# Arquivo.pt behind the curtains

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 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)

# Arquivo.pt supports 13 services

**Services catalog of the Arquivo.pt web archive**

Last updated on February 14th, 2023 at 03:35 pm

## Public services

- Search and access web-archived data since the 1990s
- Application Programming Interfaces (APIs)
- Suggest websites to be preserved
- SavePageNow: immediately archive web pages
- Integration of historical web data collections
- Training on web preservation
- Open data listing archived web information on various topics
- CitationSaver: extracts links from documents and archives the correspondent web pages
- Arquivo404: presents web-archived pages instead of "pages not found"

## For partners

- Memorial preserves your old website information before deactivating it
- High-quality archive of websites (on-demand)
- Creation of collections and thematic exhibitions
- Itinerant exhibition of posters at your facilities

arquivo.pt/catalog

# Our **successful** official History!

2007: Start of the Portuguese Web Archive project

2008: 1st crawl of the Portuguese web

2009: Recommendations for publishers

2010: Public prototype of the search and access system

2011: Open source available

2012: Application Programming Interface

2013: Decree-Law mandating our web archiving activities

2017: Training program

2018: Arquivo.pt Award

2019: Arquivo.pt Memorial

2021: Image search service

**2022: Arquivo.pt was considered the best Digital Service of 2022!**

# Behind the curtains

# The adversities

2007: Start of the Portuguese Web Archive project

*9/2013:* Total collapse due to hardware malfunction
- Irrecoverable data loss of 17% (17 TB)
- Crawling interruptions
- Suspension of search service

*2014 - 2016:* Recovery and improving robustness
- Ground-zero recovery
- 9,2 hours of downtime in 1 year

2022: Service availability of 99.894%

High staff rotation: 23 members in 15 years with average team size of 4
- On average: each 8 months, one member leaves

# Workload (only human users)

| | Average month | Total 2022 |
|---|---|---|
| Users | 10 453 | 125 439 |
| Sessions | 14 373 | 172 478 |
| PageViews | 188 253 | 2 259 039 |

Exciting figures for an archive **but** humble for a web service.

Steady growth since 2010

# Arquivo.pt is a medium-size web archive



Arquivo.pt in numbers

Last updated on March 28th, 2023 at 04:39 pm

Statistics and curiosities regarding the Arquivo.pt service (last update in January 2022)
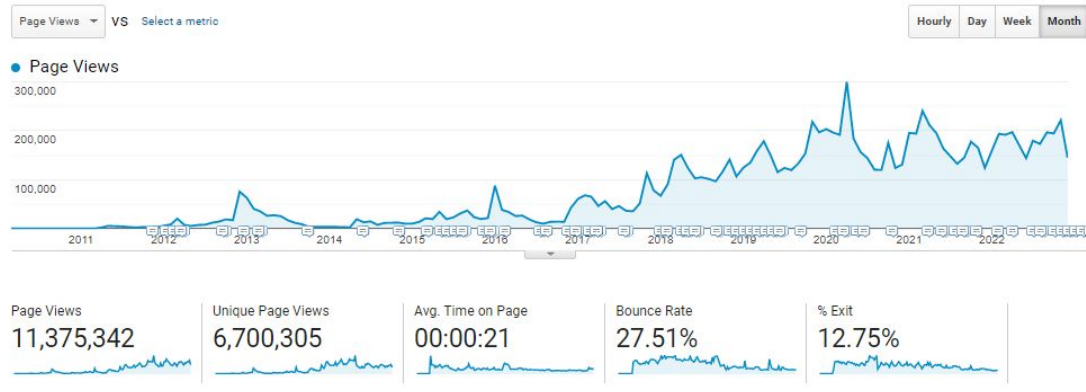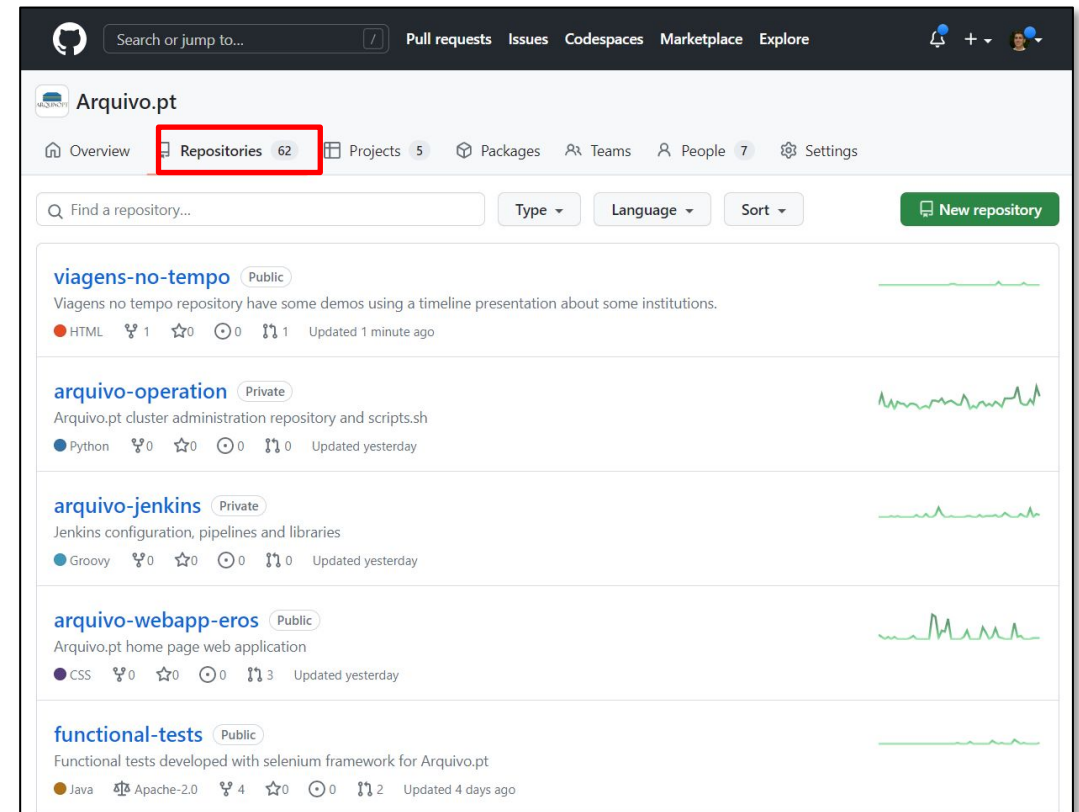
Amount of preserved data

- 17 716 million files
- 32 million websites
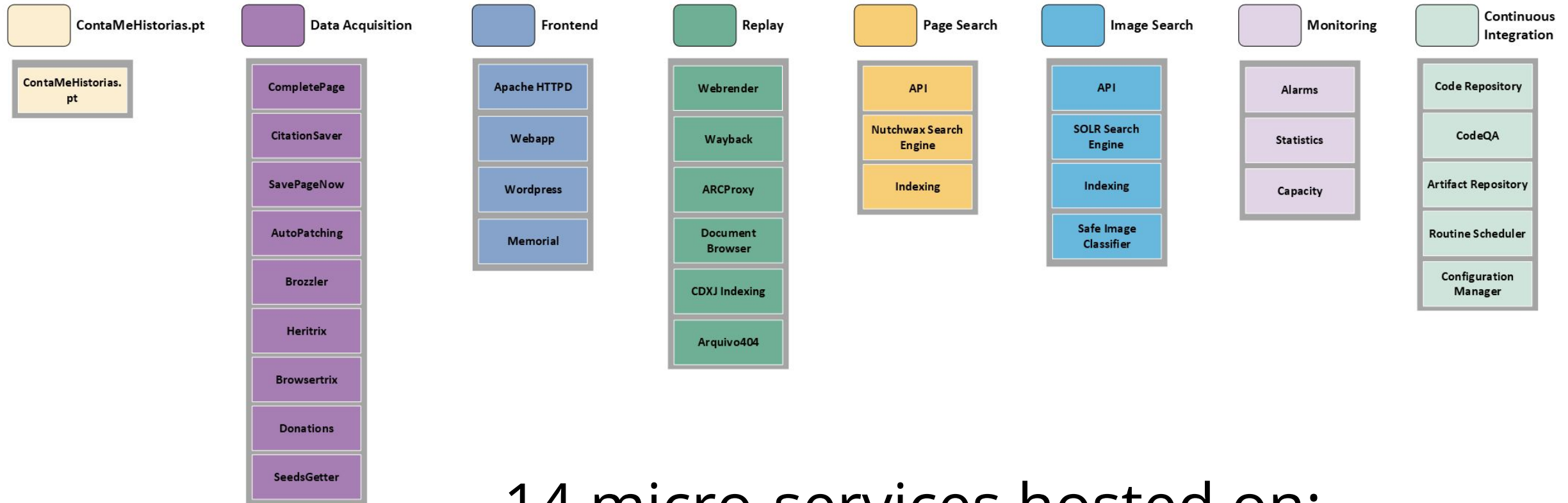- 876 TB (compressed format)

Hardware and Software

- 77 servers mounted on 167 Rack Units
- 18 TB of RAM memory
- 2 180 vCPUs
- 1 234 hard drives (5.18 PB)

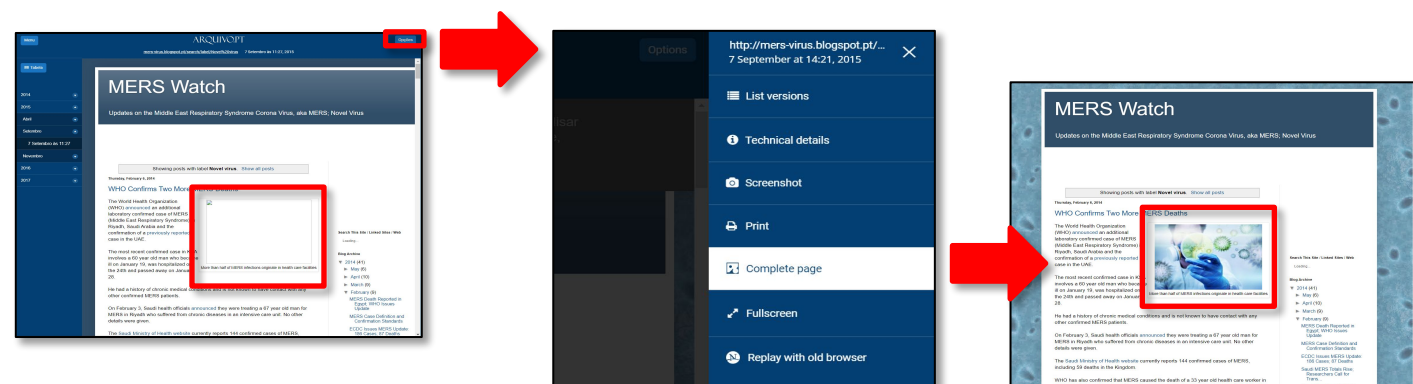## Supported by a lot of software

# Software Architecture of Arquivo.pt

**Systems**

| ContaMeHistorias.pt | Data Acquisition | Frontend | Replay | Page Search | Image Search | Monitoring | Continuous Integration |
|---|---|---|---|---|---|---|---|
| ContaMeHistorias.pt | CompletePage | Apache HTTPD | Webrender | API | API | Alarms | Code Repository |
| | CitationSaver | Webapp | Wayback | Nutchwax Search Engine | SOLR Search Engine | Statistics | CodeQA |
| | SavePageNow | Wordpress | ARCProxy | Indexing | Indexing | Capacity | Artifact Repository |
| | AutoPatching | Memorial | Document Browser | | Safe Image Classifier | | Routine Scheduler |
| | Brozzler | | CDXJ Indexing | | | | Configuration Manager |
| | Heritrix | | Arquivo404 | | | | |
| | Browsertrix | | | | | | |
| | Donations | | | | | | |
| | SeedsGetter | | | | | | |

# 14 micro-services hosted on:
-  8 systems
-  35 components
-  88 servers

**Data Acquisition**

CompletePage

CitationSaver

SavePageNow

AutoPatching

Brozzler

Heritrix

Browsertrix

Donations

SeedsGetter

pywb

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

**Data Acquisition**

- CompletePage
- CitationSaver
- SavePageNow
- AutoPatching
- Brozzler
- Heritrix
- Browsertrix
- Donations
- SeedsGetter

arquivo / **CitationSaver**

## ARQUIVO.PT

**Citation***Saver*  **Preserves citations to online resources**

Documents cite online content that quickly disappear.

CitationSaver preserves the content of cited links (e.g. web pages cited in a book) so that they can be later recovered from Arquivo.pt.

Submit a document and CitationSaver will preserve its cited links:

| URL | File | Text |
|---|---|---|

Insert the URL to the document:

Accepted formats: PDF, TXT. Max size: 100 MB

Your email (optional):

**Save cited links**

arquivo.pt/citationsaver

arquivo.pt/savepagenow

**Data Acquisition**

- CompletePage
- CitationSaver
- SavePageNow
- AutoPatching
- Brozzler
- Heritrix
- Browsertrix
- Donations
- SeedsGetter

ARQUIVO.PT

Libraries to convert files to WARC format

**ARQUIVO.PT**

sobre.arquivo.pt > Collaborate > Donate historical web content

**Donate historical web content**

Last updated on June 21st, 2021 at 12:13 pm

Contribute to preserve web knowledge by donating historical web content so that Arquivo.pt can preserve it.

Arquivo.pt periodically crawls the Portuguese web since January 2008. However, content previously published online must be gathered from external sources to be preserved by our system.

If you have web content interesting for the Portuguese community or science and technology and want to contribute for its preservation, please contact us.

We consider that all content published on sites under the .PT domain is part of the Portuguese web and must be preserved. However, contents hosted under other domains, considered of interest for the Portuguese community, will also be accepted.

**Should I supply only old content?**

We are interested in receiving contents that are no longer available online, independently from their publication date.

The web is extremely dynamic and the lifetime of most contents is very short.

Thus, many contents are lost because they become unavailable before we can gather them, even though we perform periodic crawls of the Portuguese web.

Backups of Portuguese websites are a good example of contents that may be provided.

**How can I supply web content?**

Arquivo.pt stores the archived content using the ARC format. Ideally, content should be supplied using this format.

However, it's natural that most people do not use it to keep their files. Therefore, we accept Portuguese web contents kept in any format.

Later, the Arquivo.pt team will convert them to the ARC format so that they can be integrated into our system.

## arquivo.pt/donate

**Data Acquisition**

- CompletePage
- CitationSaver
- SavePageNow
- AutoPatching
- Brozzler
- Heritrix
- Browsertrix
- Donations
- SeedsGetter

ARQUIVO.PT

Queries to external APIs and crawl log analysis

## 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

**Data Acquisition**

| | |
|---|---|
| **CompletePage** | ⌾ pywb |
| **CitationSaver** | ARQUIVO.PT    ⌨ arquivo / **CitationSaver** |
| **SavePageNow** | ⌾ pywb |
| **AutoPatching** | ⌾ pywb |
| **Brozzler** | ⌨ internetarchive / **brozzler** |
| **Heritrix** | ⌨ internetarchive / **heritrix3** |
| **Browsertrix** | Browsertrix |
| **Donations** | ARQUIVO.PT    Libraries to convert files to WARC format |
| **SeedsGetter** | ARQUIVO.PT    Queries to external APIs and crawl log analysis |

Frontend

Apache HTTPD

Webapp

Wordpress

Memorial

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

Frontend

Apache HTTPD

Webapp

Wordpress

Memorial

**Replay**

- Webrender
- Wayback
- ARCProxy
- Document Browser
- CDXJ Indexing
- Arquivo404

{.JS} JavaScript

memento

**arquivo.pt/arquivo404** is powered by Memento API
presents web-archived pages instead of "pages not found"

Page not found at live website

Page available at Arquivo.pt

# Image Search

**Apache Solr** {JSON}

## ImageSearch API v1.1 (beta)
Daniel Gomes edited this page on Dec 16, 2022 · 16 revisions

The ImageSearch API allows keyword to image search and access to preserved web content and related metadata.
The API returns a **JSON** object.
**EndPoint:** https://arquivo.pt/imagesearch

### Changelog

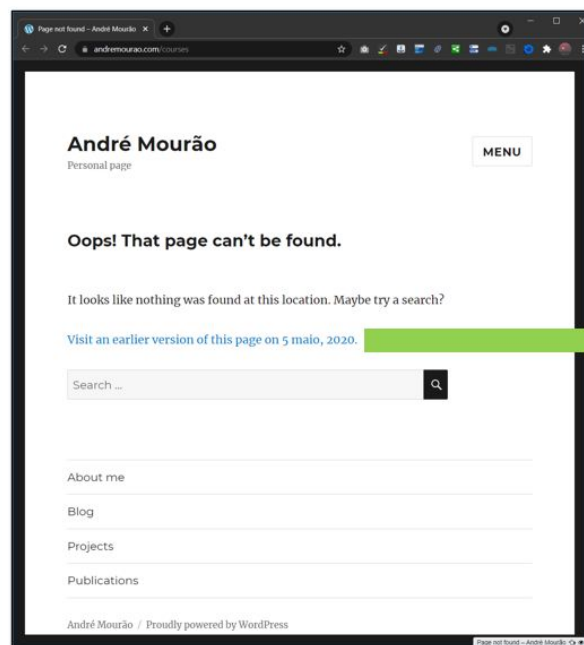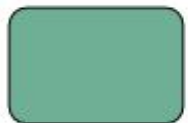Note that API version changes are **breaking**. The API endpoint does not change, but only the latest version (currently 1.1) is available. Fields removed between version are not accessible anymore due to changes in our search backend.

**1.1 (March 2021):**

- Improved search backend with 584+ million images referenced in 1800+ million pages;
- Subdomain filtering using the `siteSearch:` parameter is not expanded by default, you must use wildcards to add subdomains to the filter. Check how to do this `siteSearch:` line of the Request Parameter section;
- Added the `imgCaption` multivalued field;
- `imgAlt`, `imgTitle` and `collection` fields are now multivalued;
- Added additional image metadata fields (`imageMetadataChanges`, `pageMetadataChanges`, `matchingImages`, `matchingPages`);
- Some secondary fields were removed (`pageProtocol`, `imgThumbnailBase64`, `imgSrcURLDigest`).

**Pages** 30

Find a page...

- Home
- APIs
- Arquivo.pt API
- Arquivo.pt image search reposi...
- Arquivo.pt in a nutshell: overvi...
- Arquivo404: broken link fixer
- Compile
- ConfigureSearch
- ImageSearch API v1 (beta)
- ImageSearch API v1.1 (beta)
    - Changelog
        - 1.1 (March 2021):
        - 1.0:
    - Request Parameters
        - Search for terms

## API

## SOLR Search Engine

## Indexing

## Safe Image Classifier

### Search **images** from the past

**Apache Solr**

Image Search

API

SOLR Search Engine

Indexing

Safe Image Classifier

Apache Solr

{JSON}

Apache Solr

hadoop

arquivo / image-gpu-classifier

494 internal monitors/alarms
140 external monitors/alarms

**Continuous Integration**

| | |
|---|---|
| Code Repository | GitHub |
| CodeQA | sonarcloud / SAUCE LABS |
| Artifact Repository | APACHE HTTP SERVER PROJECT |
| Routine Scheduler | Jenkins |
| Configuration Manager | ANSIBLE |

## 58 regression tests

- 812 environments tested on Sauce labs (combinations of Firefox, Chrome and Safari on Windows, MacOS, IOS and Android)

## 50 operational routines (Jenkins)

## 50 orchestration playbooks (Ansible)

# Third party service: "Tell me stories"



Narrative: new Arquivo.pt function for free!

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

Winners of the Arquivo.pt Award 2018.

☐ Advantage of exposing public APIs!

# Communication and documentation tools



## More tools, less communication

There are over 25 official communication tools used within our organization.

We use 5 tools in Arquivo.pt and it is more than enough

- GitHub
- Email
- Cell phone
- WhatsApp
- Google Drive

# Hardware Architecture of Arquivo.pt

**Third Party Management**
- Alarms (external)

**HTTP Server**
- Wordpress
- Memorial

**Systems**
- Data Acquisition
- Replay
- Page Search
- Image Search
- Frontend
- Monitoring
- Continuous Integration
- ContaMeHistorias.pt

**Broker**
- Wayback
- ARCProxy
- Webrender
- Webapp
- Apache HTTPD
- API
- API

**Crawler**
- Browsertrix
- Brozzler
- Heritrix

**HTTP Server**
- ContaMeHistorias.pt

**Recorder**
- CompletePage
- CitationSaver
- SavePageNow
- AutoPatching

**Query Server**
- Nutchwax Search Engine
- SOLR Search Engine

**Document Server**
- Document Browser
- CDXJ Indexing
- Indexing
- Indexing

**GPU**
- Safe Image Classifier

**DevOps**
- Statistics
- Artifact Repository
- Configuration Manager
- CodeQA
- Routine Scheduler

**Third Party Management**
- Capacity
- Alarms (internal)

# 9 hardware profiles for servers to facilitate acquisitions

| Configuration Profile | CPU | RAM (GB) | Storage | Network | GPU | Description |
|---|---|---|---|---|---|---|
| Crawler | 20 | 256 | 50TB | PUBLIC | NO | Crawls information from the live-web in batch. |
| Broker | 40 | 256 | 40TB SSD | PUBLIC | NO | Hosts CDXJ indexes and routes queries to Page and Image indexes. |
| Recorder | 40 | 256 | 40TB | PRIVATE | NO | On-demand selective browser-based crawling. |
| Query Server | 20 | 512 | 50TB | PRIVATE | NO | Hosts Page and Image indexes to answer queries. |
| Document Server | 6-20 | 256 | 150TB | PRIVATE | NO | Hosts ARC and WARC files and indexing processes. |
| DevOps | 2 | 32 | 250GB (slave) / 8TB (master) | PRIVATE | NO | Continuous Integration tools: artifacts, builds, logs, and jenkins data. |
| GPU | 10 | 512 | 50TB | PRIVATE | YES | Not-Safe-For-Work image classifiers. |
| HTTP Server | 4-8 | 16 | 128GB | PUBLIC | NO | Routes HTTP requests and hosts informative website. |
| 3rd Party Management | N/A | N/A | N/A | N/A | N/A | Hired cloud services. |

| Demand: |
|---|
| HIGH |
| MEDIUM |
| LOW |

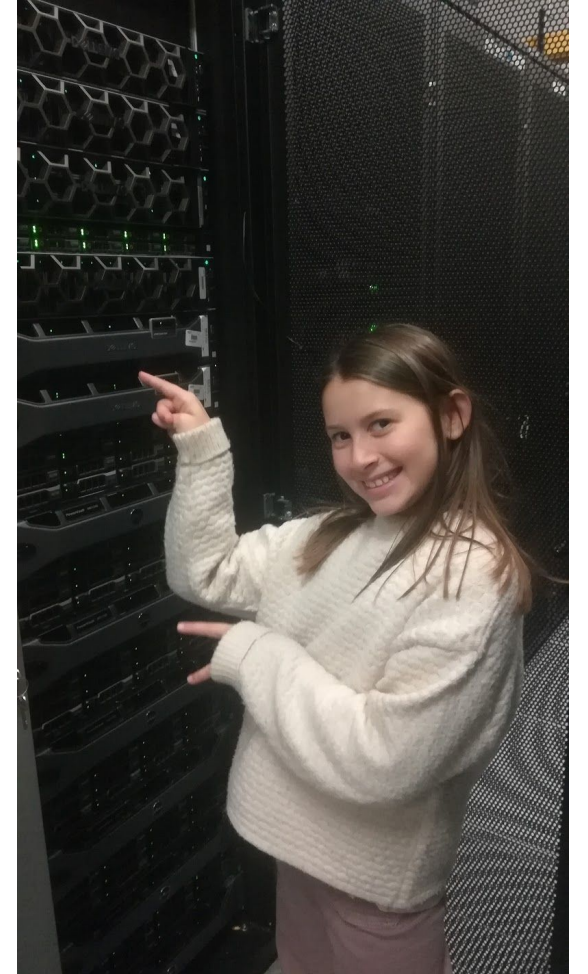# Own data-center shared with other public administration services

77 servers mounted on 167 Rack Units

- 7.4 meters in linear space

18 TB of RAM memory

2 180 vCPUs

1 234 hard drives (5.18 PB)

# The people behind the computers

# Suggested publications

- Web archives as research infrastructure for digital societies: the case study of Arquivo.pt (2022), Daniel Gomes
- The Anatomy of a Web Archive Image Search Engine (2022), André Mourão
- Information Search in Web Archives (2014), Miguel Costa
- Learning Temporal-Dependent Ranking Models (2014), Costa et al.
-  [arquivo.pt/publications](arquivo.pt/publications)

# Recommendations to manage a web archive

**No part-time web archivists**
- Web archiving is complex, requires training and full dedication
- Document thoroughly, prepare the leave when you hire someone

**Listen to your users**
- Don't guess what they need, hire User Experience researchers to find it

**Use existing tools and micro-services**
- Webrecorder.net, Archive-it
- Micro-services facilitate technology replacements

**Design-to-fail architecture**
- Service must resist to any component fault
- Test, test, test: maintaining is harder than developing