



Trends in the Evolution of the Web

João Miranda, Daniel Gomes

{joao.miranda,daniel.gomes} @ fccn.pt

<http://arquivo.pt>

- The Web is a huge source of information
 - Information published exclusively on the Web
 - Information disappears
 - Only 20% of the URLs still reference a valid content after 1 year ([Ntoulas, 2004](#))
- Preservation started by the Web Archives
 - Access for future generations
 - Enables research on finding trends (temporal dimension)

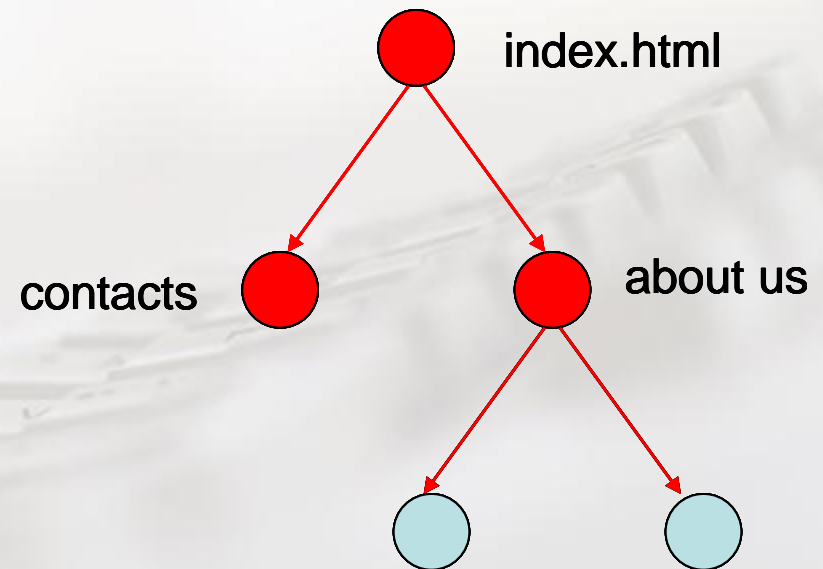
- Altavista across time

2004



The screenshot shows the Altavista search engine interface. At the top center is the Altavista logo, featuring a red stylized 'A' above the word 'altavista' in blue. Below the logo is a navigation bar with tabs for 'Web', 'Images', 'MP3/Audio', 'Video', 'Directory', and 'News'. The 'Web' tab is selected. Below the navigation bar is a search bar with a white input field and a red 'FIND' button. To the right of the search bar are links for 'Advanced Search' and 'Settings'. Below the search bar are search options: 'SEARCH: Worldwide U.S.' and 'RESULTS IN: All languages English, Spanish'. Below these options are links for 'Translate', 'Toolbar', 'Yellow Pages', 'People Finder', and 'More >>'. In the center, there is a link for 'Search for Products'. At the bottom, there are links for 'Business Services', 'Submit a Site', 'About AltaVista', and 'Help'. The footer contains the copyright notice '© 2004 Overture Services, Inc.'

- How does a crawler work?
 - Collects contents from the Web, starting from an initial set of addresses
 - Iteratively downloads contents and extracts links to find new ones



Methodology

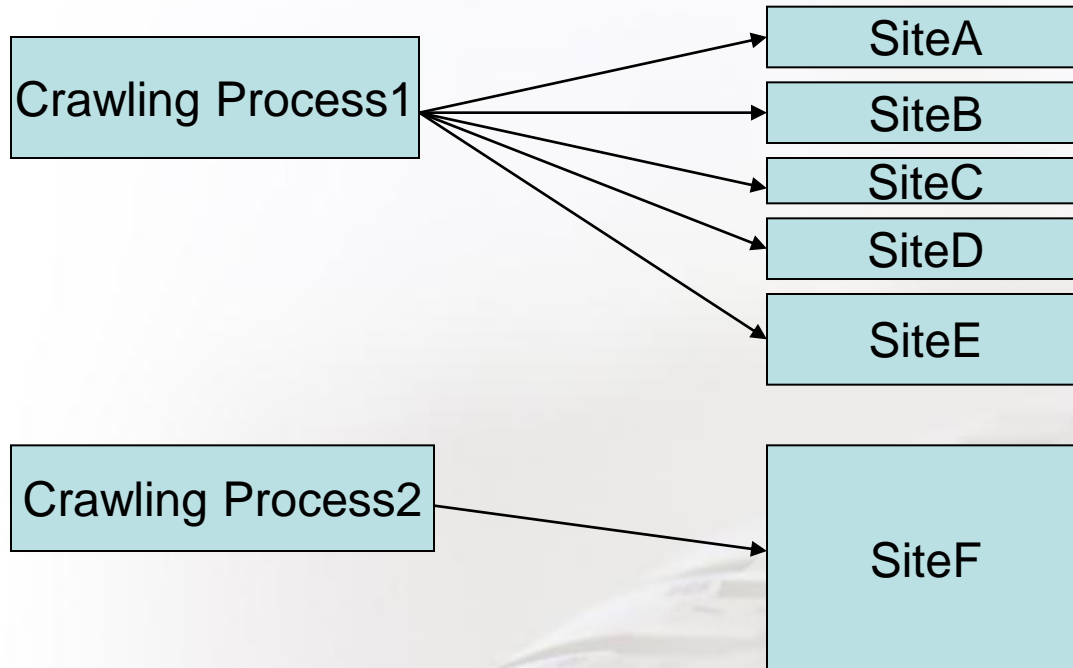
- Previous works
 - 2003 - textual media types
 - 2005 - all media types
- Ours
 - 2008 - all media types, textual media types

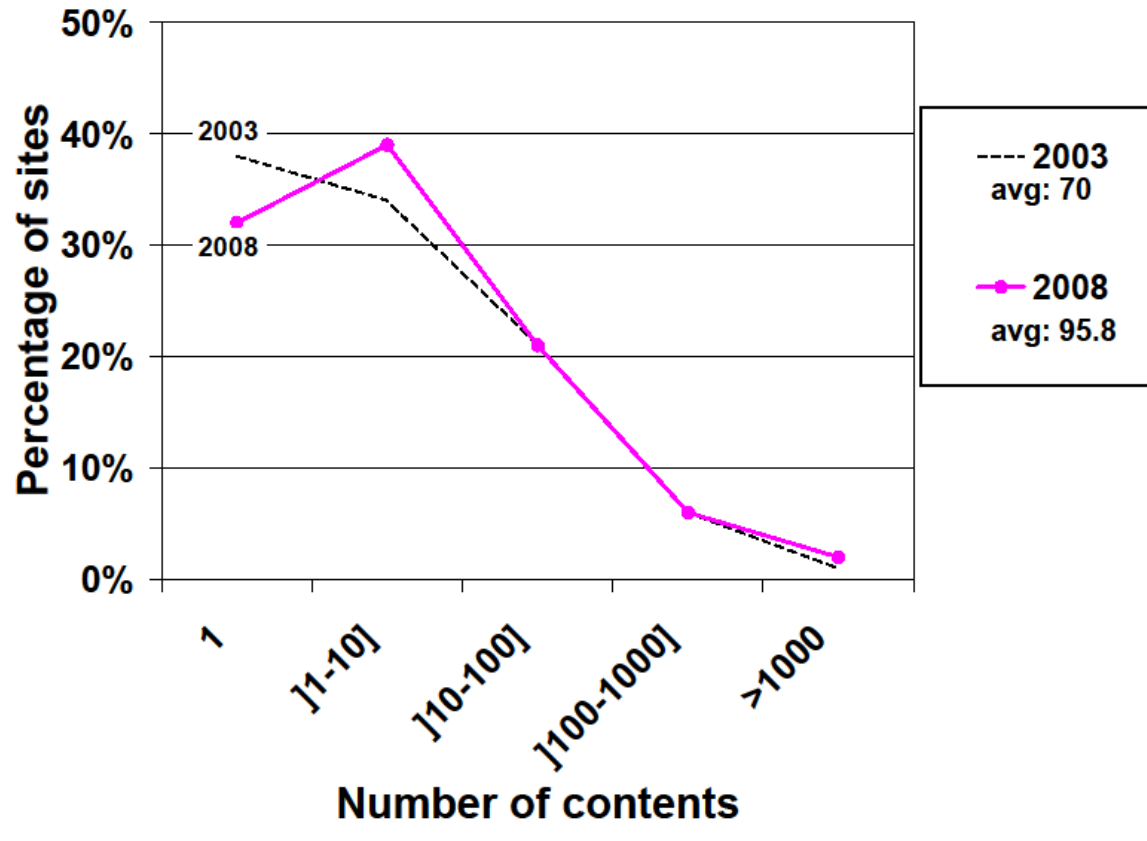
Trends

- Reflect the changes in the characteristics of contents and sites:
 - Are sites larger?
 - Which media types are prevalent?
 - Has the average size of contents grown?

- Why is it useful to analyse trends?
 - Improve the processment of web data

- Efficiently partition a large data set of URLs across several crawling processes





- Average number of contents per site increased from 70 to 96

- Browsers or document viewers for cellphones



All Media types - Results

	Media type	% contents 2005	% contents 2008	Trend $\left(\frac{\% \text{ contents 2005} - \% \text{ contents 2008}}{\% \text{ contents 2005}}\right) \times 100$
			95.2%	
			90.1%	
1	Text/html	61.2%	57.8%	-5.5%
2	Image/jpeg	22.6%	22.8%	+1.2%
3	Image/gif	11.4%	9.4%	-17.4%
4	Text/pdf	1.6%	1.9%	+18.5%
-	Other	3.2%	8.1%	-

- Decrease in *text/html* and *image/gif*
- Increase in *image/jpeg* and *text/pdf*
- Media type prevalence is more spread

Textual Media types - Results

	Media type	% contents 2003	% contents 2008	Trend $\left(\frac{\% \text{ contents 2003} - \% \text{ contents 2008}}{\% \text{ contents 2003}}\right) \times 100$
1	Text/html	96.0%	93.9%	-2.1%
2	App'n/pdf	1.9%	3.0%	+57.6%
3	Text/plain	1.0%	1.6%	+58.5%
4	App'n/x-shockwave-flash	0.5%	1.2%	+115.8%
-	Other	0.7%	0.3%	-

- *Text/html* lost presence to other formats
- Increase in *app'n/pdf*, *text/plain*, *app'n/shockwave-flash*

- Estimate the storage resources required to create Web data repositories

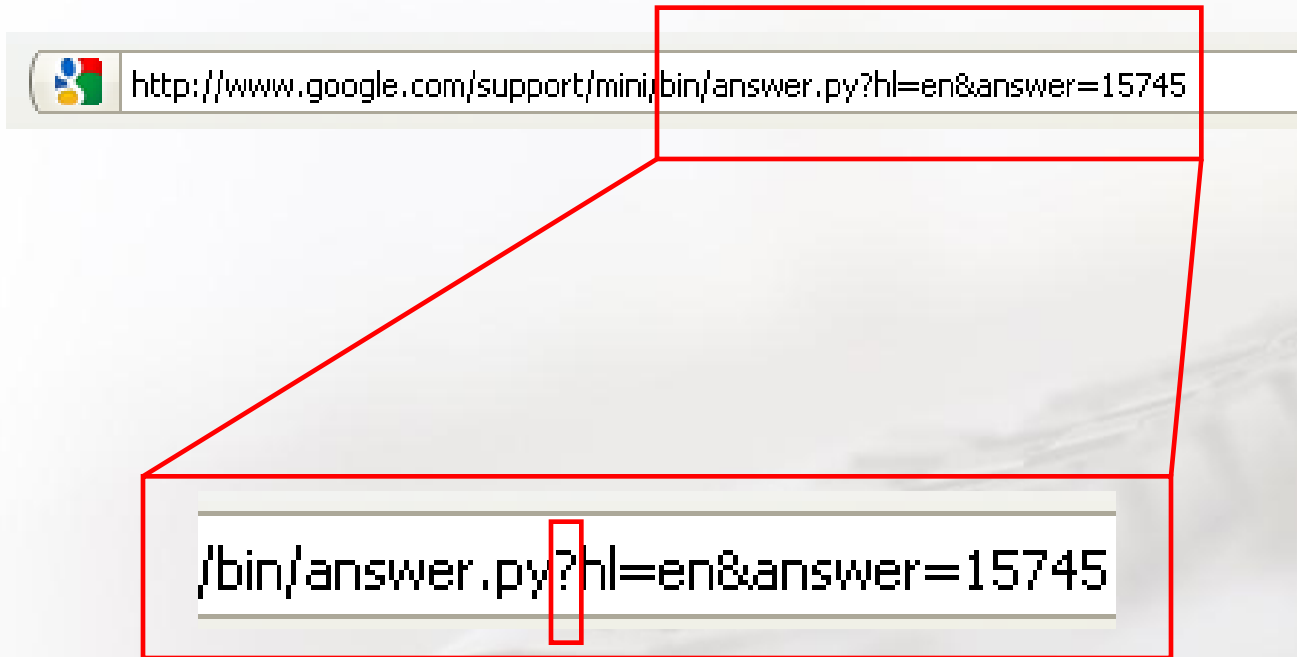


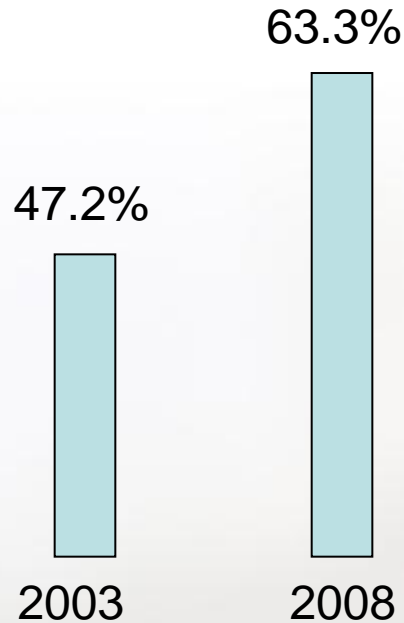
Content size - Results

Media type	Avg Size 2003	Avg Size 2008	Trend $\left(\frac{\text{Avg size 2003} - \text{Avg size 2008}}{\text{Avg size 2003}} \right) \times 100$
Text/html	21 KB	30 KB	+45.9%
App'n/pdf	207 KB	252 KB	+21.6%
Text/plain	11 KB	44 KB	+58.5%
App'n/x-shockwave-flash	44 KB	90 KB	+115.8%
...			
powerpoint	1055 KB	500 KB	-52.6%
Text/rtf	476 KB	143 KB	-70.0%

- Content size grew but for some types decreased

- Identify technological trends in Web publishing





- The percentage of URLs containing parameters increased

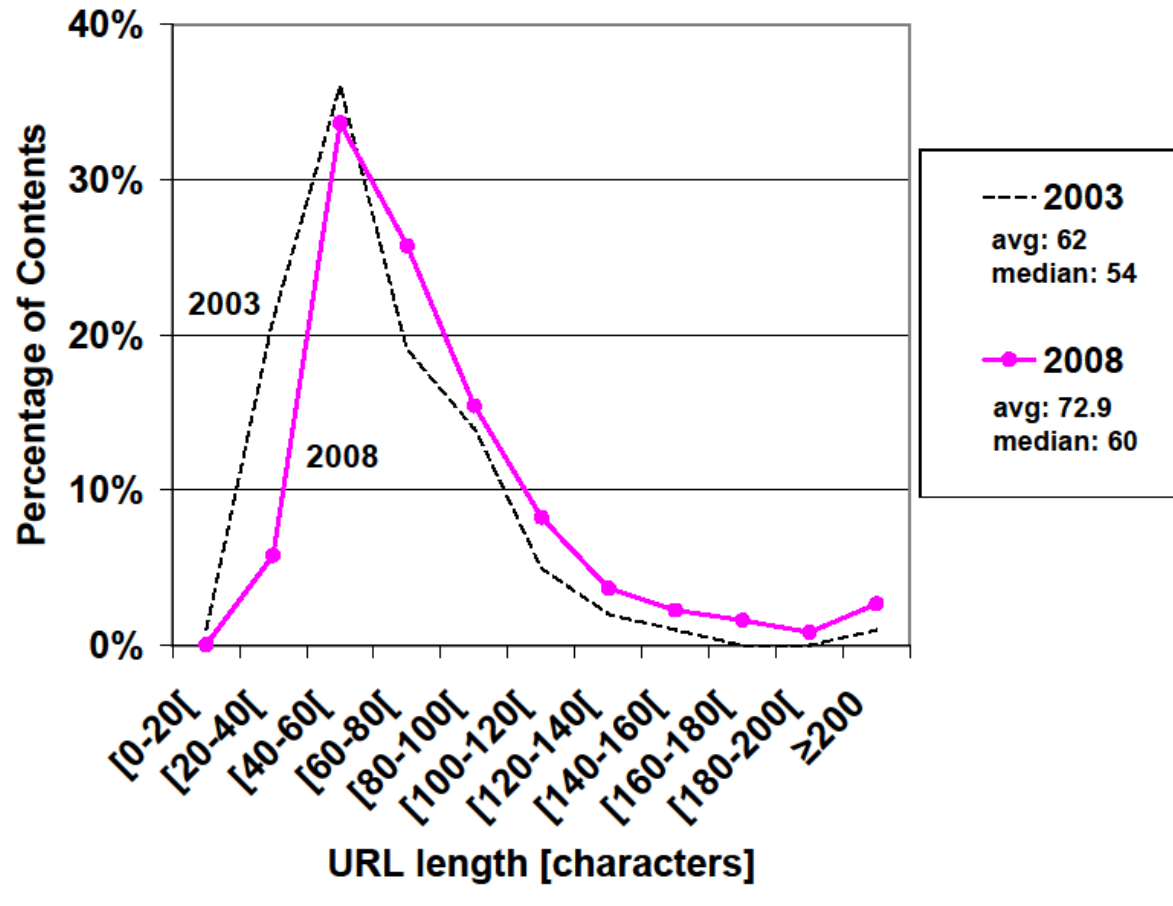
- Influences interaction design
- Determine adequate length for input boxes that receive URLs



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

[Portuguese Web Archive — Arquivo da Web Portuguesa](#) - 3 visits - Jul 13
 The **Portuguese Web Archive** (PWA) is a National Foundation for Scientific Computing (FCCN) project whose main objective is to preserve the information ...
[arquivo-web.fccn.pt/portuguese-web-archive-2?set...en](#) - [Cached](#) - [Similar](#) - [🗨](#) [📄](#) [🗑](#)

URL length - Results



- Average URL length increased from 62 to 73 after 5 years

Conclusions

- After 5 years what were the changes?
 - URL length increased slightly but the average content size increased significantly
 - Sizes did not grow for all media types

- After 5 years what were the changes?
 - Dynamically generated contents became widely used
 - Number of contents per site increased
 - HTML, GIF or JPEG became prevalent

- Testing of developed systems
- Research and development projects
- Crawl logs available for research purposes
<http://archive.pt>



Thank you.

<http://archive.pt>