



Open WebSearch

# Build Your Own (Conversational) Search Engine

Hackathon Guide

Open Web Search 



Funded by  
the European Union

SUPPORTED  
BY



# Overview



- OWS Partner TU Graz (CoDiS Lab)
- Concept of Search Applications
- Technical Details of the Prototype Search Application
- Create Your Own (Conversational) Search Application

# OWS Partner: TU Graz, CoDiS Lab



- Graz University of Technology, Austria
- **Cognitive and Digital Science Lab**
- Interdisciplinary research by connecting computer science with cognitive psychology



## Contributions to ows.eu

- Search applications
- Ethical, legal, and social aspects

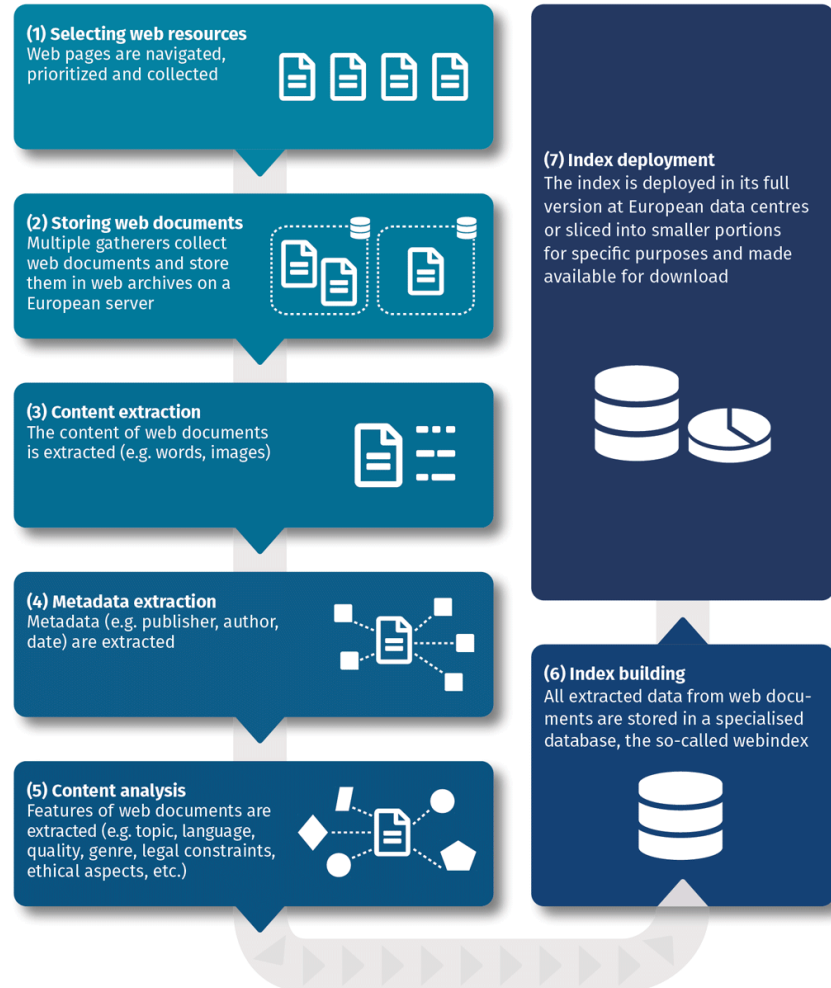


# OWS Workflow



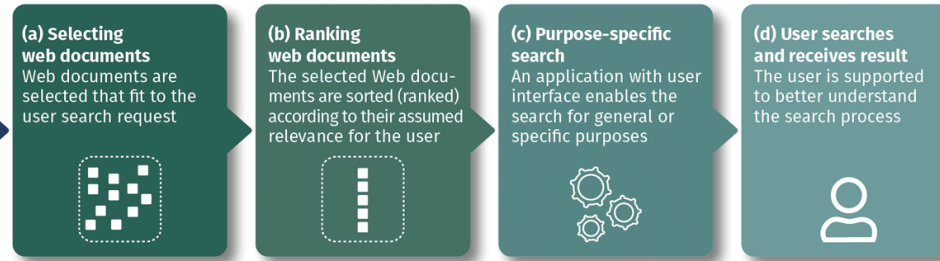
## Index Generation

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



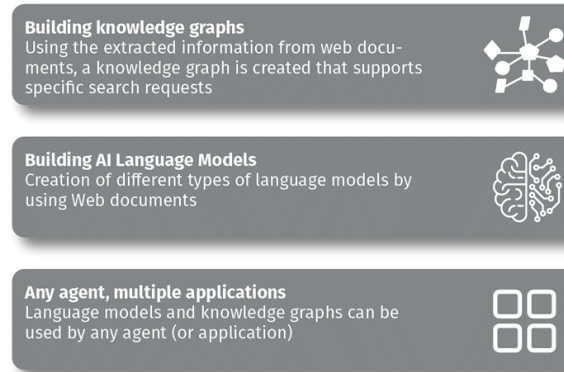
## Search Applications

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



## Data Products

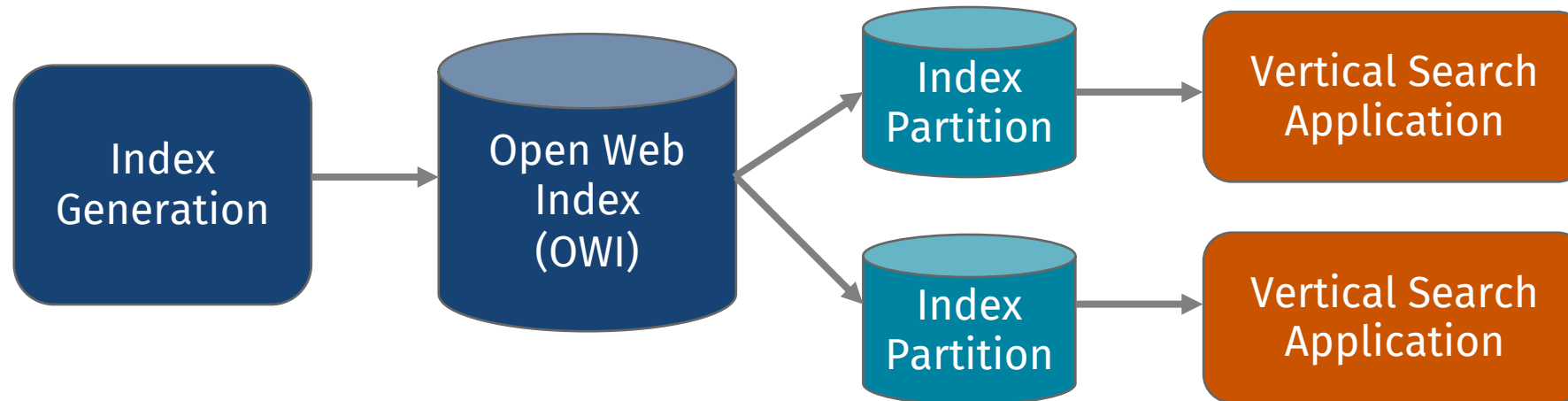
Knowledge representation models will be created using the open web index, in order to be used by any agent and for many applications



# Vertical Search Engines in OWS



- Vertical Search Engines serve special purposes, topics, or domains, e.g. product search
- Search applications are a key concept in OpenWebSearch.eu
- Based on Open Web Index and index partitions

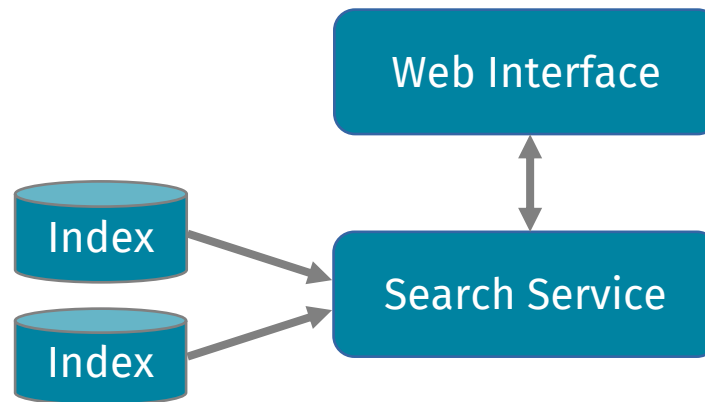


# Prototype Search Application



## Minimal implementation of an OWS vertical search engine

- Index partitions in reasonable size on a local server
- Search service that searches the index and provides REST API

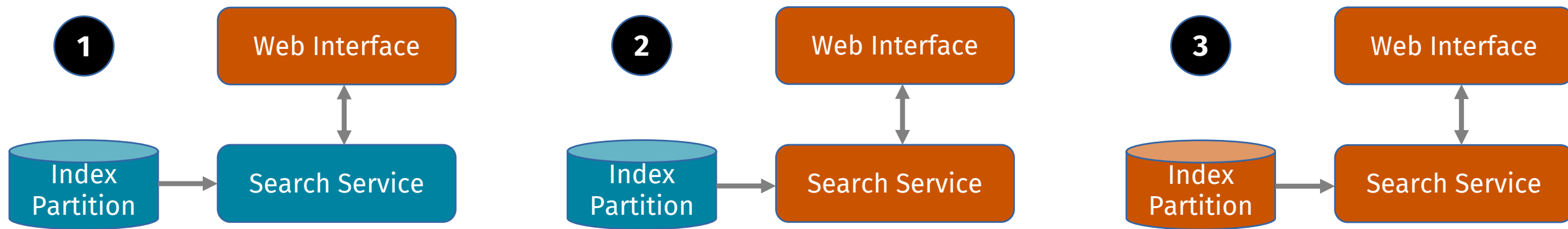


# Prototype Search Application



## Guides the development of an OWS search application

- Using the search service and development of a web interface (1)
- Using the concept and development of the whole application (2)
- Using the concept to explain the architecture of search applications (3)



# Demonstration 1: Basic Search



- Simple web interface to the Search Service
- Operates two indices
- Performs re-ranking and filtering
- Displays metadata
- Available at:  
<https://qnode.eu/ows/prosa/front-end/>

Prototype Search Application

Search term:

<b>Index:</b>	<b>Ranking:</b>	<b>Language filter:</b>	<b>Limit:</b>
<input type="radio"/> default	<input type="radio"/> default	<input checked="" type="radio"/> default	<input type="radio"/> default
<input type="radio"/> Demo Graz	<input type="radio"/> ascending	<input type="radio"/> English	<input type="radio"/> 10 items
<input checked="" type="radio"/> Demo Snapshot	<input checked="" type="radio"/> descending	<input type="radio"/> German	<input checked="" type="radio"/> 50 items
			<input type="radio"/> 1,000,000 items

Search URL: `http://localhost:8000/search?q=Graz&index=demo-snapshot&ranking=desc&limit=50`

**Search result for term: "Graz"**

*Number of retrieved items: 37 • Retrieval time: 59 ms*

**Center for Mind & Cognition**  
Schnelleinwahl mobil +496950502596,,91939791848# Deutschland +496971049922,,91939791848#  
Deutschland Einwahl nach aktuellem Standort +49 695 050 2596 Deutschland +49 69 7104 9922  
Deutschland +49 30 5679 5800 Deutschland Meeting-ID: 919 3979 1848 Orts  
*language:eng, word count:18129, index date:2023-07-01 23:52*  
<https://philosophy-cognition.com/cmcevents/event>

**[Title missing]**  
Mi perplime quali guasti possono causare taluni professori della scuola media superiore. Personalmente non ascoltavo alcuna lezione al Liceo, se non quelle del Prof. Meroni, di filosofia, e ovviamente questo non mi preclude alcuna crescita (presi 60/  
*language:eng, word count:7125, index date:2023-07-01 23:51*  
<https://www.fabiopalma.net>

**TransCoding - The Book**  
Artistic research, and with it the process of doing art and the artwork itself, are at the heart of this investigation. For that reason, I provide, as an augmented way of reading, this complementary website to



# Demonstration 2: Sightseeing Application



- Integration into existing Sightseeing Application Spots
- Clicking on a sightseeing marker triggers a search with the title of the spot as search term
- Available at:  
<https://qnode.eu/ows/prosa/spots/?mid=owsmeetinggraz>

The screenshot shows a web application interface titled "Sightseeing Graz". It features a map of Graz, Austria, with several orange location markers. A pop-up window is open over one of the markers, displaying the following information:

**Opera Building**  
Opera house and a copy of the New York Liberty Statue

Below the map, the search results for the term "Opera Building" are displayed:

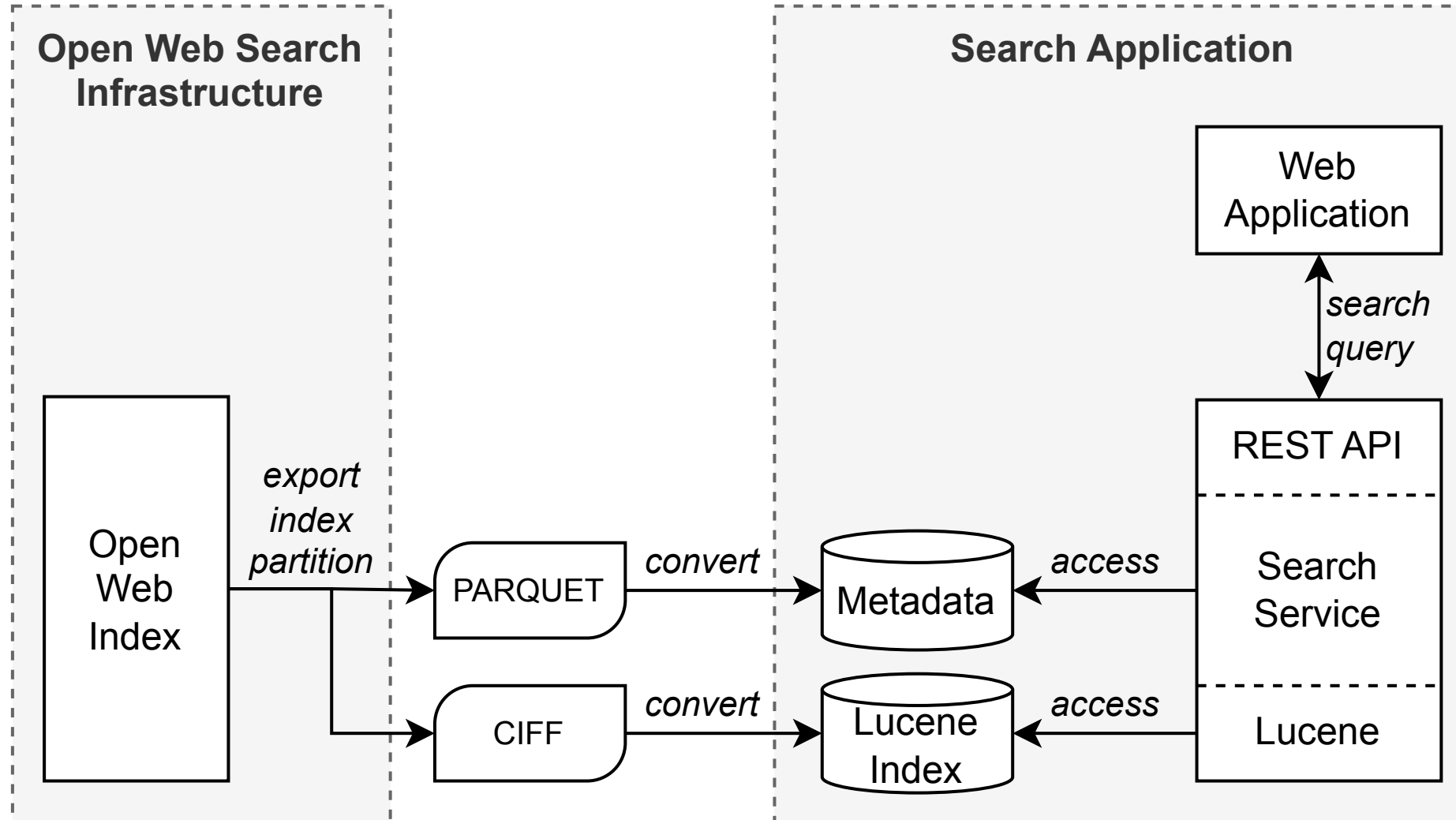
**Search result for term: Opera Building**

**Best attractions in Graz**  
The cafe, opened in a hotel with the same name, can be called a full-fledged Austrian attraction, because it is here that they serve the real Sacher cake, cooked according to the original ancient recipe. Keep in mind that there are always a lot of vi  
<https://www.tripzaza.com/en/destinations/top-attractions-in-graz>

**Opera House Graz**  
Something fascinating about Graz - the interplay of modernism and tradition - is illustrated by the sculpture "light sword" next to the opera house. It was originally made for the festival "steirischer herbst" in 1992. To celebrate the 500th annivers  
[https://www.graztourismus.at/en/sightseeing-culture/sights/opera-house\\_shg\\_1472](https://www.graztourismus.at/en/sightseeing-culture/sights/opera-house_shg_1472)

**Graz Opera**  
Past general music directors (GMD) of the company have included Niksa Bareza (1981-1990), Philippe Jordan (2001-2004), Johannes Fritsch (2006-2013), and Dirk Kaftan (2013-2017) [7] In

# Conceptual Design of the Prototype Search Application



# Index Partition: CIFF/Lucene and Parquet



- CIFF: Index exchange format
- Lucene: Inverted index
- Parquet: Tabular format

CIFF/Lucene		Parquet				
Term	List of (Document, Frequency)	Document	Lang	Full text	WARC date	...
Term 1	(docID A, 1) (docID B, 3) (docID C, 1)	docID A	en	... text ...	2023-07-01	
Term 2	(docID D, 2)	docID B	de	... text ...	2023-07-01	
Term 3	(docID E, 4) (docID F, 2)	docID C	de	... text ...	2023-07-01	
Term 4	(docID H, 2) (docID I, 1)	docID D	en	... text ...	2023-07-01	

# Parquet File / Metadata Fields



→ `id`

→ `record_id`

→ `title`

→ `plain_text`

→ `warc_date`

→ `warc_ip`

→ `language`

→ `http_content_type`

→ `http_server`

→ `url`

→ `url_scheme`

→ `url_path`

→ `url_params`

→ `url_query`

→ `url_fragment`

→ `url_subdomain`

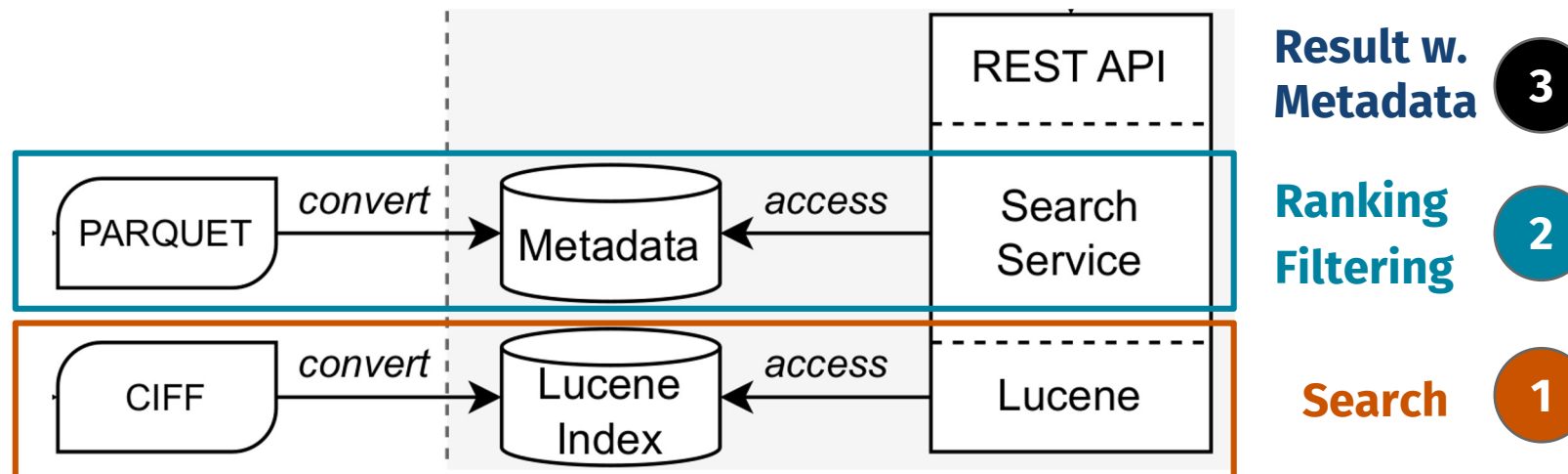
→ `url_domain`

→ `url_suffix`

# Search Process



- CIFF file is converted to a Lucene index that can be searched with the Lucene library
- Search result is filtered and ranked using metadata coming from the Parquet file



# REST API



End point: `host:port/search?q=...`

## Request parameters:

Parameter	Description
<code>q</code>	Search term(s)
<code>index</code>	Lucene index to be searched in
<code>lang</code>	Restrict searches to pages in language
<code>ranking</code>	Specify order of search result
<code>limit</code>	Set maximum number of returned results

## Example:

<https://qnode.eu/ows/prosa/service/search?q=austria>

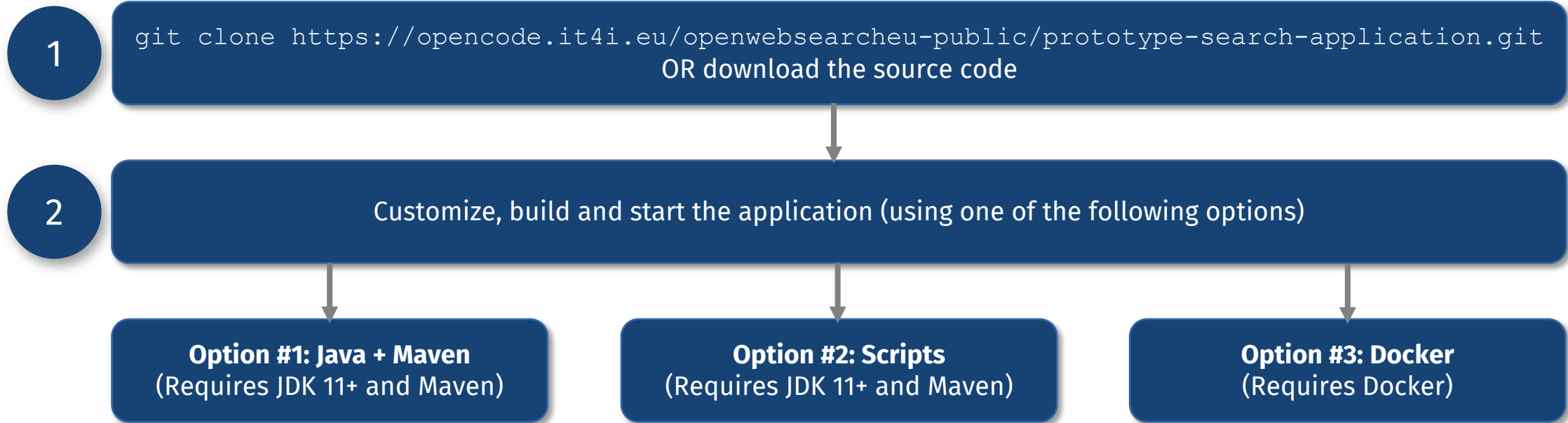
## Response data:

```
{
  results: [
    {
      id: ...,
      url: ...,
      title: ...,
      textSnippet: ...,
      language: ...,
      wordCount: ...,
      warcDate: ...,
    }, ...
  ]
}
```

# Next Steps Towards Creating Your Own Search Application



→ Find the source code of the Prototype Search Application on GitLab:  
<https://opencode.it4i.eu/openwebsearcheu-public/prototype-search-application>



# Option #1: Java + Maven



→ Prerequisites: Java (JDK 11+) + Maven

```
# change the directory  
cd search-service  
# clean the project and compile the source code  
mvn clean compile  
# build the executable  
mvn package  
# run the executable  
java -jar target/app.jar -d <YOUR_DEFAULT_INDEX_NAME> -p <YOUR_API_PORT>
```

→ **Example:** `java -jar target/app.jar -d demo-graz -p 8000`



# Option #2: Scripts



→ Prerequisites: Java (JDK 11+) + Maven

```
# change the directory  
cd scripts  
# convert indexes(optional), clean the project and create a packaged JAR  
./build.sh  
# run the executable  
./start.sh
```

→ Example: `./start.sh demo-graz 8000` (default index and port are optional)

→ Batch files as alternative for Windows

# Option #3: Docker



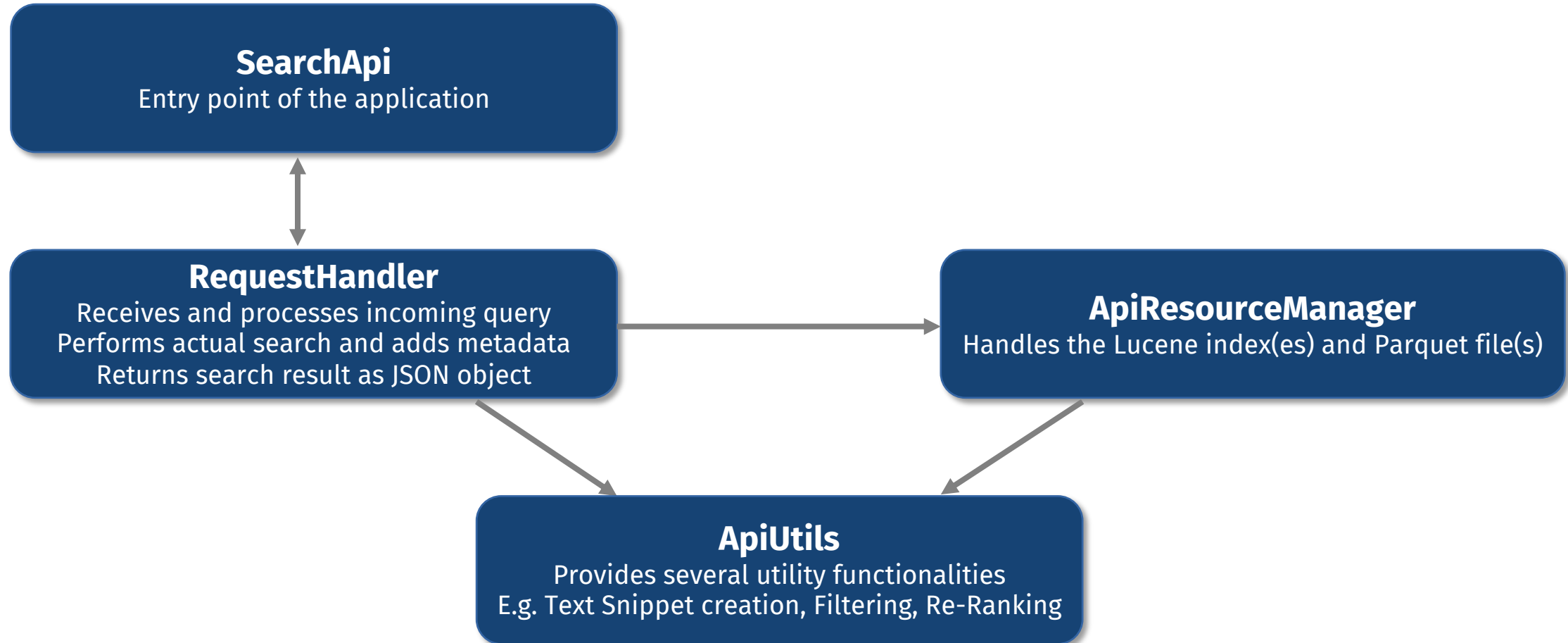
## → Prerequisites: Docker

```
# create an image
cd scripts docker build -t prototype-search-application .
# start a container
docker run -p 8000:8000 prototype-search-application \
    -d <YOUR_DEFAULT_INDEX_NAME> -p <YOUR_API_PORT>
```

→ **Example:** `docker run -p 8000:8000 prototype-search-application \`  
`-d demo-graz -p 8000`

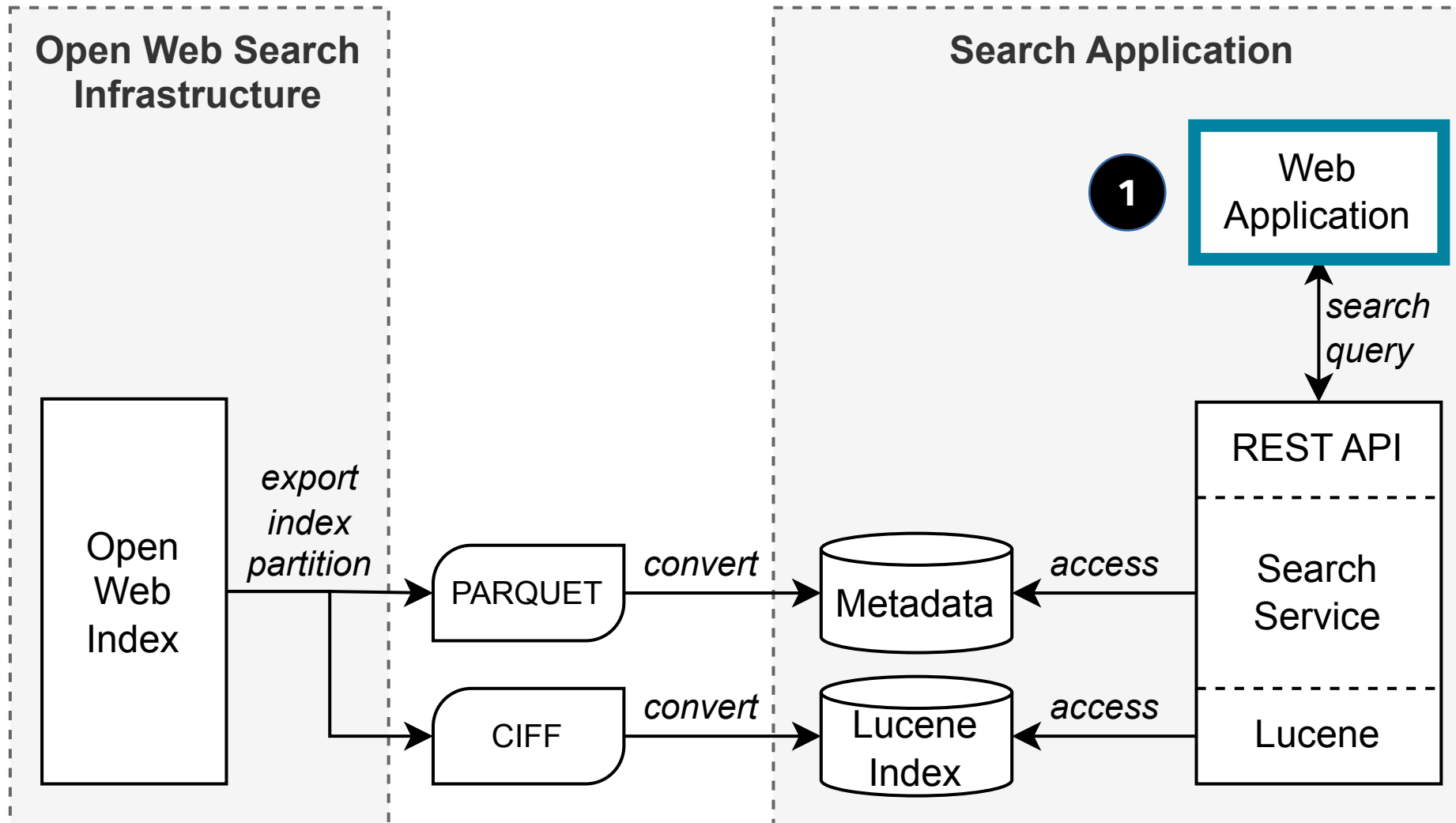
→ **Re-build of image necessary after every code change**

# Code Architecture of the Prototype Search Application



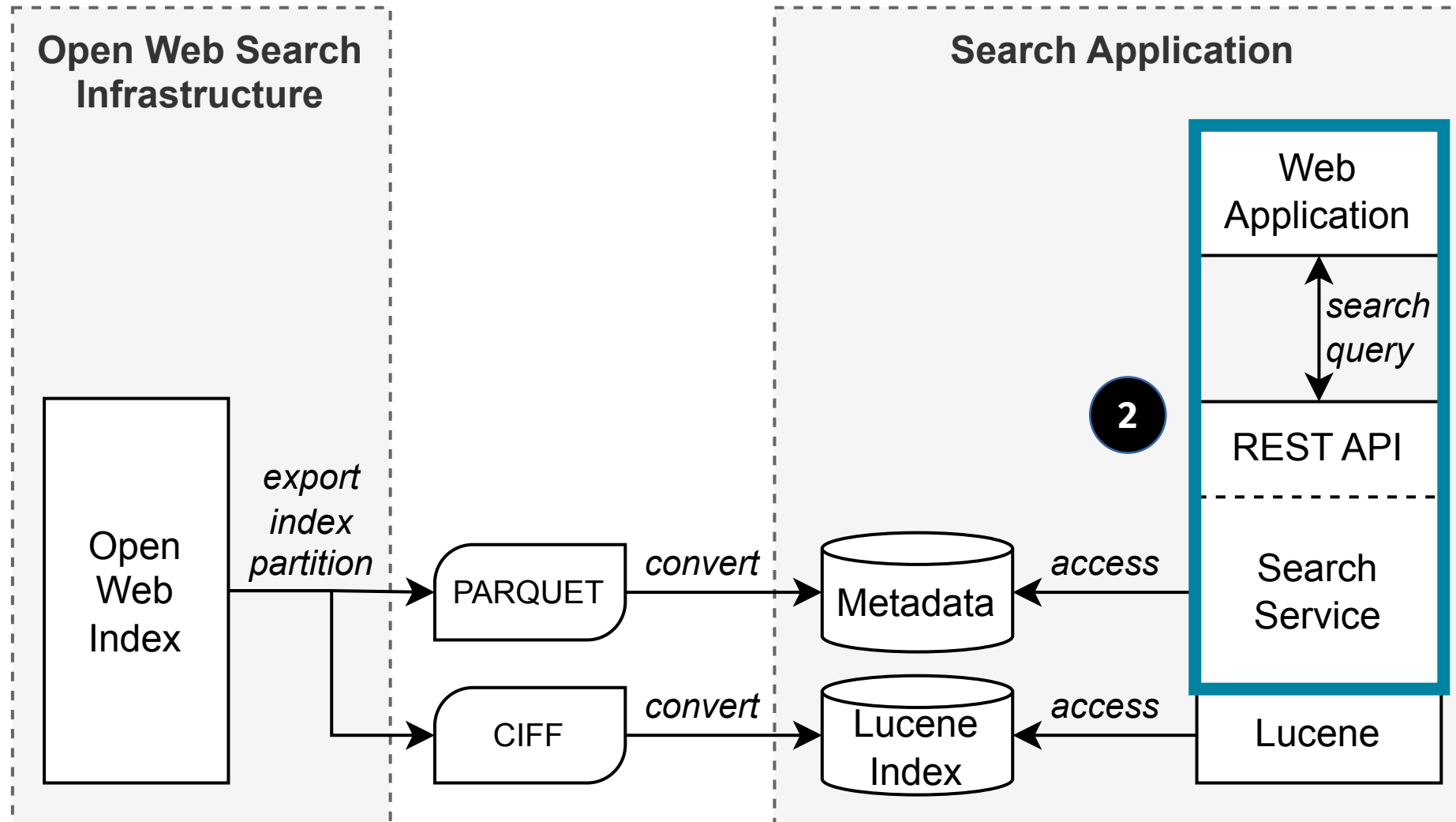
# Ideas & Suggestions on What You Can Change

## (1) Using the search service and development of a web interface



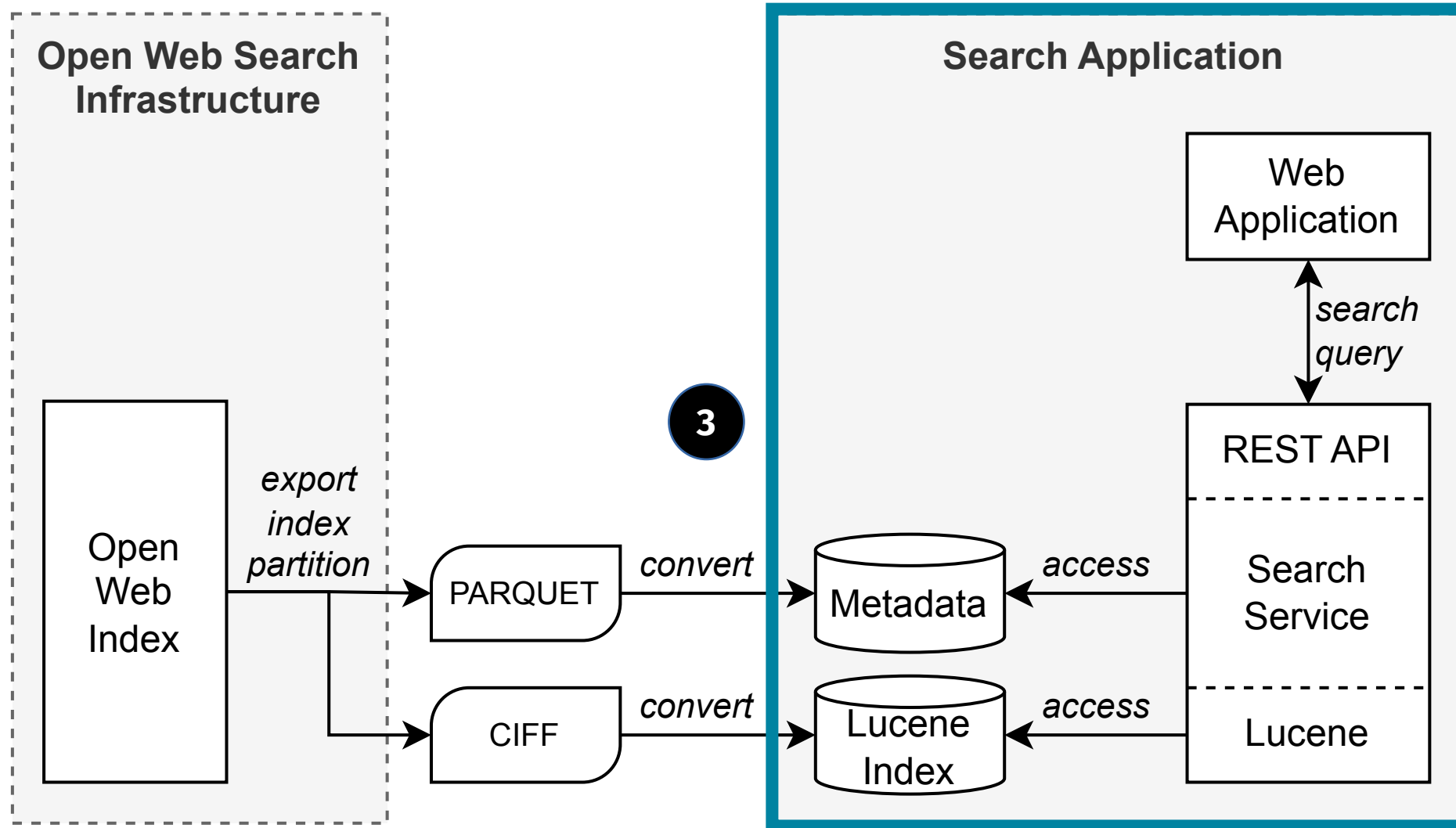
# Ideas & Suggestions on What You Can Change

## (2) Using the concept and development of the whole application



# Ideas & Suggestions on What You Can Change

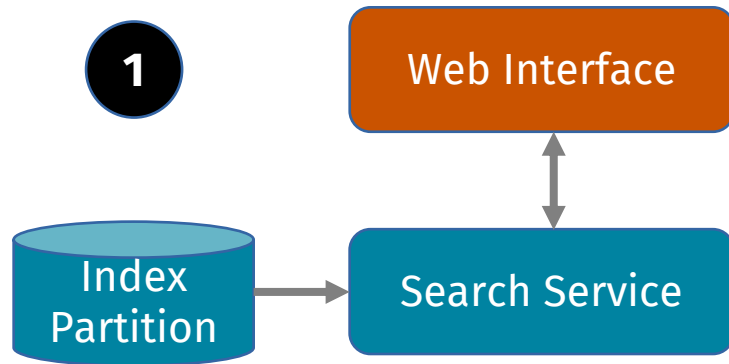
(3) Using the concept to explain the architecture of search applications



# Ideas & Suggestions on What You Can Change



(1) Using the search service and development of a web interface

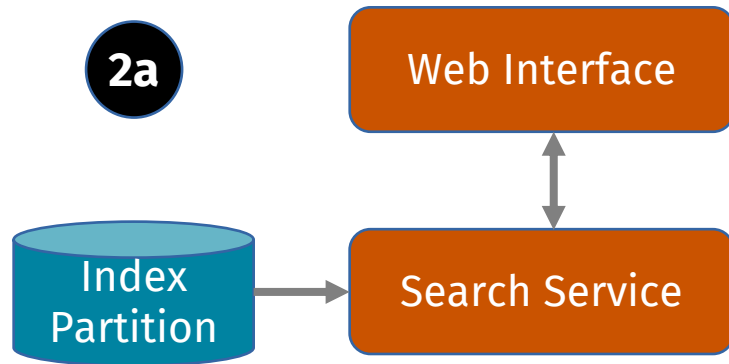


## Web / User Interface:

- Change the existing web search interface
- Create a new web search interface from scratch
- Example: Replace list representation of search results
- Example: Re-build a sightseeing application for Passau using OpenStreetMap

# Ideas & Suggestions on What You Can Change

(2) Using the concept and development of the whole application



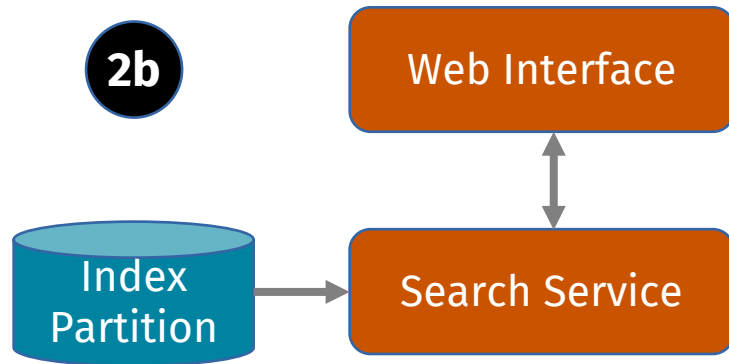
## API Parameters

- Add additional parameters or change the functionality of existing parameters
- Example: Filter by Top-Level-Domains (as API parameter and in web interface)



# Ideas & Suggestions on What You Can Change

(2) Using the concept and development of the whole application

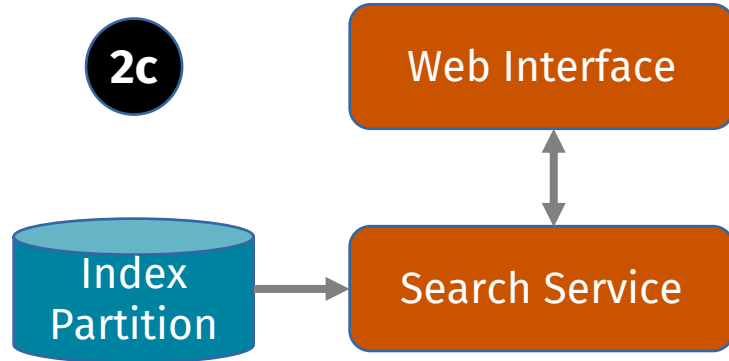


## Re-Ranking

- Perform re-ranking not (only) based on word count of the document
- Example: Add additional re-ranking mechanism(s) that can be selected for re-ranking in web search interface
- Example: Further analysis of full text (e.g., perform a sentiment analysis)

# Ideas & Suggestions on What You Can Change

(2) Using the concept and development of the whole application



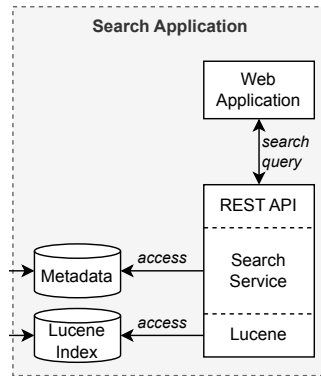
## Text Snippet

- Use a more meaningful text snippet than just the beginning of the document
- Example: Use a text snippet from the full text that contains the queried term at least once
- Example: Automatically generate a summary of the full text

# Ideas & Suggestions on What You Can Change

(3) Using the concept to explain the architecture of search applications

3a



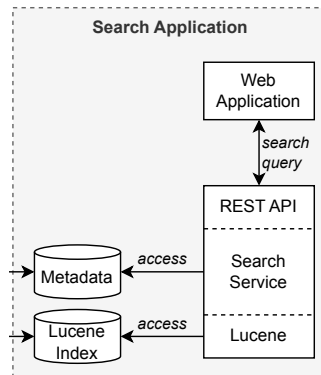
## Conceptual Design (Software Architecture)

→ Design an alternative concept based on the existing software architecture or from scratch

# Ideas & Suggestions on What You Can Change

(3) Using the concept to explain the architecture of search applications

3b



## Conceptual Design (Use Cases)

- Design (and implement) a concept of a use case
- Example: Conversational Search Application using an LLM API (e.g., ChatNoir Chat [1], OWLer GPT [2])

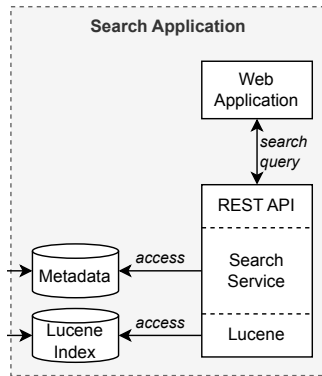
[1] <https://chat.web.webis.de/>

[2] <https://owlergpt.pads.fim.uni-passau.de/>

# Ideas & Suggestions on What You Can Change

(3) Using the concept to explain the architecture of search applications

3c



## Conceptual Design (Use Cases)

- Design (and implement) a concept of a use case
- Example: Geo-annotation based on full text (e.g., using the API of [geonames.org](https://www.geonames.org/))
- Example: Topic extraction based on full text (e.g., search and/or filter by topic)

# Entry Points In The Code



- `MARKER-PARAMETERS`
- `MARKER-LUCENE-SEARCH`
- `MARKER-QUERY-PARSING`
- `MARKER-METADATA-ENRICHMENT`
- `MARKER-LANGUAGE-FILTER`
- `MARKER-RERANKING`
- `MARKER-TEXT-SNIPPET`

```
/*  
 * MARKER-<ID> *  
 * Some description. *  
***/
```

A description for each marker can be found in the source code

# Any Questions?



## Contact:

- Sebastian Gürtl <[sebastian.guertl@tugraz.at](mailto:sebastian.guertl@tugraz.at)>
- Alexander Nussbaumer <[alexander.nussbaumer@tugraz.at](mailto:alexander.nussbaumer@tugraz.at)>
- Graz University of Technology, Austria

