# Search the Past Web with Arquivo.pt

Forum on the archiving of the
Polish Internet resources

daniel.gomes@fccn.pt

# Who are we?

**Free online** service to research the Past Web

Preserves **publicly accessible** information related with:

- Portugal
- **Research** and **Education** (international)

Governmental service provided by Foundation for Science and Technology (Portugal)

A digital research infrastructure

# Arquivo.pt begun in 2007

# Arquivo.pt is used word-wide

| | Country | Users | % Users |
|---|---|---|---|
| 1. | Portugal | 46,891 | 46.56% |
| 2. | United States | 26,373 | 26.19% |
| 3. | Brazil | 2,266 | 2.25% |
| 4. | Russia | 2,234 | 2.22% |
| 5. | United Kingdom | 2,231 | 2.22% |
| 6. | Japan | 2,172 | 2.16% |
| 7. | Canada | 1,237 | 1.23% |
| 8. | Mozambique | 1,213 | 1.20% |
| 9. | India | 902 | 0.90% |
| 10. | Germany | 894 | 0.89% |

- 53% of users are international
- User Interfaces and documentation also in English
- Combined with Google Translate enables **cross-lingual access** to preserved content

# Arquivo.pt preserves national and international historical web content



**nautilus.fis.uc.pt**- **1993**
(oldest page)



**spacelink.nasa.gov** – **1992**
(oldest image)

# How to **search** and **access** the preserved information?

# Search **texts** from the past in any language

# Search the **history** of a web address

# Browse [bbc.co.uk](bbc.co.uk) 10 years ago

# Search **images** from the past

# **A**pplication **P**rogramming **I**nterfaces
# to develop new applications

**Arquivo.pt APIs**

Your new app



## APIs

Daniel Gomes edited this page 11 days ago · 33 revisions

### APIs specific to Arquivo.pt that enable the full exploration of its functions

- Arquivo.pt API (Full-text & URL search)
- Image Search API v1.1 (beta version)

### APIs based on international standards to enable interoperability among web archives and code reuse

- CDX-server API (URL search)
- Memento API (URL search)

arquivo.pt/api

# Arquivo.pt API in action: **TellMeStories.pt**



**Automatically** generates **narratives** about any subject based on online news from the Past.

Winners of the Arquivo.pt Award 2018.

# Collections and automatic selection

# Public list of Arquivo.pt **collections**



List of the collections preserved by Arquivo.pt (publicly available)

File   Edit   View   Insert   Format   Data   Tools   Add-ons   Help   Last edit was made on 10 September by SAW - FCCN Serviços e Arquivos Web

A2   =IMPORTRANGE("https://docs.google.com/spreadsheets/d/1hfLyq9N2ZJEd1hra2V0OPI24jTy42ByR6M667rG-oHo/edit","Summary!A2:G")

| | Collection ID | Description | Collection Start date | Collection End date | Total files (Number of URLs collected) | Total seeds (Preserved sites) | Total volume of WARCS files (TB) | Collection in Production? (URL, Page and Image indexed) |
|---|---|---|---|---|---|---|---|---|
| 2 | AWP1 | 1st Complete crawl of the Portuguese web, mainly from the .PT domain, in 2008. | 2008-02-12 | 2008-03-06 | 56,046,288 | 154,787 | 1.60 | ☑ |
| 3 | AWP2 | 2nd Complete crawl of the Portuguese web, mainly from the .PT domain, in 2008. | 2008-03-11 | 2008-05-30 | 48,718,404 | - | 1.60 | ☑ |
| 4 | AWP3 | 3rd Complete crawl of the Portuguese web, mainly from the .PT domain, in 2008. | 2008-10-21 | 2008-12-10 | 51,863,006 | 193,294 | 2.00 | ☑ |
| 5 | AWP4 | 4th Complete crawl of the Portuguese web, mainly from the .PT domain, in 2009. | 2009-05-01 | 2009-05-31 | 68,776,707 | 366,880 | 2.50 | ☑ |
| 6 | AWP5 | 5th Complete crawl of the Portuguese web, mainly from the .PT domain, in 2009. | 2009-10-01 | 2009-10-31 | 119,135,566 | 373,323 | 3.80 | ☑ |
| 7 | AWP6 | 6th Complete crawl of the Portuguese web, mainly from the .PT domain, in 2009. | 2009-12-01 | 2009-12-31 | 118,810,364 | 340,018 | 3.50 | ☑ |
| 8 | AWP7 | 7th Complete crawl of the Portuguese web, mainly from the .PT domain, in 2010. | 2010-05-01 | 2010-05-31 | 87,988,812 | 389,957 | 2.90 | ☑ |
| 9 | AWP8 | Incremental crawl of the Portuguese web web, mainly from the .PT domain, in 2010. The AWP8 crawl is incremental because it was performed using DeDuplicator (http://landsbokasafn.github.io/DeDuplicator/) taking the content of AWP7 as baseline. Thus, the files that remained unchanged from the AWP7 complete crawl were not archived (duplicated) on the AWP8 incremental crawl. | 2010-08-01 | 2010-08-31 | 75,771,317 | 411,562 | 1.90 | ☑ |

[arquivo.pt/collections](arquivo.pt/collections)

# Developed a low-cost methodology to automatically **select** and preserve online information about given topics

**Preserving Websites Of Research & Development Projects**

Daniel Bicho
Foundation for Science and Technology:
Arquivo.pt
Av. do Brasil, 101
1700-066 Lisboa, Portugal
daniel.bicho@fccn.pt

Daniel Gomes
Foundation for Science and Technology:
Arquivo.pt
Av. do Brasil, 101
1700-066 Lisboa, Portugal
daniel.gomes@fccn.pt

International Conference on Digital Preservation 2016, available at arquivo.pt/publications

# Experiments to preserve European Union web content

# Tools to support web preservation

# **Suggest** websites to be preserved

- Public form
- Any domain
- Any language
- Submit home page
➢ arquivo.pt/suggest

## Suggest websites to be preserved

International websites are welcome!

sawfccn@gmail.com Switch accounts

*Required

Email *

Your email address

Website (1 per line) *

Your answer

Description of the Websites

Your answer

# Missing content? **Complete the page**!



Missing files obtained from the live-web/archives are then integrated to improve the quality of the archived page

➢ Collaborative curation: **Give it a try!**

# Don't kill your historical website!
## Preserve it in the **Arquivo.pt Memorial**: [arquivo.pt/memorial](arquivo.pt/memorial)

# [arquivo.pt/arquivo404](arquivo.pt/arquivo404)
## **fix** the **broken links** in your website!



Page not found at live website

Page available at Arquivo.pt

# Explore the Past Web

# Meet the Winners of the Arquivo.pt **Annual** Awards





Support of the President of the Republic of Portugal

➢ arquivo.pt/awards

# Create Web Design museums

# Attend Training courses about web preservation and research

**New ways of searching the past** (module A)

– Any Internet user

**Publish preservable information on the web** (B)

– Web authors

**Automatic processing of information preserved from the Web** (C)

– Developers, Researchers

**Do-it-yourself web archiving** (D)

– Researchers, Librarians, Archivists

**Videos in English also available!**



arquivo.pt/training

# Challenges and recommendations to start a web archive

# Challenges

**Hire and maintain skilled human resources**

- Web archiving requires specialists but is not taught at any course
- Web Archiving is Big Data, must compete with the Internet giants (e.g. Google, Facebook) to hire

**Costs of dissemination**

- Difficult "to sell" something for free in the attention economy

**Lack of awareness about the value of online information**

- Information that rules our current World is published exclusively online
- "Information is power", "Data is the new oil", "Look into the past to see the future", …
- However, societies passively lose most of their online information

# Lack of awareness about the value of online information

Let's measure the value of preserving online information or

the waste of not doing it.

Investment made on creating the online information preserved in Arquivo.pt

Design/development: 27 Msites x 5 600€/website = 151 000 M€

Content: 8 000 Mpages x 68€/page = 544 000 M€

Value of the preserved heritage =
	**695 000 M€**

Gross Domestic Product of Portugal in 2020 =
	**202 000 M€**

Can we afford to waste all this investment?

# Recommendations to start a web archive

## No part-time web archivists

- Web archiving is complex, requires full dedication
- Start with a small but autonomous web archiving team so that it does not have to permanently compete for resources

## Use existing tools and services
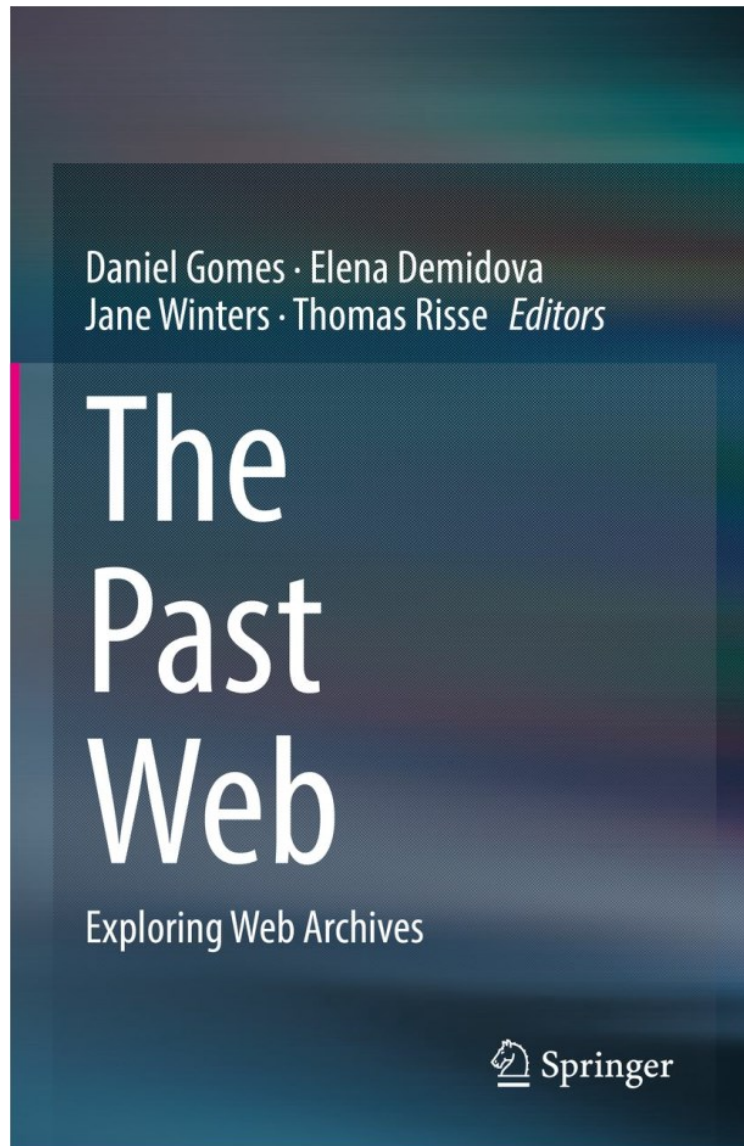
- Webrecorder.net, Archive-it

## Listen to your users

- Don't guess what they need, hire User Experience researchers

## Start small but perform the full worklow

- Short list of websites and perform the full preservation cycle: collect, store, provide access and disseminate service
- Doing a little is infinitevely more than to do nothing

# Read this book ☺



Preprint version available in open access at:
[arquivo.pt/book](arquivo.pt/book)

# Part 6:
# The future of web archiving

Julien Masanès: the "father" of web archiving in Europe

# Keep in touch!

**Subscribe our channels**:

➢ [arquivo.pt/subscribe](arquivo.pt/subscribe): mailing list in English

➢ [arquivo.pt/news](arquivo.pt/news): social networks and videos