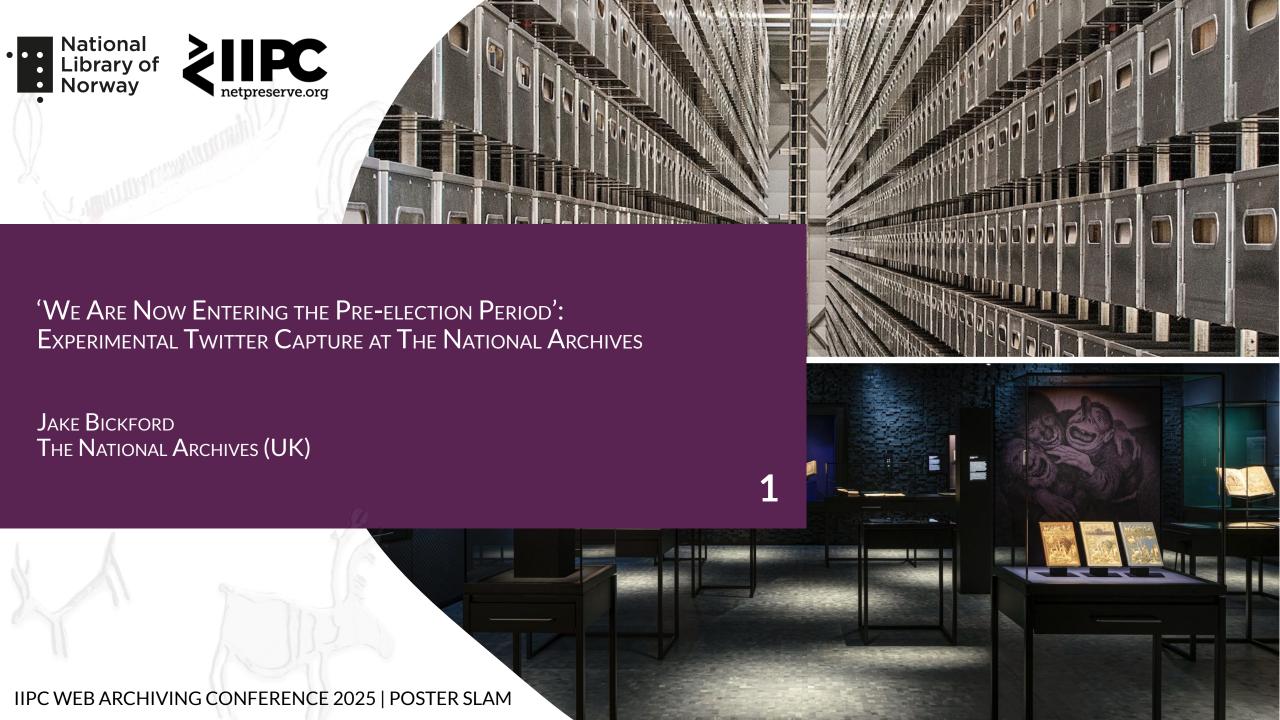
INSTRUCTIONS:

Find your title slide below. Use the blank slide following your title slide for your 1-minute lightning talk.

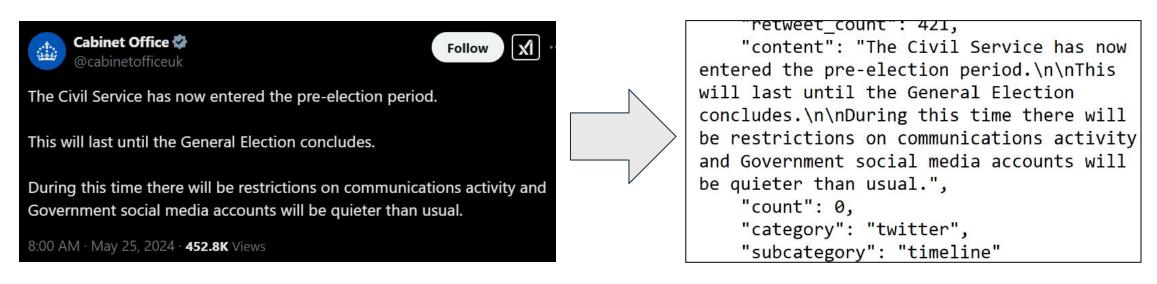
- Do not include any audio, video, or clickable links in your slide.
- If you wish to use a different template for your slide, click "LAYOUT" in the toolbar to see your options.
- You may add other images, text boxes, or tables/diagrams if needed, but take care not to overload your slide.

If you don't need a slide for your talk, we will use the slide with the title of your poster. Please remove your "blank" slide if you choose this option.

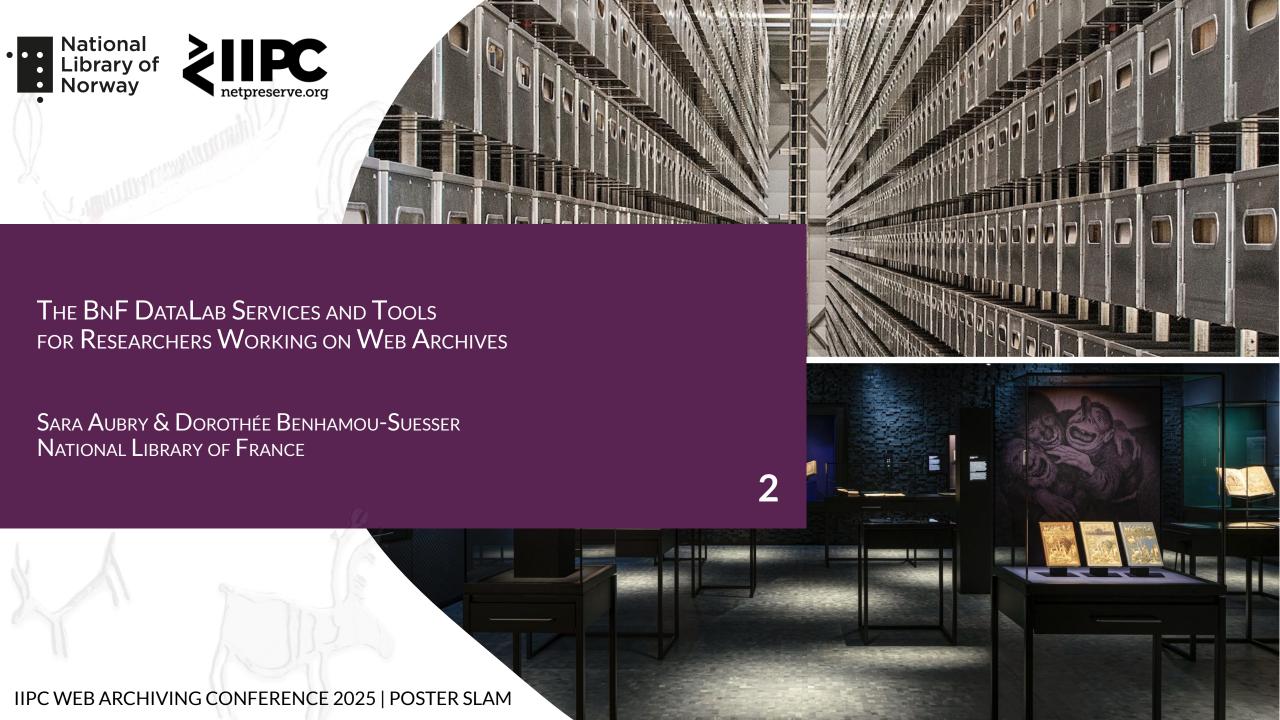




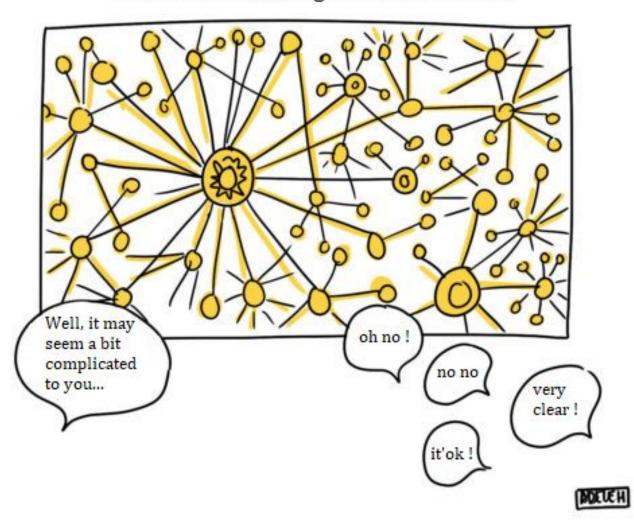
'We Are Now Entering the Pre-election Period': Experimental Twitter Capture at The National Archives (UK)



FOLLOWING THE UNEXPECTED ANNOUNCEMENT OF THE 2024 UK GENERAL ELECTION, WE USED CLIENT-SIDE CAPTURE WITH GALLERY-DL TO RAPIDLY FILL GAPS IN OUR SOCIAL MEDIA COLLECTION THAT HAD BEEN CAUSED BY CHANGES TO THE TWITTER API



Librarians working with researchers



Source: https://ffl.hypotheses.org

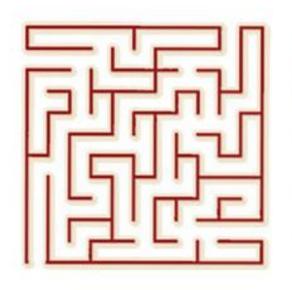






Designing Art Student Web Archives

RESULTS FROM A 2 YEAR STUDY OF WORKFLOWS



Staff support

Project Tracking Database

Archive-It organization

Finding Aid description

Accession documentation (accruals?!)







WY WY OUT OUT TOTAL TOTA

Next Steps Towards a Formal Registry of Web Archives for Persistent and Sustainable Identification

Why do we need a registry?

A registry can bridge between web archives & current/future researchers



Web Archives Registry

Merge

Minimum information!

Example

ISNI

Web Archives	(cgistry	
Formal Id A-Z	Name A-Z	Domain A-Z Valid to A-Z
1000 0000 0168 4170	sh Web Archive	nli.ie 01-01-2011 🗸
1000 0000 0168 4177	NLA Australia	pandora.nla.gov.au 01-01-2004 🖳 01-01-2012
1000 0000 0168 4177	Australian Web Archive	webarchive.nla.gov.au 01-01-2009
1000 0000 0168 4182	Netarkivet	netarkivet.dk 01-07-2005 01-07-2035
1000 0000 0168 4182	Danish Web Archive	kb.dk 01-07-2035 🗸



What is the challenges?

- Web archives will merge
 National Library of Australia
- Web archives will change URL domain (National Library of Ireland 2019)
- Web archives will change access URLs (Royal Danish Library 2030)

Where to use it?

- when a reference points to an archive URL that does not work
- to update resolver to PWIDs

Other information elsewhere?

Alternative domains, Service path, e.g. for Archive.Today https://arquiv

https://arquivo.pt/wayback Service type, e.g. Access

- archive.ph

archive.md Trailing characters to date,

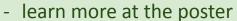
e.g. im_ for image

Case sensitive lookup

https://archive-it.org/home/nli/{4-14-digit YYYYMMDDhhmmss}/{url}







- give feedback

Next steps

- come to IIPC webinar this year







New

domain





about the directory

▼ Search the ELO Direct

Search

contact administrator

create/edit entries **▼** Browse the Directory

Genre/Length

Home About Log in My account ELO home

ELECTRONIC LI

The Electronic Literature Directory is a resource for readers and writers of born-digital literature. Created by the Electronic Literature Organization, it provides an extensive database listing electronic works, their authors, and their publishers. The descriptive entries are drafted by a community of e-lit authors

These changes in approach reflect shifting attitudes seen within scholarship.

This Directory is a proje supported by tax-deduc server today.

Individual works

Browse all works of electronic literature.

Antecedents

Browse all works marked as e-lit antecedents.

E-lit resources

Browse all works of criticism, collections, and resources.

Featured articles

Browse all articles commissioned for the Directory.

OUR PROCESS

ABOUT
MY ACCOUNT
ELO HOME
LOG IN

RECENTLY

With love, from a







Arquivo.pt Award 2025 (annual): open applications

10 000€ FOR THE WINNER FOT

Works that make use of Arquivo.pt

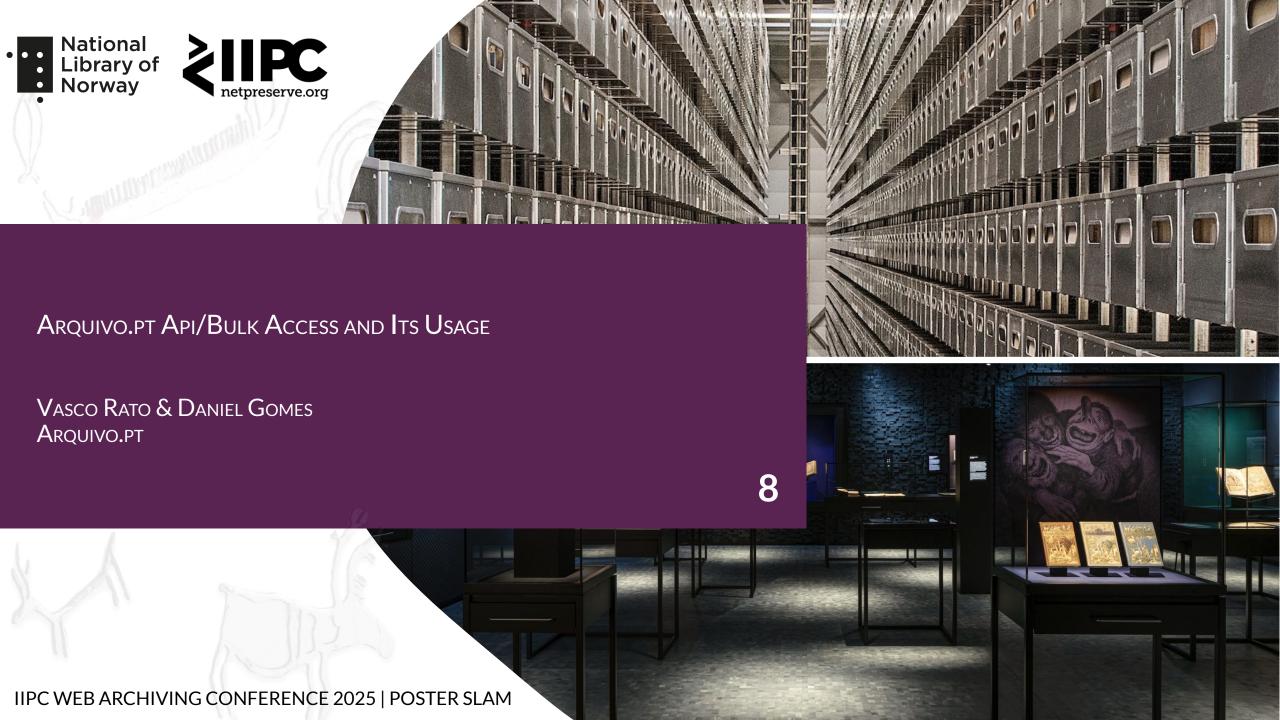
Applications until 6 May 2025

HIGH PATRONAGE OF PRESIDENT OF PORTUGAL

SHARE!

KNOW MORE: ARQUIVO.PT/AWARD

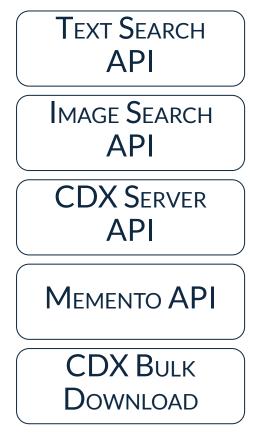




Arquivo.pt Api/Bulk Access and Its Usage

Data is only useful if it's accessible and searchable!











Glória LLM



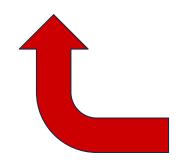


Breaking the cycle

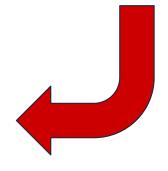
Experiment with a new capture technology



Export WARC and upload it to our Archive-It



Capture does not appear to work when viewed















Community of Practice Social Media Archiving

Praktijkwijzer Sociale Media Archiveren

Toolbox

https://kennisbank.meemoo.be/toolbox/praktijkwijzer-sociale-media-archiveren

Deze praktijkwijzer leidt naar praktische handleidingen en workflows voor socialemedia-archivering die nog relevant zijn en zo goed als mogelijk up-to-date worden gehouden door de archiverende community. We moedigen aan om eigen handleidingen of relevante links toe te voegen aan...

Sociale media archiveren met Browsertrix

Toolbox

https://kennisbank.meemoo.be/toolbox/sociale-media-archiveren-met-browsertri

In het kader van het project Best practices voor de archivering van sociale media in Vlaanderen en Brussel werden verschillende tools getest om socialemediaplatformen te archiveren. Deze handleiding beschrijft de tool Browsertrix die gebruikt kan worden voor Facebook en...

Sociale media archiveren met Webrecorder

Toolbo

https://kennisbank.meemoo.be/toolbox/sociale-media-archiveren-met-webrecorder

In het kader van het project Best practices voor de archivering van sociale media in Vlaanderen en Brussel werden verschillende tools getest om socialemediaplatformen te archiveren. Deze handleiding beschrijft de tool Webrecorder voor het archiveren van sociale media.

tutorials on tools and workflows https://meemoo.be/en/knowledge-base

Package level

TABLE OF CONTENTS

- 1 METS.xml (file)
 - a Elements and internal references
 - b <mets> section
 - c <metsHdr> section
 - d <dmdSec> section
 - e <amdSec> section
 - f <fileSec> section
 - g <structMap> section
- 2 /metadata (directory)
 - a /descriptive (directory)
 - b /preservation (directory)
 - a Describing Intellectual Entities
 - b Adding provenance of representations
- 3 /representations (directory)

Representation level

TABLE OF CONTENTS

- 1 /representation_1 (directory)
- 2 METS.xml (file)
 - a Elements and internal references
 - b <mets> section
 - c <metsHdr> section
 - d <structMap> section
- 3 /data (directory)
- 4 /metadata (directory)
 - a /descriptive (directory)
 - b /preservation (directory)

e-ARK based SIP specification https://developer.meemoo.be

Supporting best practices for archiving social media by heritage institutions in Flanders (and beyond)



Ellen Van Keer <u>ellen.vankeer@meemoo.be</u> Katrien Weyns <u>katrien.weyns@kuleuven.be</u>



Planning Web Archiving Within a Four-Year Scope

We have made a new collection plan for years 2025-2028

COVERS BOTH WEB ARCHIVING AND E-LEGAL DEPOSIT

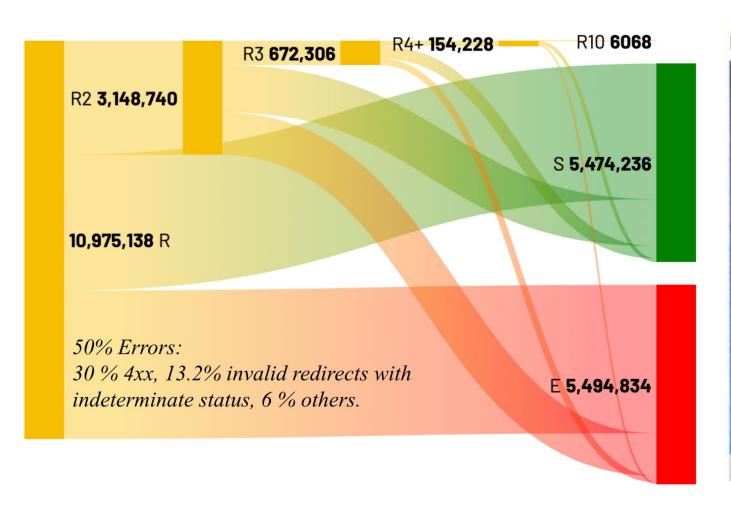
Also, a new policy document (what we should collect) is in the works ***

INTERESTED ABOUT THE COMPARISON RESULTS OR THE ACTUAL PLAN? COME CHAT WITH ME!





Redirects Unraveled: From Lost Links to Rickrolls



Rick Astley-Never Gonna Give You Up



Step 1

Capture

- Screenshots are taken from dynamic or ephemeral content sources
- Ensuring contextual accuracy and relevance

Step 2

Data Extraction with AI

- Image Recognition: Identifying key visual elements
- Natural Language Processing (NLP):
 Extracting and interpreting text
- Optical Character Recognition (OCR):
 Converting embedded text into searchable data

Step 3

Data Utilization

Extracted data can be categorized and stored for:

- Metadata analysis
- User interaction studies
- Archival documentation





A Prototype Pipeline for WARC Annotation

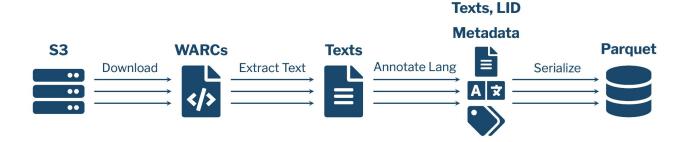
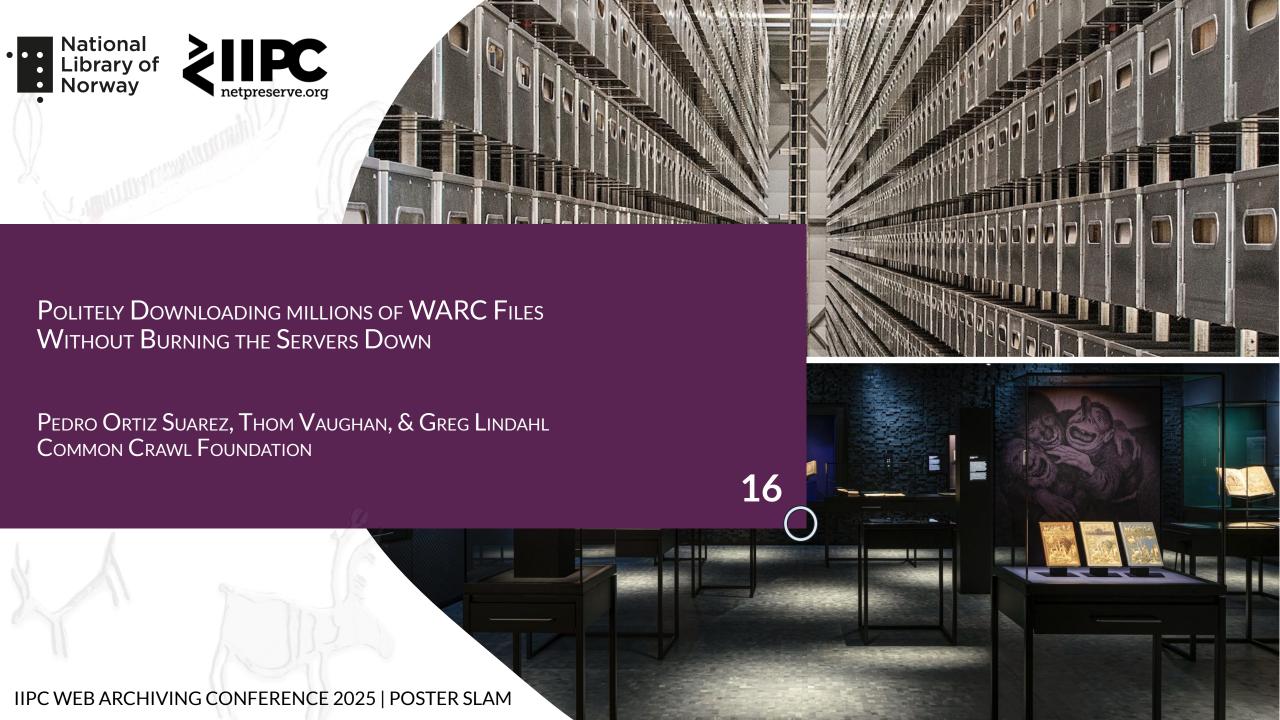


Figure 1. A schema of the first prototype for our pipeline. represents the AWS S3 bucket where the Common Crawl Data is stored, represents WARC files, represents text extracts of WARCs, represents language tags, represents metadata coming from WARC headers and represents Parquet files. Finally the arrows represent parallel processing.

- WE INSPIRE OURSELVES ON EXISTING PIPELINES FROM THE NLP AND THE WEB ARCHIVING COMMUNITIES, AND DEVELOP OUR OWN PROTOTYPE OF A PIPELINE FOR WARC ANNOTATION
- The first experimental prototype aims to be efficient, modular, open-source and user-friendly so that little knowledge of large-scale data processing is needed
- This first prototype focuses on organizing and compiling some of the hectic developments of data pipelines in the NLP/ML space and bringing them to the web archiving community
- Our prototype focuses on textual data, as Common Crawl is a text-only archive
- We note that the architecture of our prototype is extensible and borrows ideas from NLP pipelines that have already been extended to non-textual archives
- FOR THIS FIRST PROTOTYPE WE ONLY IMPLEMENT TEXT EXTRACTION AND LANGUAGE ANNOTATIONS.







cc-downloader



- We use **Exponential Backoff** and **Jitter** to develop our own download client for **Common Crawl**, **cc-downloader**.
- We develop it in the **Rust Programming Language** and release pre-compiled binaries for Linux, Mac and Windows (x86-64).
- We release cc- downloader under Apache 2.0 and MIT licenses
- THE DEFAULT CONFIGURATION CAN BE TWEAKED BY THE USER
- CC-DOWNLOADER IS DISTRIBUTED AS A CLI TOOL, BUT WE ARE PLANNING ON RELEASING IT AS A LIBRARY AND ADDING PYTHON BINDINGS



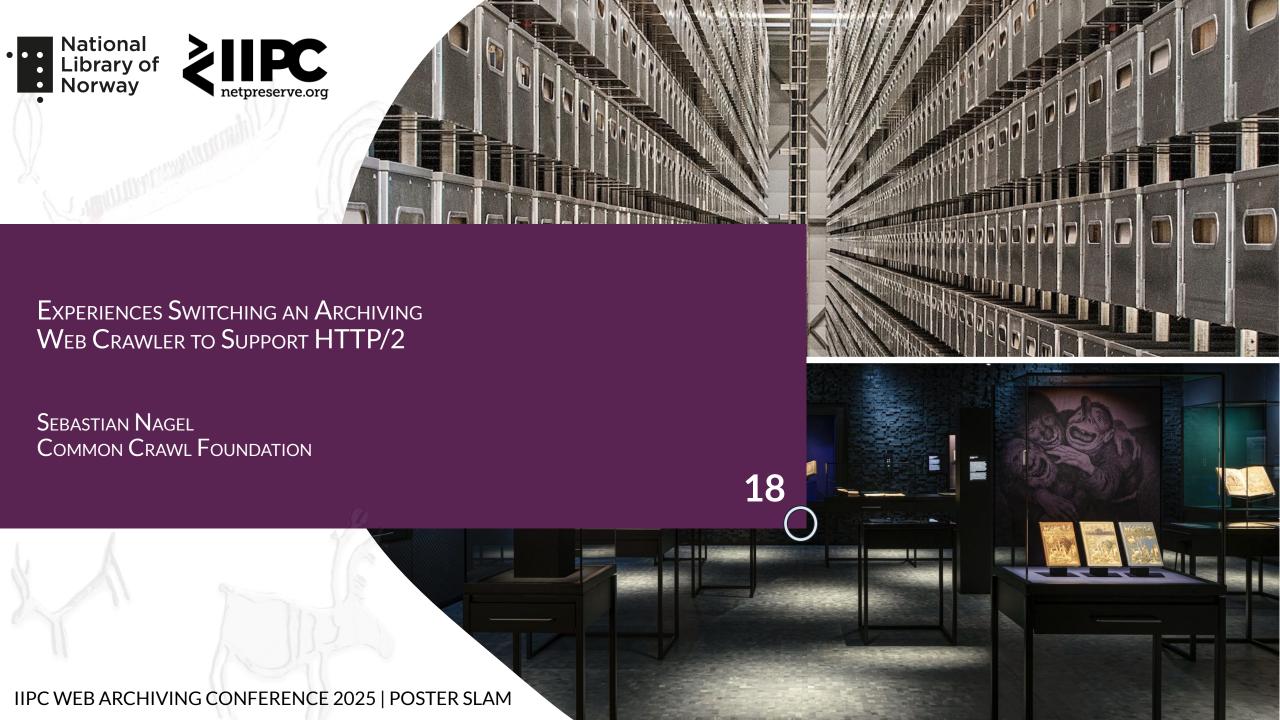




The Robots Exclusion Protocol

- 1994 ROBOTS.TXT DISCUSSED ON MAILING LIST 1996 INOFFICIAL RFC PROPOSAL
- 2022 RFC 9309 STANDARDIZES THE REP
- 2025 ROBOTS.TXT IS THE MOST WIDELY USED PROTOCOL TO BLOCK CRAWLING BY GENERATIVE AI CRAWLERS AND AGENTS
 - WHICH CRAWLERS ("USER-AGENTS")?
 - Does this impact web archiving?





HTTP/2 Support in a Web Archive Crawler

- HTTP/2 INTRODUCED IN 2015 AS AN IMPROVEMENT OVER HTTP/1.1 (1997)

- 70% of the traffic in the web are already on HTTP/2 HIGH TIME TO ADD SUPPORT FOR IT TO OUR WEB CRAWLER THE NEXT VERSION - HTTP/3 - IS OUT SINCE 2022

IMPLEMENTATION AND PERFORMANCE OF THE UPGRADED CRAWLER
 WARC CAPTURES OF HTTP/2 REQUESTS AND RESPONSES

