



Open WebSearch

European Web Crawls and LLMs: the OpenWebSearch.eu Project

Prof. Dr. Michael Granitzer
Chair of Data Science University of Passau and
Coordinator OpenWebSearch.eu

Open Web Search 



Funded by
the European Union

SUPPORTED
BY

NGI

OpenWebSearch.eu will create an open European infrastructure for internet search, based on European values and jurisdiction



What?

Restore an open search ecosystem / market as a basis for a new Internet Search

- lay a foundation for a new Internet search
- contribute to Europe's digital sovereignty
- empower Europe's researchers, innovators and businesses to systematically tap into the Web as business and innovation resource



Why?

1. Web search is dominated and limited by a few gatekeepers like Google, Microsoft, Baidu, Yandex.

Resulting situation:

- unilateral, biased, opaque access to information
- locked-in effects

2. Tapping the Web as resource is challenging for innovators and researchers



Who?

14 renowned European universities + institutions will pool their expertise and resources.

- including some of the largest research and computing centres in Europe
- e.g. IT4Innovations, Leibniz Supercomputing Centre, CSC, European Organisation for Nuclear Research CERN



How?

Develop the core of a European Open Web Index

Four Objectives

1. Open Technology Stack
2. Resource provision by a network of infrastructure providers
3. Added value services
4. Bootstrapping the ecosystem

14 Partners plus Third Party Calls



Webis.de



Research



Infrastructure



NGOs



Businesses

First Goal

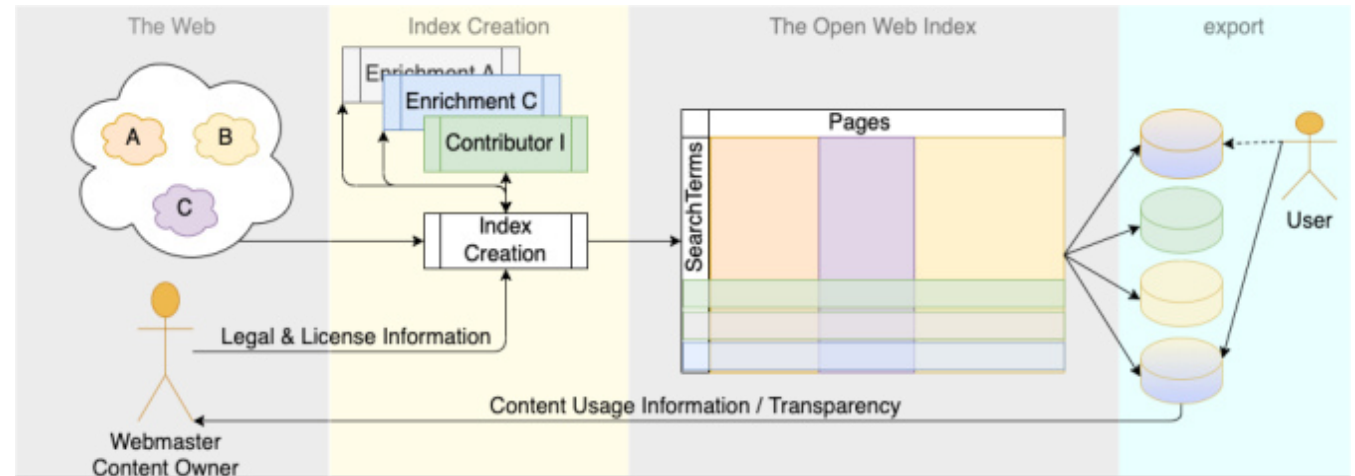
Build an Open Web Index Collaboratively



- Build an Open Web Index, i.e. a data structure for searching the Web
- Empower users, researchers & innovators to build on top of the Index

Principles for an Open Web Index

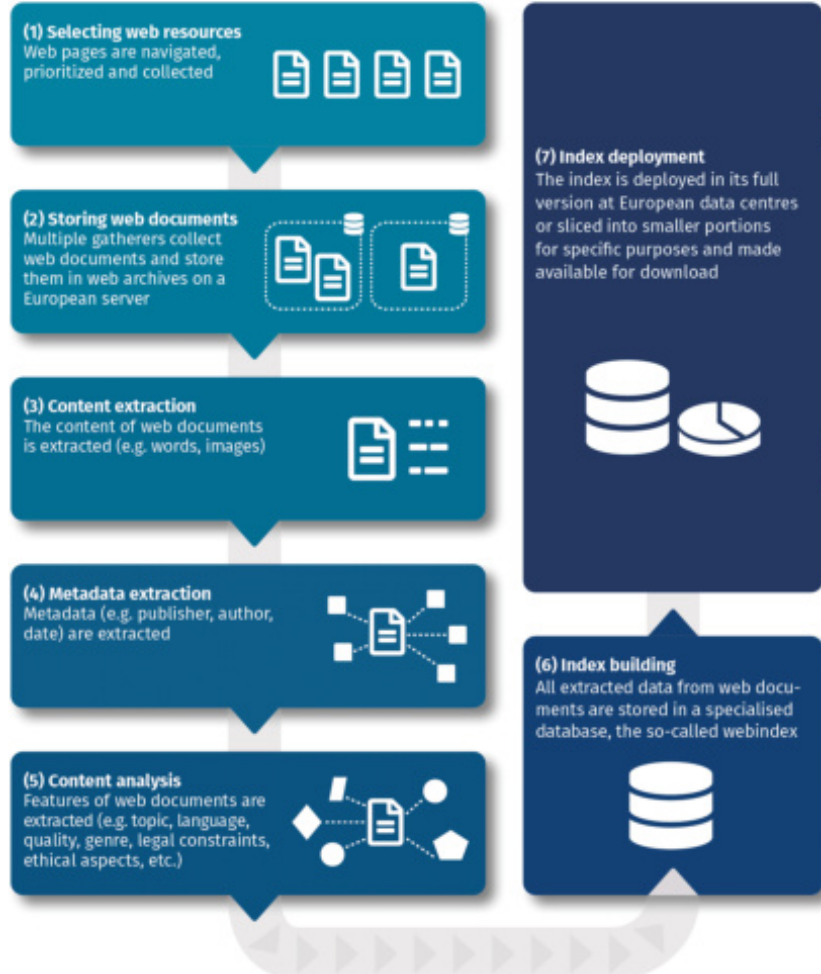
- Open Data for maximum re-use
- Open Source / Open Configuration for transparency
- Open Resources for using your own compute
- Open to contributions from third parties - empower innovators
- Collaborative management of a Web Index - connect existing infrastructure organisations
- Content control - respect legal, societal and ethical frameworks



Second Goal: Creating an Open Infrastructure

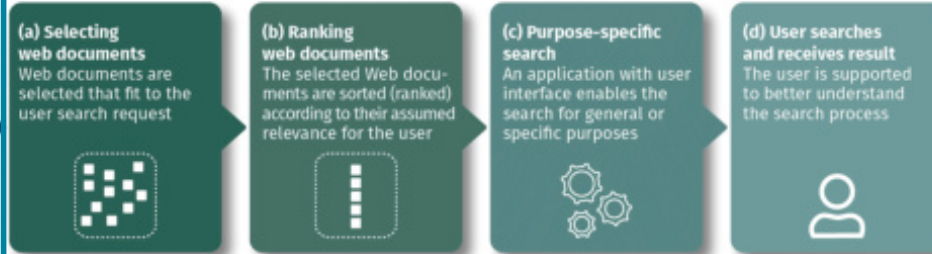
Index Generation

Web resources are selected and retrieved, their content and metadata are analysed, and all data stored in the index database.



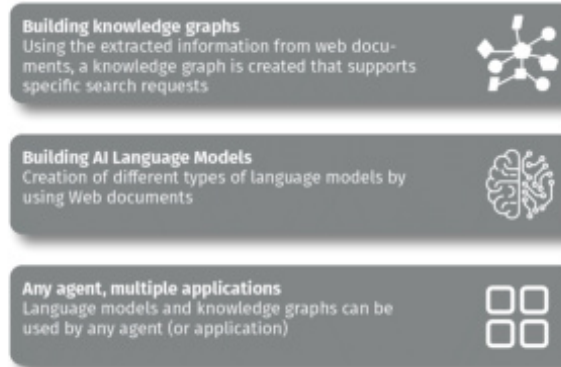
Search Applications

A user search request will be answered by a search application that makes use of the open web index.



Data Products

Knowledge representation models will be created using the open web index, in order to be used by any agent and for many applications



LUMI@CSC



KAROLINA@IT4I



Third Goal: Grow an Open Web Search Ecosystem



Index Generation

Web resources are selected and retrieved, their content and metadata are analysed, and all data stored in the index database.

(1) Selecting web resources
Web pages are navigated, prioritized and collected



Data Contributions

(e.g. provision of crawls, content push)

Content Curation

(e.g. science, education, languages)

Technology Contributions

(e.g. enrichment,)

Standards and ELSA Clearance

(e.g. license, metadata formats)

(5) Content analysis
Features of web documents are extracted (e.g. topic, language, quality, genre, legal constraints, ethical aspects, etc.)



(7) Index deployment

The index is deployed in its full version at European data centres or sliced into smaller portions for specific purposes and made available to users



Search Applications

A user search request will be answered by a search application that makes use of the open web index.

(a) Selecting web documents
Web documents are selected according to user search criteria

(b) Ranking web documents
The selected web documents are ranked according to relevance for the user

(c) Purpose-specific search
An application with user-specific requirements enables the user to perform general or specific searches

(d) User searches and receives result
The user is supported to better understand the search process

New Search Paradigms

(e.g. Argument search, conversational search)

Vertical Search Engines

(e.g. Open Science / Mobile Location Search)

Language Models

(e.g. language specific, search specific)

Benchmarking

(e.g. search engines, language models)

Web Analytics

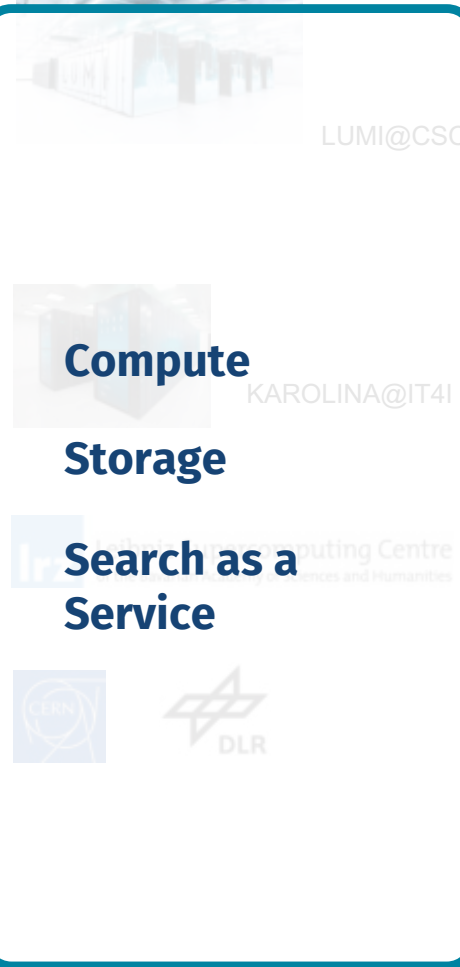
(e.g. Content distribution, social media)

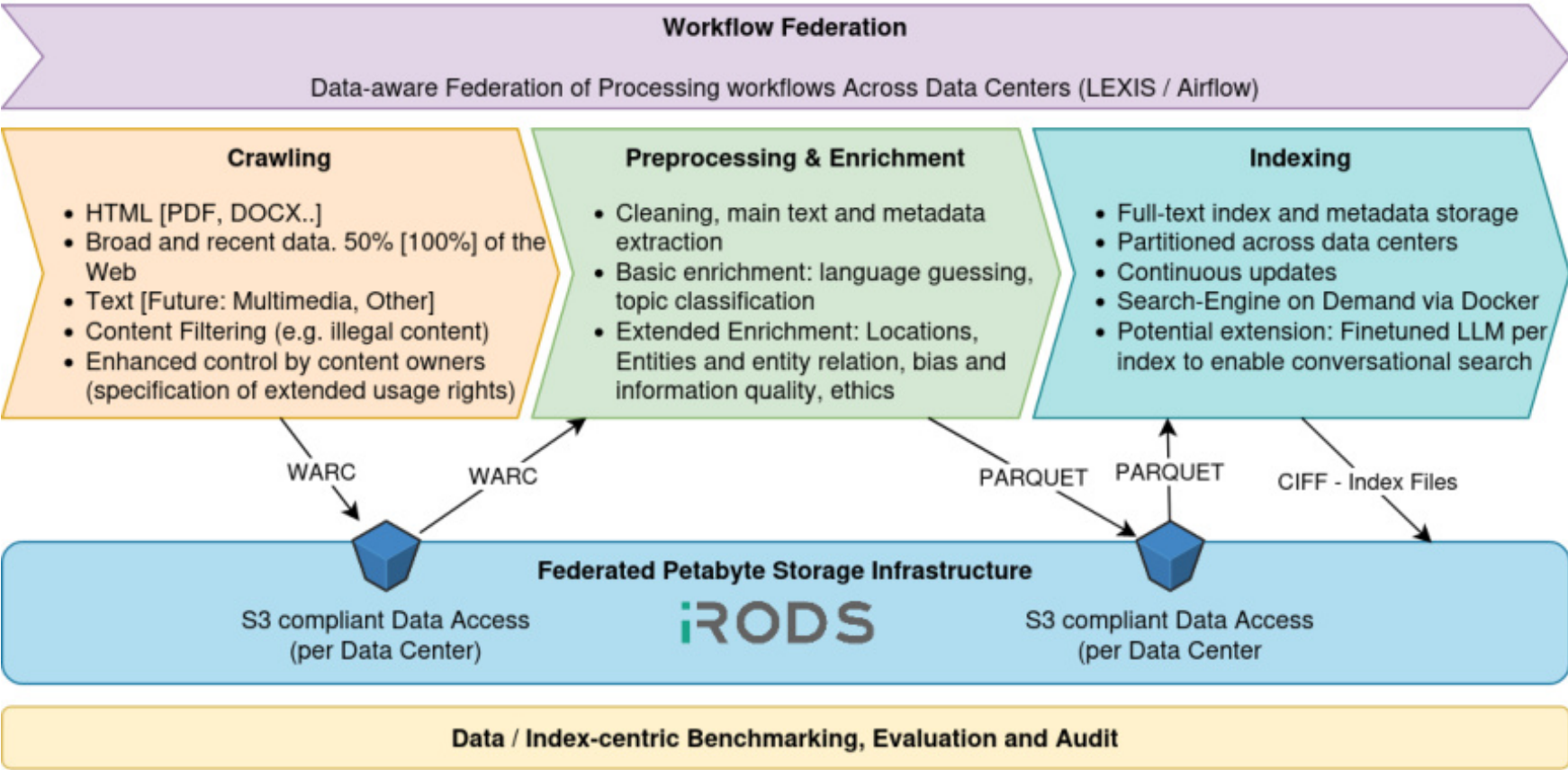
Data Products

Knowledge graphs can be created using the open web index in order to be used by any agent and for many applications

Building AI Language Models

Large language models can be built by training on the open web index to support specific search requests





Local / Vertical/ Federated Search
1. Re-use parts of the index
2. Domain-specific processing
3. Merge with local data
4. Local LLMs / Knowl. Graphs
5. Continuous Update

.....

Web-scale Analytics

Web 3.0 search

VR Web

Training and Finetuning of LLMs
1. Filter relevant web data
2. Finetune / train in same infrastructure
3. Or re-use in local training
4. Merge with up-to-date, cleaned, legally compliant web-data
5. Index plus LLMs for conversational search

Data Preparation

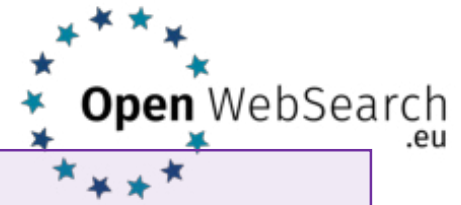
Storage

Application



More technical details on processes and infrastructure

The OWler – an Open Cooperative Web Crawler



OWler is an incremental and focused web crawler that extends the StormCrawler. It collects web pages that satisfy some specific criteria, e.g., URLs that belong to a given domain or that contain a user-specified topic, and documents its fetch activities in the WARC file format.

Code: <https://openwebsearcheu.pages.it4i.eu/wp1/owseu-crawler/owlers/>

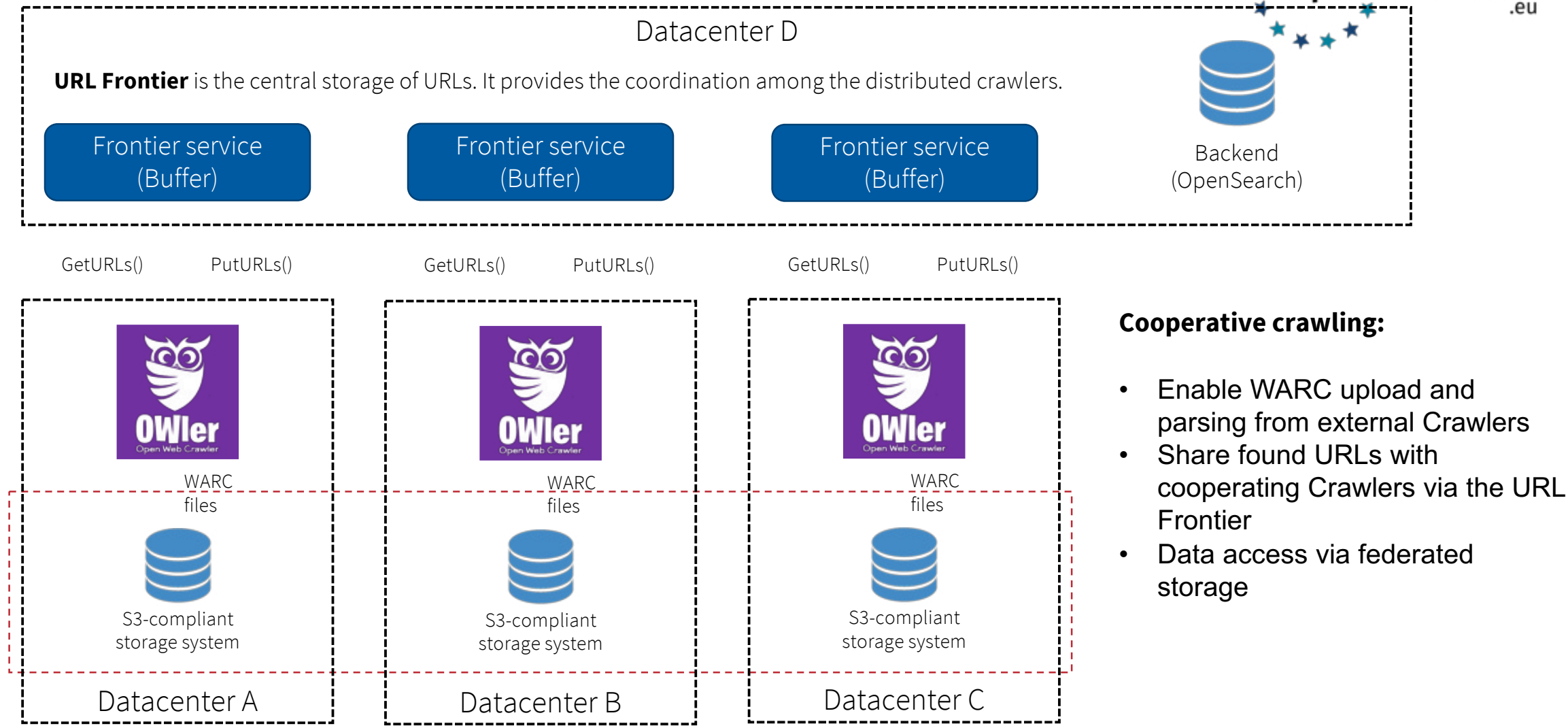


StormCrawler is a popular and mature open source web crawler. It is written in Java and is both lightweight and scalable, thanks to the distribution layer based on Apache Storm.



Apache Storm is a distributed system for processing streams of data. The work is delegated to different types of components that are each responsible for a simple specific processing task.

Cooperative Crawling



- Pipeline implementation and deployment until September 23
- Continuous Crawling from September onwards
- **Results**
 - Partitioned Index for download (Top Level Domain / topic based organisation)
 - Additional: Metadata, cleaned main text
 - Resiliparse as Preprocessing and Cleaning Library
- **Future Challenges:** legal compliance in crawling (e.g. IPR, GDPR), data quality, integration of LLMs, support training for LLMs, scaling-up, legal and ethical clearance

Advantages of an Open Web Search Infrastructure



- A strategic contribution to **digital sovereignty**
- Restoring an open, **human-centric search engine market** in Europe - for diverse and unbiased access to knowledge and information.
- **Lowering entry barriers for tapping the Web as a resource** at scale: for researchers, innovators, and businesses
- Providing **web services for other digital infrastructures** and data spaces, e.g. the European Open Science Cloud, GAIA-X, EDICs/MCPs, Copernicus, and many more.
- **Increasing transparency** of web content usage
- **Increase of control** for content owners and users
- Enable Europe to develop its own **large-scale language models and generative AI**



Conclusion



- Web and Web Search critical for Europe's digital sovereignty
- Reducing entry barriers for researchers, innovators and business to increase competitiveness with Big Tech
 - Opening the core element in Web Search: an Open Web Index
 - Open software stack and open pipelines
 - A distributed, data-centric and Web-scale compute and storage analytics infrastructure
- Build an ecosystems around data, software and infrastructure
- Transparency and control for considering legal, ethical and societal values

Thanks. Questions?



Contact us:

To keep in touch with these possibilities or to join us send an email to **join@openwebsearch.org**



We are looking for ...

→ **help hosting a distributed Open Web Index**

Data centres

Industry & business partners

→ **discover the business models of an Open Web Index**

→ **develop new search & retrieval paradigms and content analysis algorithms**

Researchers & tech innovators

Policy makers

→ **help shaping the governance of an open search ecosystem**

→ **We will also offer small grants for potential contributors.**