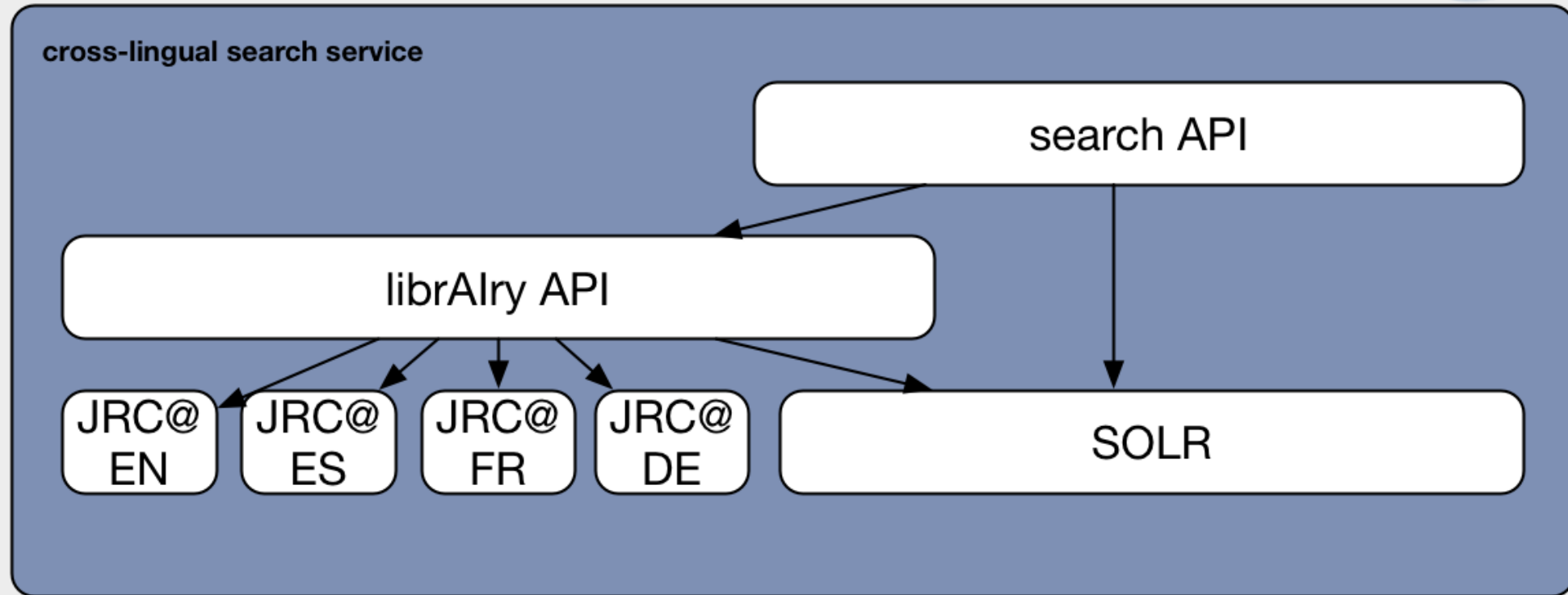# Search-API

**Carlos Badenes-Olmedo,** Francisco Yedro, Oscar Corcho

Ontology Engineering Group (OEG)

Universidad Politécnica de Madrid (UPM)

Spain

# Motivation / Proposal

- We got a *huge* collection of (un-labelled) documents and we would like to *explore* the knowledge inside.

- Imagine that we could run an *unsupervised*, *automated* pipeline to generate *connections* between them

- state-of-the art techniques to programmatically generating *annotations* for each of the texts inside big collections of documents

- in a way that is *computationally affordable*

# Architecture

# Deployment

- https://github.com/TBFY/search-API

- the search-API, librAIry, Solr and all the annotation models are distributed as ***Docker images***

- ***Docker-Engine*** and ***Docker-Compose*** are only required to deploy the whole system

- The service descriptor to run it is:

```
$docker-compose up
```

```yaml
01. version: '3'
02. services:
03.   nlp:
04.     image: librairy/nlp:1.3
05.     environment:
06.       - REST_PATH=/
07.       - JAVA_OPTS=-Xmx1024m
08.   jrc-en-model:
09.     image: librairy/jrc-en-model:1.3
10.     ports:
11.       - "8085:7777"
12.     environment:
13.       - REST_PATH=/
14.       - NLP_ENDPOINT=nlp
15.       - JAVA_OPTS=-Xmx128m
16.   ...
17.   solr:
18.     image: solr:7.7
19.     ports:
20.       - "8983:8983"
21.     volumes:
22.       - ./target/solr-data:/opt/solr/server/solr/mycores
23.       - ./src/main/banana:/opt/solr/server/solr-webapp/webapp/banana
24.     entrypoint:
25.       - docker-entrypoint.sh
26.       - solr-precreate
27.       - documents
28.     environment:
29.       - SOLR_JAVA_MEM=-Xms512m -Xmx512m
30.   librairy-api:
31.     image: librairy/api:1.3
32.     ports:
33.       - "8081:7777"
34.     environment:
35.       - LIBRAIRY_API_USERS=oeg:oeg2018
36.       - JAVA_OPTS=-Xmx128m
37.       - REST_PATH=/librairy-api
38.     volumes:
39.       - ./tmp:/librairy
40.       - /var/run/docker.sock:/var/run/docker.sock
41.   search-api:
42.     image: librairy/search-api:1.1
43.     ports:
44.       - "8080:7777"
45.     environment:
46.       - LIBRAIRY_API_USERNAME=oeg
47.       - LIBRAIRY_API_PASSWORD=oeg2018
48.       - LIBRAIRY_API_ENDPOINT=http://librairy-api:7777/librairy-api
49.       - MODEL_ENDPOINT=http://jrc-%%-model:7777
50.       - SOLR_ENDPOINT=http://solr:8983/solr/documents
51.       - JAVA_OPTS=-Xmx128m
52.       - REST_PATH=/search-api
```

# Deployment

| SERVICE | LOCAL-ENDPOINT | REMOTE-ENDPOINT |
|---|---|---|
| Search-API | http://localhost:8080/search-api | http://tbfy.librairy.linkeddata.es/search-api |
| librAIry-API | http://localhost:8081/library-api | http://library.linkeddata.es/api |
| Solr | http://localhost:8983/solr | http://library.linkeddata.es/data |
| Dashboard | http://localhost:8983/solr/banana | http://library.linkeddata.es/data/dashboard |
| JRC-English-Model | http://localhost:8085 | http://library.linkeddata.es/jrc-en-model/ |
| JRC-Spanish-Model | http://localhost:8086 | http://library.linkeddata.es/jrc-es-model |
| JRC-French-Model | http://localhost:8087 | http://library.linkeddata.es/jrc-fr-model |
| JRC-German-Model | http://localhost:8088 | http://library.linkeddata.es/jrc-de-model |

# Actions

- **Add new documents**
  - (1) from structured text files
  - (2) from existing Solr indexes
  - (3) from remote APIs
  - (4) from raw files
- **Annotate existing documents**: *with main topics*

**librAlry-API**
http://librairy.linkeddata.es/api

- **Search documents**: *by name, (or) source, (or) language, (or) words*
- **Read documents**: *by id*
- **Retrieve similar documents**: *by name, (or) source, (or) language, (or) words*
  - (1) from a document id
  - (2) from free-text

**Search-API**
http://tbfy.librairy.linkeddata.es/search-api

# A1: Add New Documents
*from structured text files*

```
01.   Company;;Bergan Mercy Medical Center;; Bergan Mercy
02.   Company;;The Unsigned Guide;; The Unsigned Guide is
03.   Company;;Rest of the world;; Within sports and games
```

- ***librAIry-api*** *is the responsible for ingesting documents*

- ***CSV*** or ***JSON-L*** text files (even in ***{tar}.gz*** format) can be used

- A HTTP_POST request to:
  http://librairy.linkeddata.es/api/documents

```
01.   {
02.       "contactEmail": "cbadenes@fi.upm.es",
03.       "dataSink": {
04.           "format": "SOLR_CORE",
05.           "url": "http://library.linkeddata.es/data/tbfy"
06.       },
07.       "dataSource": {
08.           "name": "entities",
09.           "dataFields": {
10.               "id": "1",
11.               "labels": ["0"],
12.               "text": ["2"]
13.           },
14.           "filter": ";;",
15.           "format": "CSV_TAR_GZ",
16.           "offset": 0,
17.           "size": -1,
18.           "url": "https://bit.ly/2Y7dkoY"
19.       }
20.   }
```

# A1: Add New Documents
## *from Solr indexes*

- ***librAIry-api*** *is the responsible for ingesting documents*

- The documents can already be hosted in a ***Solr collection***

- A HTTP_POST request to:
  http://librairy.linkeddata.es/api/documents

```
01.    {
02.        "contactEmail": "cbadenes@fi.upm.es",
03.        "dataSink": {
04.            "format": "SOLR_CORE",
05.            "url": "http://librairy.linkeddata.es/data/tbfy"
06.        },
07.        "dataSource": {
08.            "name":"jrc",
09.            "dataFields": {
10.                "id": "id",
11.                "name": "name_s",
12.                "labels": [
13.                    "labels_t"
14.                ],
15.                "text": [
16.                    "txt_t"
17.                ]
18.            },
19.            "filter":"size_i:[100 TO 10000] && source_s:jrc
20.                && lang_s:es && labels_t:[* TO *]",
21.            "format": "SOLR_CORE",
22.            "offset": 0,
23.            "size": -1,
24.            "url": "http://librairy.linkeddata.es/solr/documents"
25.        }
26.    }
```

# A1: Add New Documents
## *from remote APIs*

- documents can be available from a remote **HTTP_Restful API** (e.g OpenOpps-API)

- or stored on remote **storage systems** (e.g. Amazon S3)

- a **Harvester** client has been created to handle these sources

- It **download**, **parse** and **store** the resources from external sources into Solr collection with a predefined format

- GitHub project: https://github.com/TBFY/harvester

# A1: Add New Documents
## *from raw files*

- There may be exceptional circumstances that do not conform to the above scenarios

- Since the search service uses Solr internally to store its documents, the **Solr API** can be used directly to add new documents: http://librairy.linkeddata.es/data/#/tbfy

- The only requirement in this scenario is to use the following **fields** to describe the documents:

| Field | Description |
|---|---|
| id | Unique identifier |
| name_s | Document title |
| txt_t | Textual content |
| size_i | Number of characters |
| format_s | Original format (e.g. json, PDF, xml..) |
| lang_s | Language code (ISO 639-1) |
| source_s | Origin identifier (e.g. ted, jrc) |
| date_dt | Publication date (ISO 8601) |

# A2: Annotate Documents

- **librAIry-api** *is the responsible for annotate documents*

- The documents should be hosted in a **Solr collection**

- The similarity calculation is based on these annotations.

- A HTTP_POST request to:
  http://librairy.linkeddata.es/api/annotations

*json request*

```
01.  {
02.      "contactEmail": "cbadenes@fi.upm.es",
03.      "dataSink": {
04.          "format": "SOLR_CORE",
05.          "url": "http://librairy.linkeddata.es/data/tbfy"
06.      },
07.      "dataSource": {
08.          "dataFields": {
09.              "id": "id",
10.              "name": "name_s",
11.              "labels": ["labels_t"],
12.              "text": ["txt_t"]
13.          },
14.          "filter": "-topics0_t:[* TO *] && lang_s:en",
15.          "format": "SOLR_CORE",
16.          "offset": 0,
17.          "size": -1,
18.          "url": "http://librairy.linkeddata.es/data/tbfy"
19.      },
20.      "modelEndpoint": "http://librairy.linkeddata.es/jrc-en-model"
21.  }
```

# A3: Search Documents

- ***search-api*** *is the responsible for searching documents*

- Swagger documentation is available at:
  http://tbfy.librairy.linkeddata.es/search-api

- Document retrieval can be done via HTTP_GET requests filtered by query parameters:

| Parameter | Filtered by |
|-----------|-------------|
| lang | language |
| name | words contained in the document title |
| source | document origin |
| text | words contained in the document |
| size | maximum number of documents |
| cursor | first index |

```
1.  [
2.    {
3.      "id": "18-583998-001-es",
4.      "language": "es",
5.      "name": "Provisión de servicios de consultoría relativos a servicios de gestión
         de instalaciones (FM)",
6.      "source": "ted"
7.    },
8.    {
9.      "id": "18-587627-001-es",
10.     "language": "es",
11.     "name": "Batimetría: trazado de mapas de alta resolución del fondo marino",
12.     "source": "ted"
13.   }
14. ]
```

# A4: Read Documents

- **search-api** *is the responsible for reading documents*

  *http://tbfy.librairy.linkeddata.es/search-api/documents/18-590660-001-en*

- Swagger documentation is available at:
  http://tbfy.librairy.linkeddata.es/search-api

- Following the Restful principles, a document can be read by making a HTTP_GET request to the URI containing the document identifier

```
1.  {
2.      "id": "18-590660-001-en",
3.      "name": "Occupational Safety and Health (OSH) Services",
4.      "text": "The European Central Bank (ECB) is seeking through this open procedure a
         supplier for the provision of Occupational Safety and Health Services (OSH) and in
         tends to award contract to the supplier offering the best value for money. OSH matt
         ers are regarded as integral components of every task and function at the ECB. In t
         his regard, the ECB aims at providing a modern, ergonomic and healthy working envir
         onment that meets the relevant requirements and generally accepted technical and st
         ructural OSH standards and, in doing so, aims at minimising occupational accidents
         and injuries. The European Central Bank (ECB) is seeking through this open procedur
         e a supplier for the provision of Occupational Safety and Health Services (OSH) and
         intends to award contract to the supplier offering the best value for money. OSH m
         atters are regarded as integral components of every task and function at the ECB. I
         n this regard, the ECB aims at providing a modern, ergonomic and healthy working en
         vironment that meets the relevant requirements and generally accepted technical and
         structural OSH standards and, in doing so, aims at minimising occupational acciden
         ts and injuries.",
5.      "format": "xml",
6.      "language": "en",
7.      "source": "ted",
8.      "date": "Mon Jan 07 23:00:00 GMT 2019",
9.      "tags": "82 1016 929"
10. }
```

# A5: Retrieve Similar Documents
## from a document

- **search-api** *is the responsible for retrieve similar documents*

- documents are sorted according to their content similarity to the reference

- the request can be filtered by:

| Parameter | Filtered by |
|---|---|
| lang | language |
| name | words contained in the document title |
| source | document origin |
| terms | words contained in the document |
| size | Maximum number of documents |

*http://tbfy.librairy.linkeddata.es/search-api/documents/18-590660-001-en/items*

```
01.   [{
02.           "id": "19-052340-001-en",
03.           "name": "Occupational Safety and Health (OSH) Services",
04.           "score": 5775.54638671875
05.       },
06.       {
07.           "id": "19-052340-001-de",
08.           "name": "Dienstleistungen im Bereich Sicherheit und
09.                           Gesundheitsschutz am Arbeitsplatz (Occupational
10.                           Safety and Health – OSH)",
11.           "score": 5756.74755859375
12.       },
13.       {
14.           "id": "18-590660-001-es",
15.           "name": "Servicios de salud y seguridad en el trabajo",
16.           "score": 5756.74755859375
17.       }
18.   ]
```

# A6: Retrieve Similar Documents
*from free-text*

*JSON request*

```
01.  {
02.      "text": "Fast food restaurants can also face claims over food
03.                        allergies. Ingredient lists must be comprehensive to
04.                        allow people with allergies to avoid consuming foods
05.                        that will make them sick.",
06.      "size": 10,
07.      "lang": "es",
08.      "source": "jrc"
09.  }
```

- **search-api** *is the responsible for retrieve similar documents*

- Now a HTTP_POST request to
  http://tbfy.librairy.linkeddata.es/search-api/items

*JSON response*

```
01.  [
02.    {
03.      "id": "jrc32002L0067-en",
04.      "name": "Commission Directive 2002/67/EC of 18 July 2002 on the labelling
05.  of foodstuffs containing quinine, and of foodstuffs containing caffeine
06.  (Text with EEA relevance)",
07.      "score": 6640.35693359375
08.    },
09.    {
10.      "id": "jrc31997R0258-en",
11.      "name": "Regulation (EC) No 258/97 of the European Parliament and of the
12.  Council of 27 January 1997 concerning novel foods and novel food ingredients",
13.      "score": 6633.22314453125
14.    }
15.  ]
```

| Parameter | Filtered by |
|-----------|-------------|
| lang | language |
| name | words contained in the document title |
| source | document origin |
| terms | words contained in the document |
| size | Maximum number of documents |

# Next Steps

- develop new data-source **adapters** if required (e.g *Elasticsearch)*

- automatic documents **uploading** (daily from OpenOpps-API)

- improve **quality** of models (mainly french and german) for short texts

- create **notebooks** that show how:
    - to search/read documents about public procurement
    - to explore a multilingual collection of legal documents
    - to retrieve similar contracts/news/.. from a given text

# Search-API

**Carlos Badenes-Olmedo,** Francisco Yedro, Oscar Corcho

Ontology Engineering Group (OEG)

Universidad Politécnica de Madrid (UPM)

Spain