# Time to explore, time to learn from the archived Web

Arquivo.pt training initiative

ricardo.basilio@fccn.pt

# Agenda

1. Arquivo.pt

2. ROSSIO Infrastructure for Digital Humanities
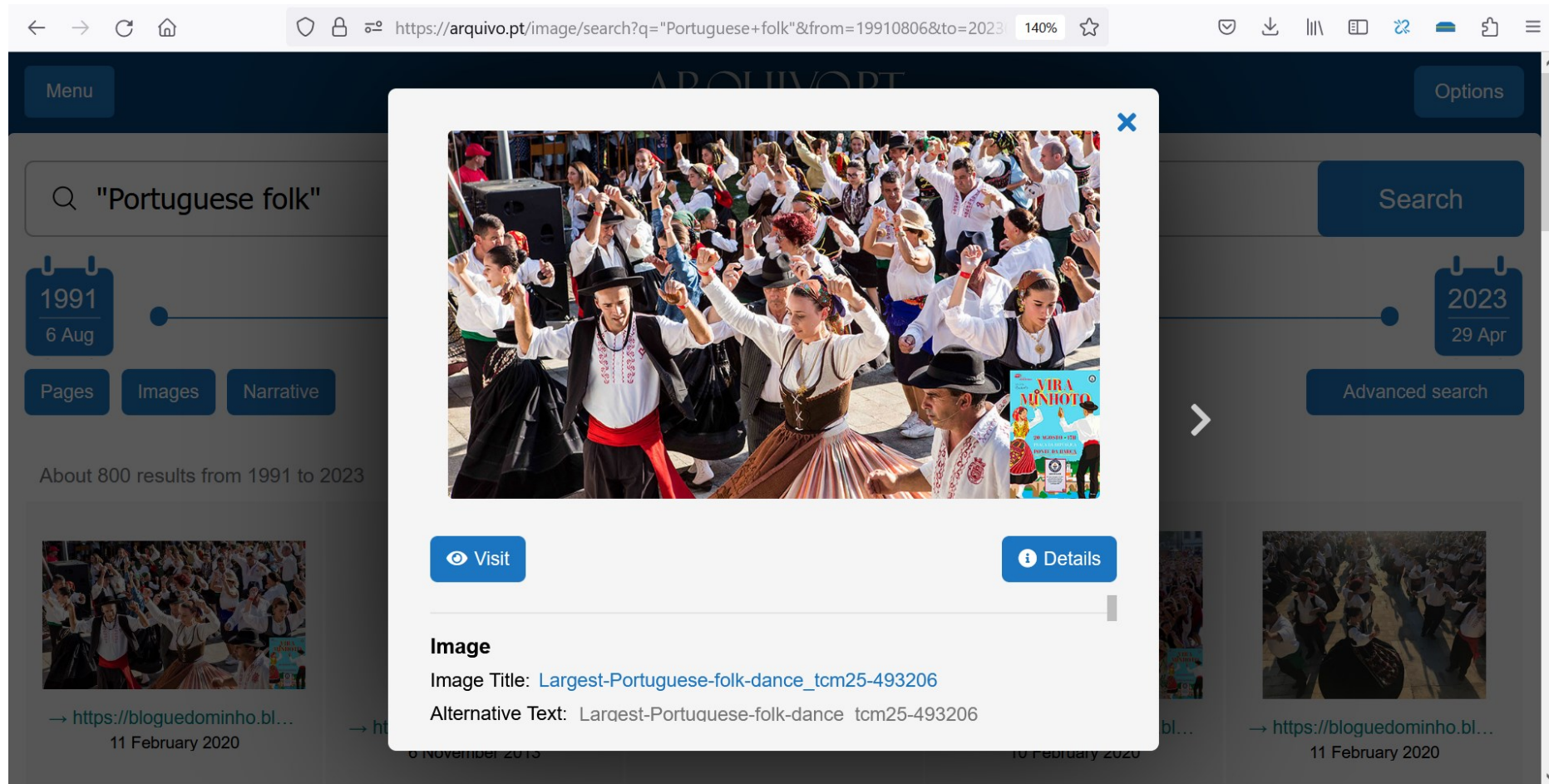
3. Arquivo.pt training initiative

# 1

# Arquivo.pt

# Arquivo.pt

- Portuguese public Web preservation service

- Provides a public search interface (by URL, text, image, and API)

- Aimed at academic research, but also at all citizens

# Arquivo.pt



[Example: a search for "Portuguese folk" on Arquivo.pt, the Portuguese Web Archive](#)

# Arquivo.pt



Arquivo.pt public interface

# Arquivo.pt



Arquivo.pt public interface and the use of the narrative button that uses the external service *Tell me stories/Conta-me histórias*
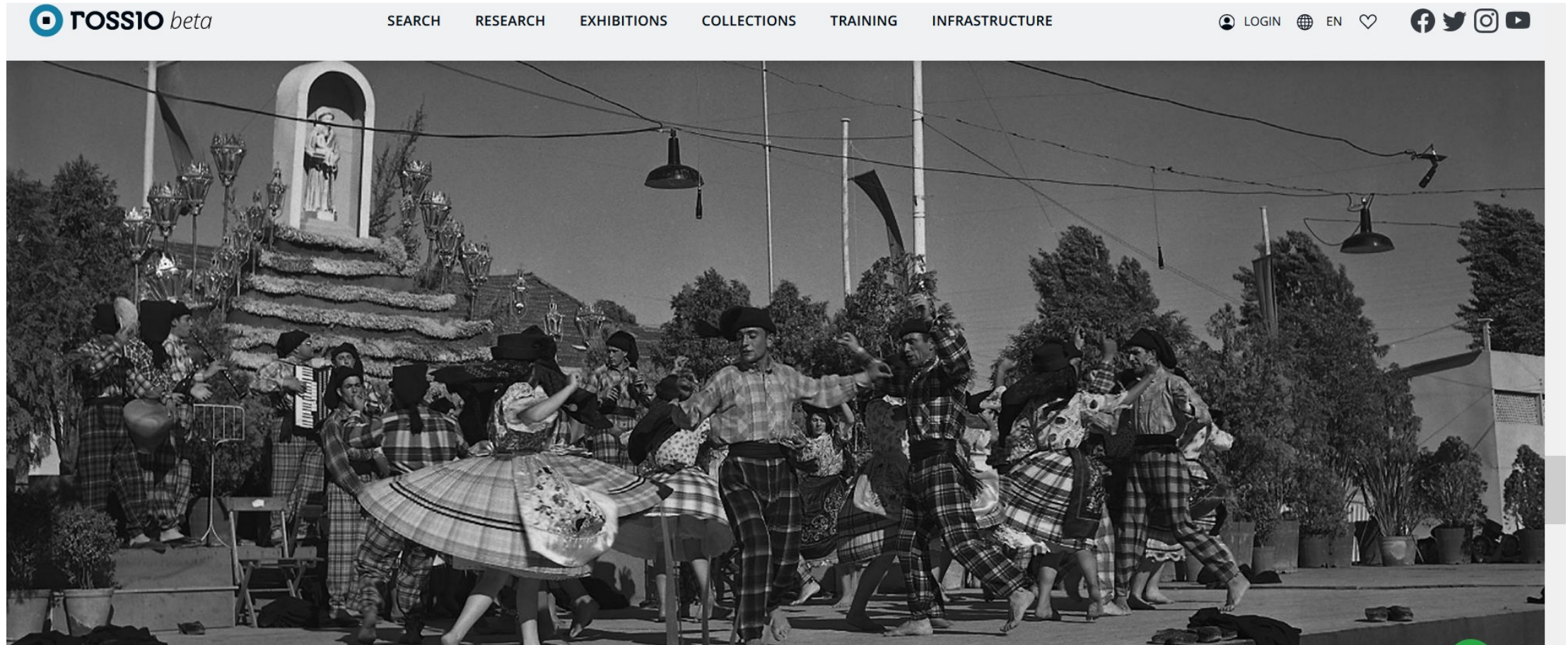
# 2

# ROSSIO

# ROSSIO

- Consortium of cultural and educational institutions

- Coordinated by the NOVA School of Social Sciences and Humanities of Lisbon

- Aggregates, contextualizes, enriches and disseminates digital content

Portuguese node of the DARIAH.EU

# ROSSIO



Web portal of the ROSSIO Infrastructure: https://rossio.pt

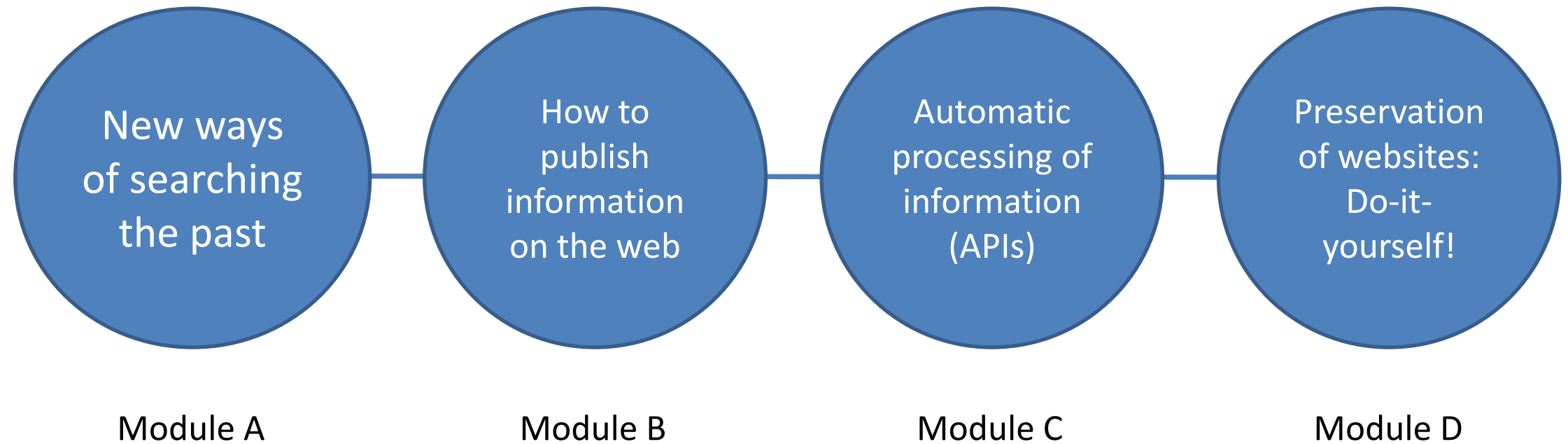# ROSSIO

# 3

# Arquivo.pt training initiative

Café with Arquivo.pt
Webinars Cycles

# Arquivo.pt training initiative

New ways of searching the past

How to publish information on the web

Automatic processing of information (APIs)

Preservation of websites: Do-it-yourself!

Module A    Module B    Module C    Module D

arquivo.pt/training

# Arquivo.pt training initiative



Module A          Module B          Module C          Module D

arquivo.pt/training

# Arquivo.pt training initiative

## New ways of searching the past

Module A

## Training contents

- The problem of 80% of Web content disappearing

- Search for historical contents in Arquivo.pt

- Everyday use cases

# Arquivo.pt training initiative

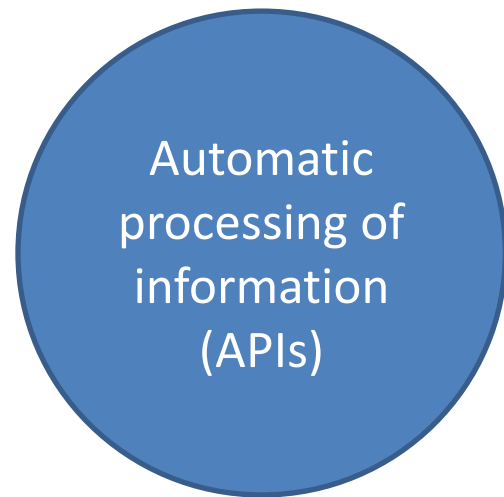How to publish information on the web

Module B

Training contents

- Recommendations for publication

- Robots.txt configuration

- Examples of poorly preservable Web publishing technologies

# Arquivo.pt training initiative

**Automatic processing of information (APIs)**

Module C

## Training contents

- Arquivo.pt APIs

- CDXJ Indexes (arquivo.pt/api)

- Use cases

# Arquivo.pt training initiative

**Preservation of websites: Do-it-yourself!**

Module D

## Training contents

- Record Web contents locally in a standard format

- ArchiveWeb.page tutorial

- Practical exercises (websites, social media like Facebook, Twitter)

# Arquivo.pt training initiative

# Café with Arquivo.pt
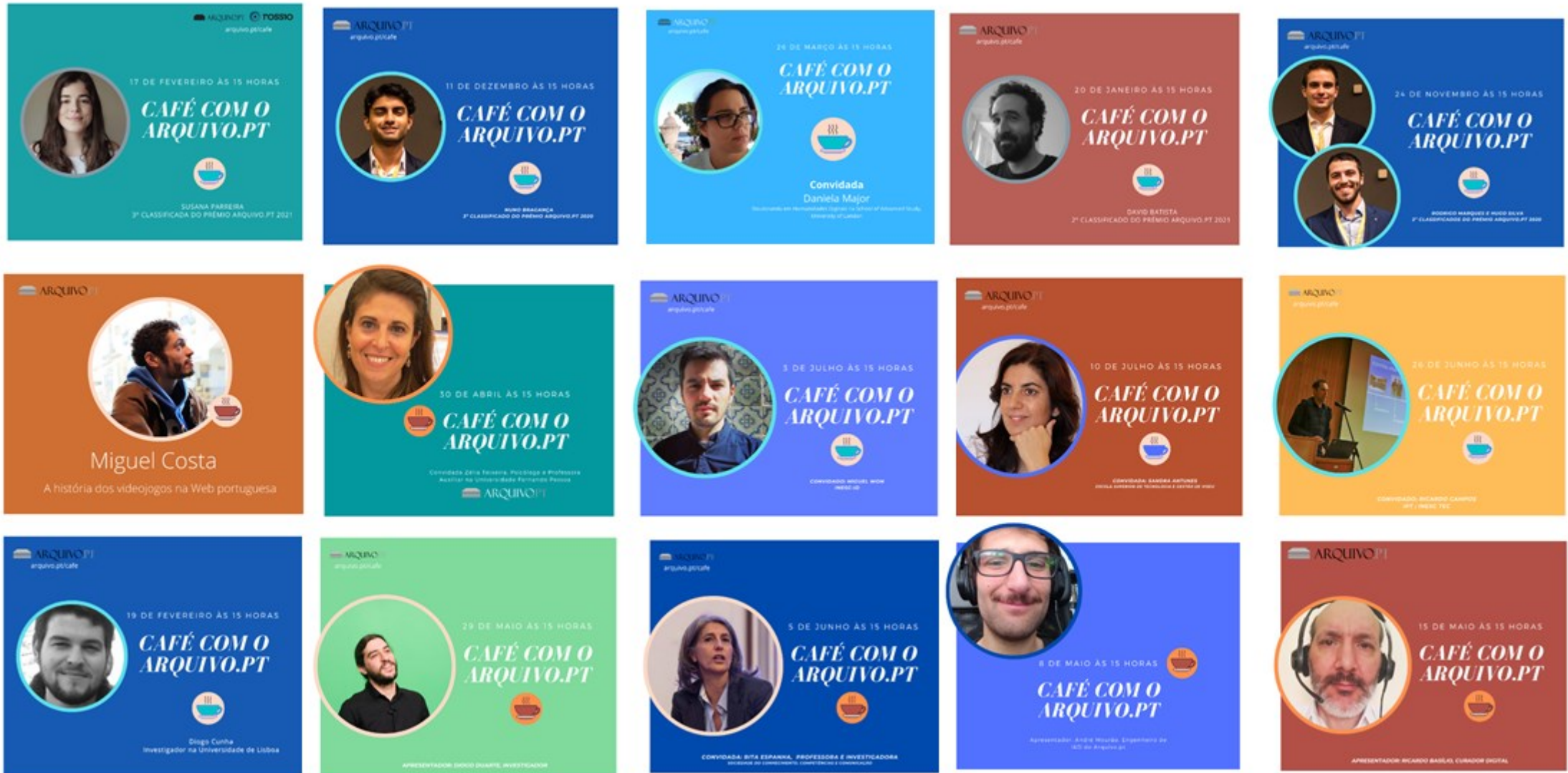
# Arquivo.pt training initiative



2019



2020

# Arquivo.pt training initiative - Café with Arquivo.pt

- <span style="color:red">680 participants</span>

- 23 sessions

- 30 participants/session
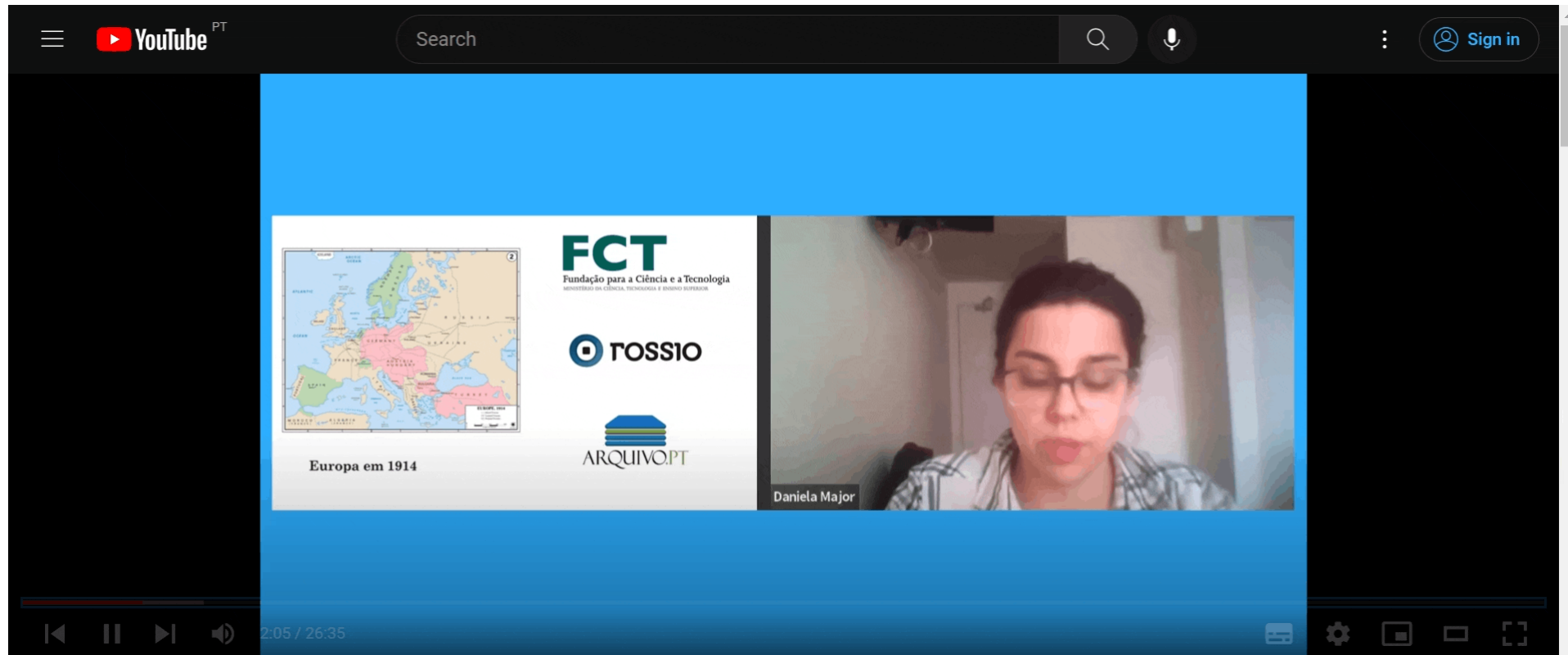
- 90% satisfaction

https://arquivo.pt/onlinecafe

# Arquivo.pt training initiative - Café with Arquivo.pt
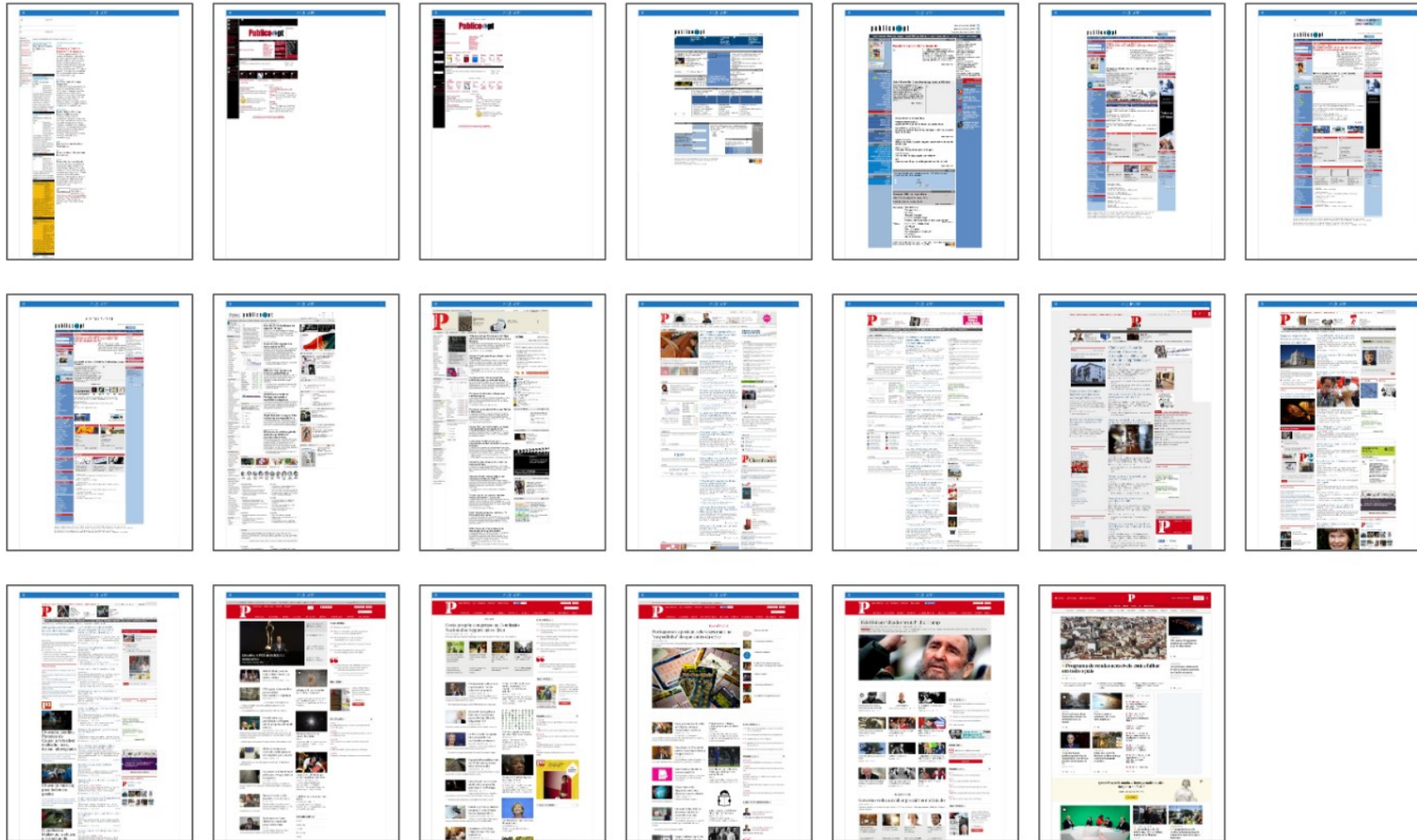
# Arquivo.pt training initiative - Café with Arquivo.pt



[The commemoration of the Centennial of World War I, by Daniel Major](#)

# Arquivo.pt training initiative - Café with Arquivo.pt

# Arquivo.pt training initiative - Café with Arquivo.pt



Front pages of Portuguese online newspapers, Arquivo.pt Award 2021 by Susana Parreira

# Arquivo.pt training initiative – Café with Arquivo.pt



Front pages of Portuguese online newspapers, Arquivo.pt Award 2021 by Susana Parreira

# Arquivo.pt training initiative - Café with Arquivo.pt



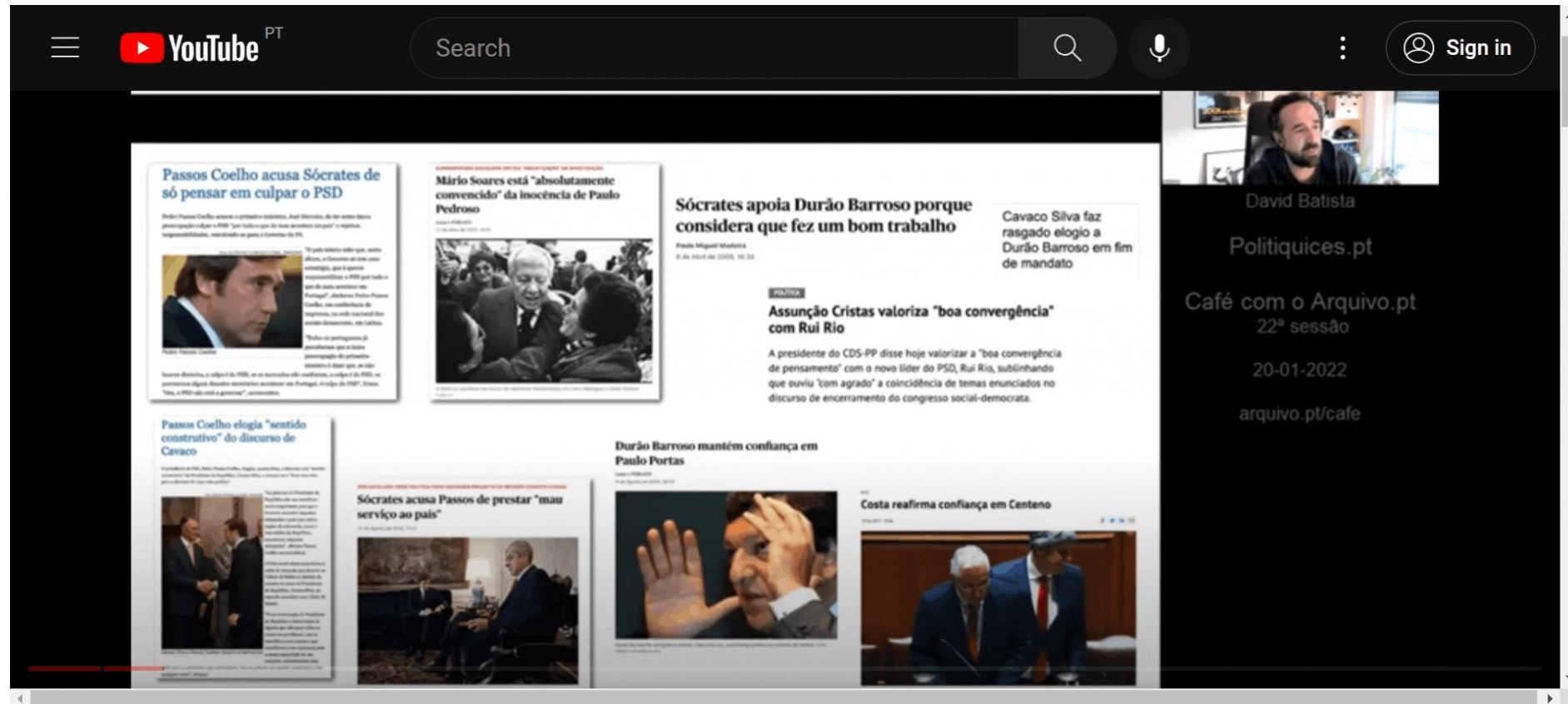[Major Minors, Arquivo.pt Award 2021 by Paulo Martins e Leandro Costa](#)

# Arquivo.pt training initiative – Café with Arquivo.pt



[Major Minors, Arquivo.pt Award 2021 by Paulo Martins e Leandro Costa](#)

# Arquivo.pt training initiative – Café with Arquivo.pt



[www.politiquices.pt](http://www.politiquices.pt), Arquivo.pt Award 2021, by David Batista

# Arquivo.pt training initiative – Café with Arquivo.pt



[www.politiquices.pt](http://www.politiquices.pt), Arquivo.pt Award 2021, by David Batista

# Arquivo.pt training initiative – Café with Arquivo.pt



Parlamientary archive (Portugal), Arquivo.pt Award 2022, by Tiago Santos - arquivo-parlamento.pt

# Arquivo.pt training initiative – Café with Arquivo.pt



Parlamientary archive (Portugal), Arquivo.pt Award 2022, by Tiago Santos - arquivo-parlamento.pt

# What a web archive is for?



Areas of study of the 160 works competing for the Arquivo.pt Award, between 2018 and 2023

# What a web archive is for?



Studies versus IT applications. Of the 160 works competing for the Arquivo.pt Awards, between 2018 and 2023

# Arquivo.pt training initiative

# Webinars cycle

# Arquivo.pt training initiative – Webinars cycle

Webinars cycle promoted by the Lisbon
City Council

Target: Common citizens



Page view through GoogleTranslate

facebook.com/competenciasdigitaiscmlx

# Arquivo.pt training initiative – Webinars cycle

Webinars cycle promoted by the
Fundação Calouste Gulbenkian
(Art Library and Archive)

Target: artists, gallery owners  and

reseachers

Page view through GoogleTranslate

gulbenkian.pt/biblioteca-arte/noticias/arte-para-sempre-na-web

# Arquivo.pt training initiative – return to face-to-face presentations

Digital Humanities,  a priority

Target: Professors, students and researchers



Workshop promoted by the Centre for Comparative Studies (CEComp) at the School of Arts and Humanities, University of Lisbon.

# Arquivo.pt training initiative - Methodological challenges

«Dealing with a greater volume of data will be a common occurrence that historians will have to face from now on. This opens up **a world of possibility for historians**, but it also forces us to ask difficult questions».

Daniela MAJOR, Early Stage Researcher in Digital Humanities at the School of Advanced Study in London. It Takes a Village to Raise an Archive. How the Use of Web Sources Fosters Collaboration, Published at Web Archiving Section of the Society of American Archivists on January 18, 2023.

# Arquivo.pt training initiative

# Conclusions

# Conclusions

- 2018 - Four modules to offer to the community

- 2020 – Café and Webinars to receive contributes from the community

- 2022 – Two-way knowledge sharing (Users and Arquivo.pt team)

There are methodological challenges for Digital Humanities using

Web archives and these can be overcome with close collaboration.

# Thank you!

contacto@arquivo.pt

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

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

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