.174 | 0.570 | 0.476 | 0.464 | **0.688** |

4x higher

All results show a statistical significance of $p<0.01$ against NutchWAX

# How much did the search effectiveness improve with the temporal features?

# With Temporal Features is Better

- ## NDCG@1
  – RankBoost 0.53  > NT RankBoost 0.44      **+9%**
  – AdaRank 0.40     > NT AdaRank 0.38      **+2%**
  – ListNet 0.45      >  NT ListNet 0.37      **+8%**
  – R. Forests 0.65  > NT R. Forests 0.55      **+10%**

- ## NDCG@10
  – RankBoost 0.57  > NT RankBoost 0.51      **+6%**
  – AdaRank 0.48     > NT AdaRank 0.47      **+1%**
  – ListNet 0.46      > NT ListNet 0.43      **+3%**
  – R. Forests 0.69   > NT R. Forests 0.65      **+4%**

# Better Results = Happier Users



sapo                                                    ×    Search the Archive

between: 01/01/1996    and: 31/12/2012                       Advanced search

Results 1 to 10 from 149,648,512

**SAPO** - Servidor de Apontadores Portugueses
**10 December, 1997** - other dates
8a2 **SAPO** - Servidor de Apontadores Portugueses    Ainda lhe restam dúvidas sobre o **SAPO** ? Esclareça-se!
c4d Novidades Novos Links , Congressos , ... Ensino e Investigação Universidades , Institutos , Escolas , ...
Comunicação Social Jornais , Rádios , Televisão , ... Entretenimento Desportos ...
http://www.sapo.pt/

**149.648.512**

**SAPO** - Portugal Online!
**8 June, 2010** - other dates
**SAPO** - Portugal Online! Saltar para: Pesquisa [1] , Lista de Serviços [2] , Notícias [3] ou Destaques **SAPO**
[4] **SAPO**.pt Pesquisa **SAPO** Web Imagens Notícias Blogs Produtos Directório PAi PBi Pesquisar: Onde:
Pesquisar Serviços Mail Blogs Carros Casas Fotos Mapas Vídeos Notícias Messenger Todo o **SAPO** ...
http://www.sapo.pt/

Eu Não Desisto: abril 2004 Archives
**17 October, 2009** - other dates
Jornal de Notícias, Minho, Braga, 17.12.2004, ou em http://jn.**sapo**.pt/2004/12/17/minho ... .blogs.**sapo**.pt/arquivo
/2004_04.html#128423  Posted by mauricio_102 at 02:46 PM | Comentários: (20 ... Portugueses ... III". 30_4_04 -
"LISTAGEM dos Artigos do Mês de Abril 2004". 28_4_04 - ""blogs.**sapo**" 25 ...
http://eunaodesisto.blogs.sapo.pt/arquivo/2004_04.html
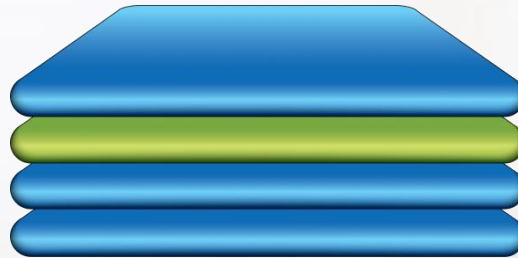
# Conclusions

# Conclusions

- Users need analytical tools for specific users + search tools for generic users.

- State-of-the-Art searching technology provides poor results.

- Temporal information intrinsic to web archives improves their search results.

- Learning-to-Rank technology greatly improves search results.

# Resources

- Public service since 2010.
  - [http://archive.pt](http://archive.pt)

- Test collection to support evaluation.
  - [https://code.google.com/p/pwa-technologies/wiki/TestCollection](https://code.google.com/p/pwa-technologies/wiki/TestCollection)

- L2R dataset for web archive IR research.
  - [http://code.google.com/p/pwa-technologies/wiki/L2R4WAIR](http://code.google.com/p/pwa-technologies/wiki/L2R4WAIR)

- All code available under the LGPL license.
  - [https://code.google.com/p/pwa-technologies/](https://code.google.com/p/pwa-technologies/)

# Thank you.



[http://archive.pt](http://archive.pt)

miguel.costa@fccn.pt