# An updated Portrait of the Portuguese Web

João Miranda, Daniel Gomes
{joao.miranda,daniel.gomes} @ fccn.pt
http://arquivo.pt

Fundação para a Computação Científica Nacional
Foundation for National Scientific Computing

UMIC
Agência para
a Sociedade
do Conhecimento

POS_Conhecimento
Programa Operacional Sociedade do Conhecimento

União Europeia
FEDER

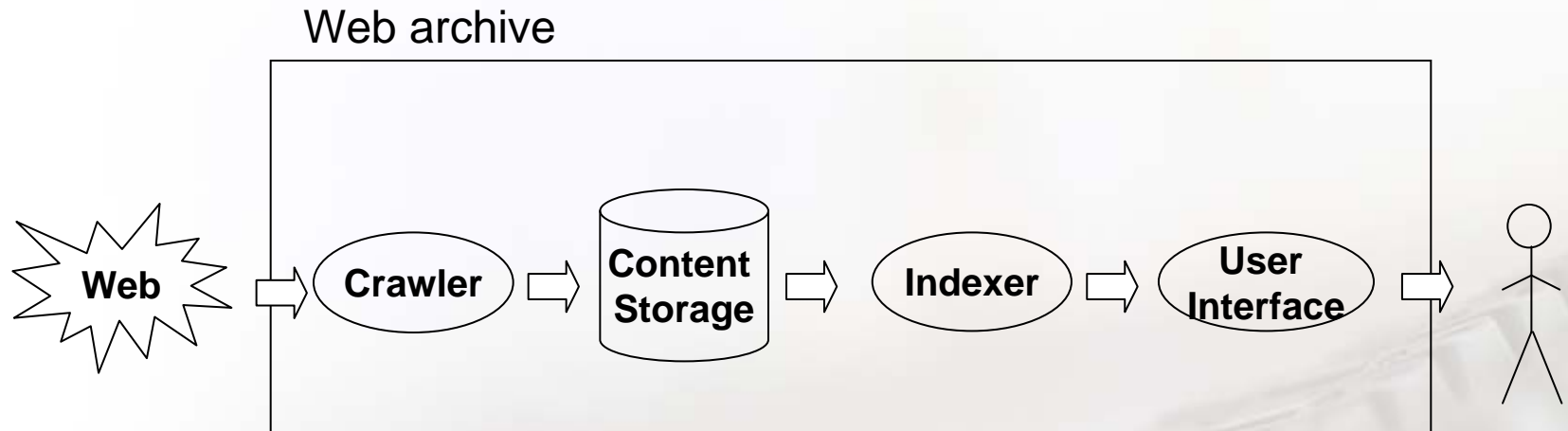- Introduction

- Methodology

- Metrics

- Conclusions

# Introduction

- The Web is a huge source of information
  - Information published exclusively on the Web
  - Information disappears

- Preservation started by the Web Archives
  - Access for future generations

- First initiative: Internet Archive

- Altavista across time

2009

- The Portuguese Web Archive

Web archive

- What is a crawler?
  - Collects contents from the Web
  - Starts from an initial set of addresses

- How does it work?
  - Iteratively downloads contents
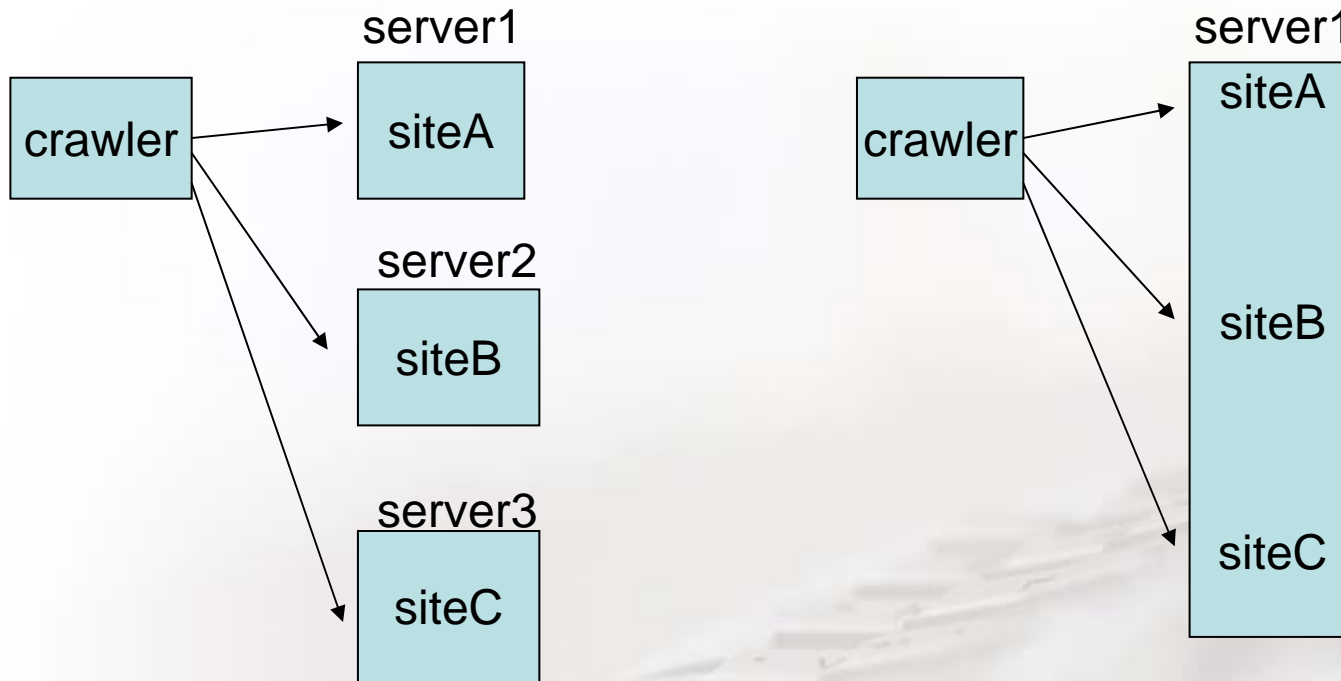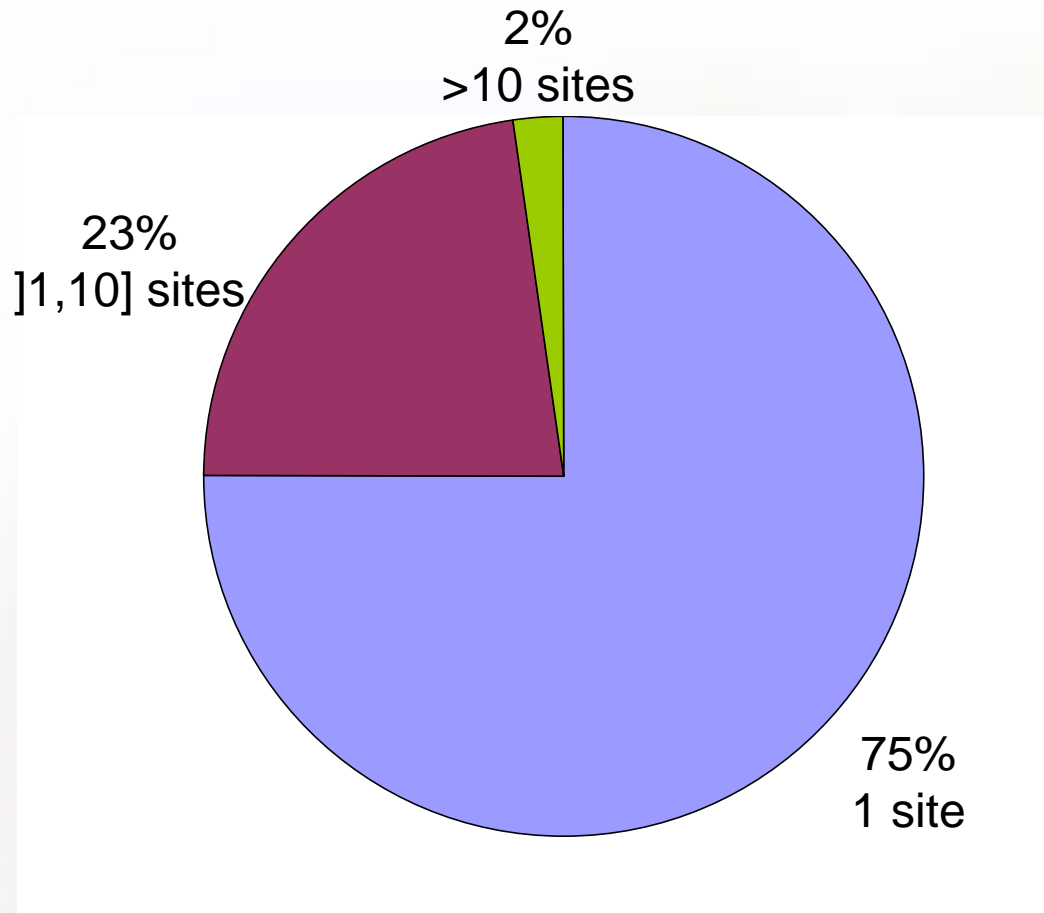  - Extracts links to find new ones

# Methodology

- Crawl of the Portuguese Web (March-May, 2008)
    - .PT domain
    - Heritrix crawler
    - 180 000 initial addresses
    - 48 million contents
    - 2.5 TB

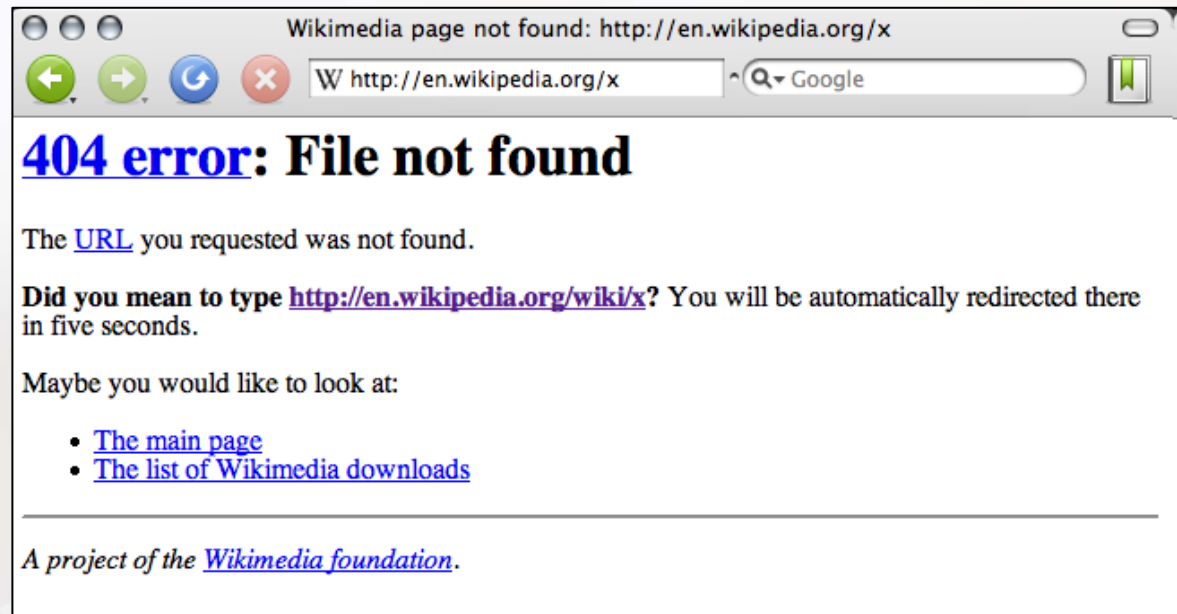- No content analysis, only log analysis

# Metrics
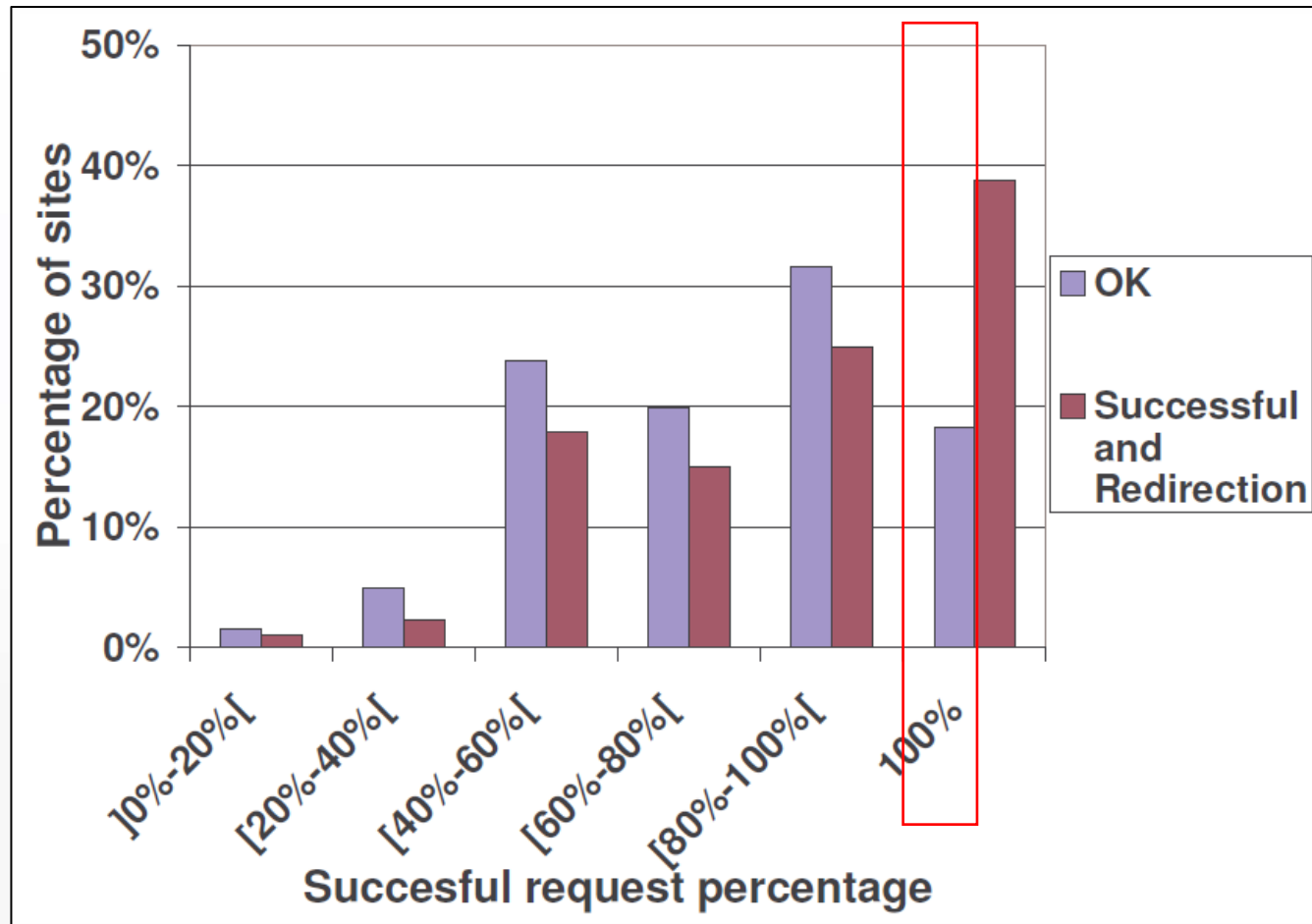
- Politeness policies for crawling

- 75% of the IP addresses host 1 site

- Quality indicator

- Large % of broken links mines trust of users

- 18% of the sites returned 100% OK responses

- Browsers or document viewers for cellphones
- Parsing and indexing for search engines

- 90% of the number of contents are html, jpeg, gif

- 69% of the amount of data are html, pdf, jpeg

|   | Media type | % contents |
|---|------------|-----------|
| 1 | Text/html  | 57.8%     |
| 2 | Image/jpeg | 22.8%     |
| 3 | Image/gif  | 9.4%      |
| 4 | Text/xml   | 1.9%      |
| - | Other      | 8.1%      |

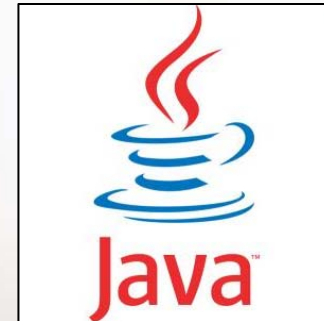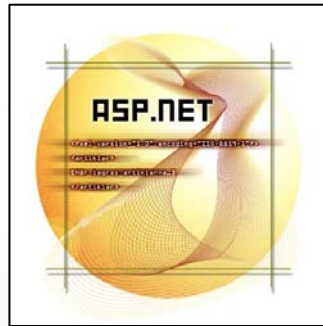|   | Media type | % amount data |
|---|------------|---------------|
| 1 | Text/html  | 35.4%         |
| 2 | App'n/pdf  | 17.9%         |
| 3 | Image/jpeg | 16.1%         |
| 4 | Text/plain | 4.2%          |
| - | Other      | 26.4%         |

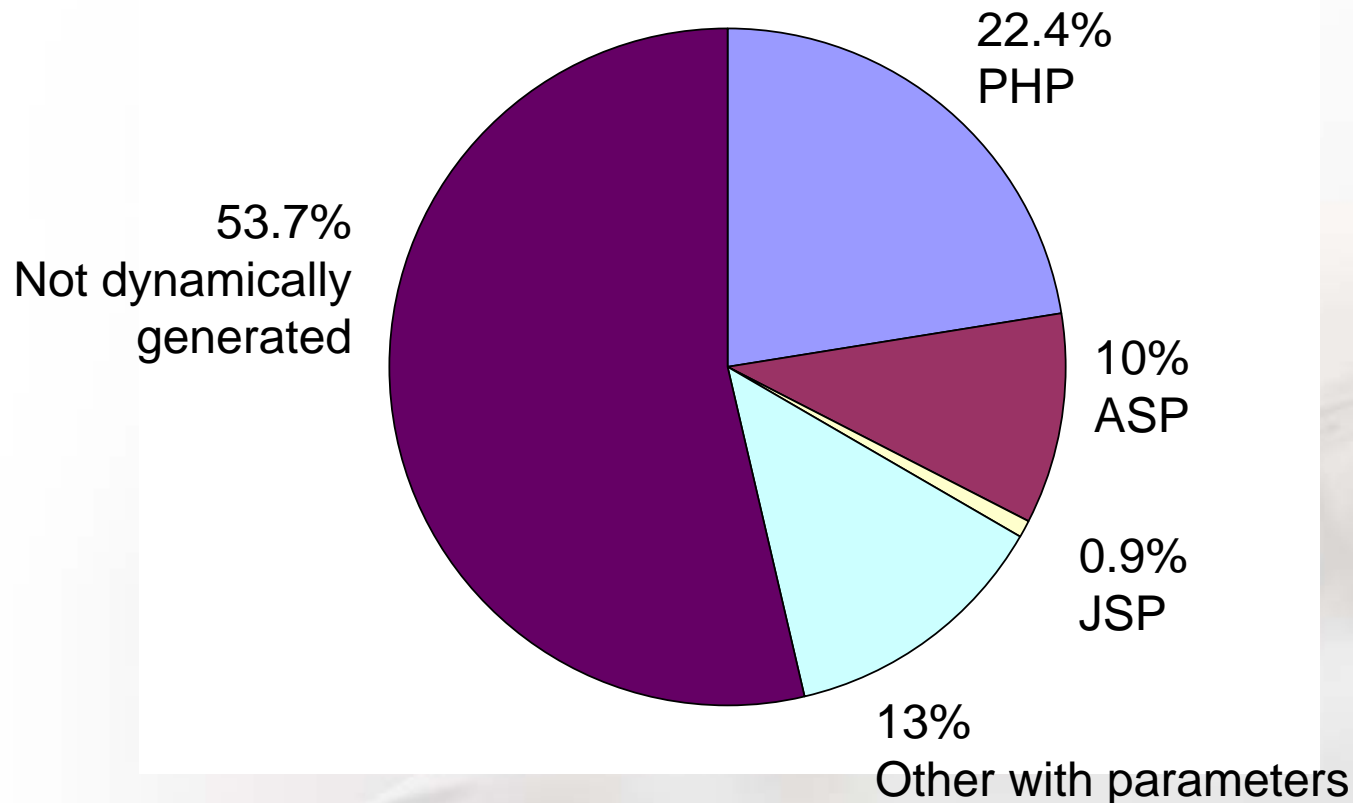- Estimate the storage resources required to create Web data repositories

- 96% lower than 128 KB

- Identify technological trends in Web publishing

- At least 46.3 % of the contents were dynamically generated



53.7%
Not dynamically
generated

22.4%
PHP

10%
ASP

0.9%
JSP

13%
Other with parameters
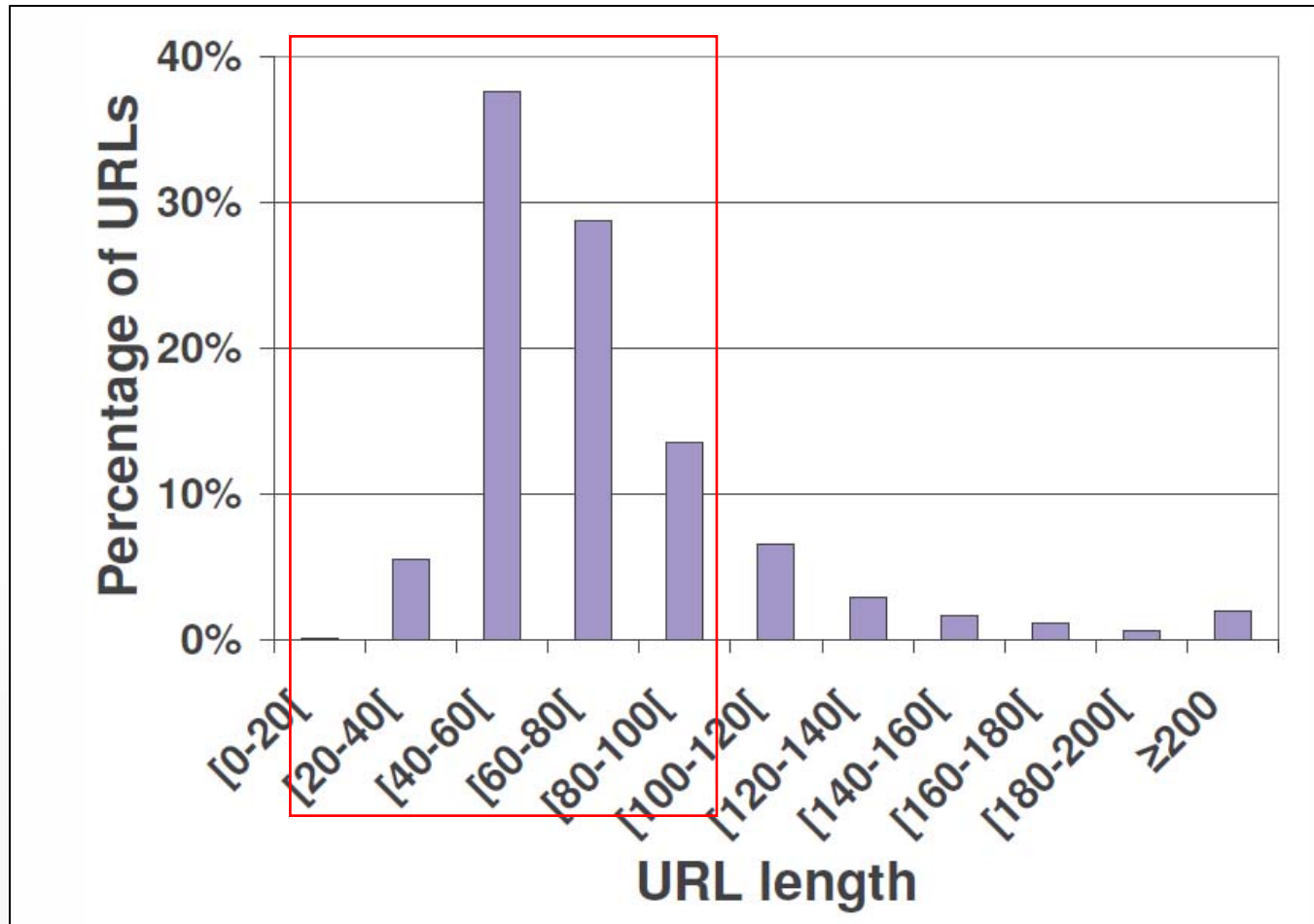
# URL length – Why?

- Influences interaction design
- Determine adequate length for input boxes that receive URLs



- How many characters should be presented on a search engine results page

- 84% lower than 100 characters

# Conclusions

- Long URL addresses

- Half of the contents are dynamically generated (mainly PHP)

- 90% of the contents are HTML, JPEG and GIF

- 69% of the amount of data are HTML, PDF and JPEG

- 96% of the contents are smaller than 128 KB

- Half of the sites present a successful response rate below 80%

- Most IP addresses host a single site

- Study trends in the evolution of web characteristics
  - João Miranda, Daniel Gomes, Trends in Web characteristics, 7th Latin American Web Congress

- Analyze metrics extracted from content and link analysis

- Anyone can contribute to preserve the Web

- Lend disk space to keep backup copies
  - Just need to install rARC
  - http://arquivo.pt/rarc

- Help required to test beta version

# Thank you.

**Logs used in this study are available for research purposes. Please contact us.**

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