# Data Mosaics: standards and prototype for a data mashup platform

## D2 – Demonstrator Documentation

**Author: Philippe Duchesne**

Joint Research Centre

# 1. Introduction

This prototype covers the implementation of the three functional elements that constitute the complete mosaic solution:

1.  a web service acting as a catalogue of mosaics;
2.  a representation format for the exchange of mosaics, based on the model described above; and
3.  a web application able to display and author mosaics in that representation format.

The Operations Manual section below documents the functionalities of the web application. The overall architecture, the implementation details and the installation instructions are described in the Technical Manual section.

# 2. Operations Manual

## 2.1. Search

The Search panel is the welcome page of the application. It can also be accessed using the 'Search' button in the top menu (Figure 1, a).

Search terms can be entered in the text input field (Figure 1, b). In the absence of a search value, the result list displays recently modified mosaics.

When a search term is entered, the result list displays separately the mosaics, the resources and the annotations whose title contains the search term.

If the search term is a URL, the search engine will query against the resource URL or the annotations source and target URLs.



**Figure 1. Search repository**

Clicking on a result will open the containing mosaic, and in the case of a resource or an annotation, focus on it.
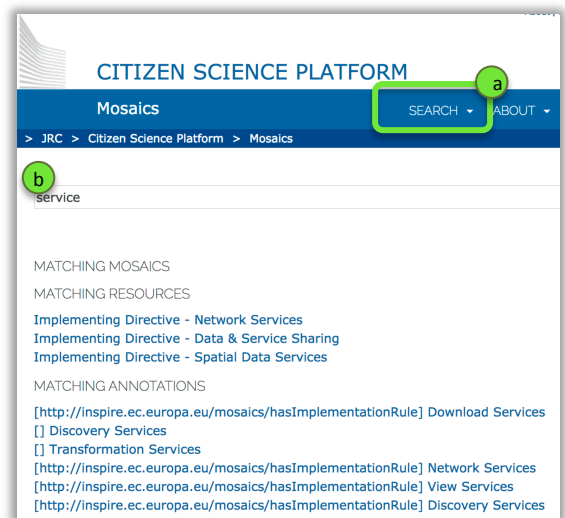
## 2.2. Browse mosaics



**Figure 2. Browse interface overview**

The 'Browse Mosaics' panel is accessible via the top menu item 'All Mosaics', in the Search dropdown (Figure 1, a). Its main components are (Figure 2):

a.  Mosaic list
b.  Current mosaic resource list
c.  Main resource view
d.  Highlighted annotated fragments with a contextual menu
e.  Scroll bar with quick links to annotations
f.  Secondary views, displaying resources linked to the main resource

The mosaic and resource lists provide edit and delete actions for every item.

Upon selection of a resource, the resource is displayed in the main view, and the annotations to/from that resource are overlaid on the displayed resource, and listed under the resource in the resource list panel (Figure 2, g).

Clicking on an annotation in the main view open a secondary view to display the linked resource, and focuses on the fragment of that resource (if any).

Hovering over a resource view displays an action bar in the top left corner (Figure 3). These actions are, from left to right:



**Figure 3. Resource action bar**

• Toggle graph view
• Toggle annotations display
• Open resource in a separate browser tab
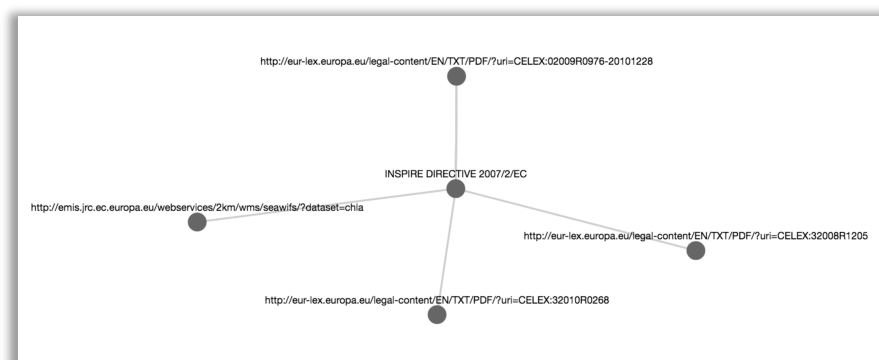
### 2.2.1. Graph view



**Figure 4. Graph View**

When switching to Graph View in the resource action bar, the main panel displays a graph view of the resource and its linked annotations (Figure 4).

## 2.3. Authoring

### 2.3.1. Steps to create a mosaic



**Figure 5. Adding a mosaic**

1. Create a new mosaic by clicking the 'Add Mosaic' action in the mosaic list (Figure 5).
2. Fill the "Edit Mosaic" dialog with the title of the mosaic.

### 2.3.2. Create a resource in a mosaic

When in the scope of a mosaic, a resource can be added to the current mosaic with the following steps:



**Figure 6. Adding a resource**

1. Click the 'Add Resource' action in the resource list (Figure 6) to bring the "Edit Resource" dialog.

   a. In the resource edition dialog (Figure 7), the URL of the resource can be typed manually, or can be picked from the search results on the catalogues that have been configured in the platform (Figure 7, a).



**Figure 7. Looking for a resource**

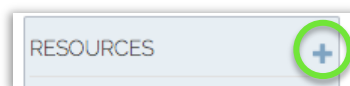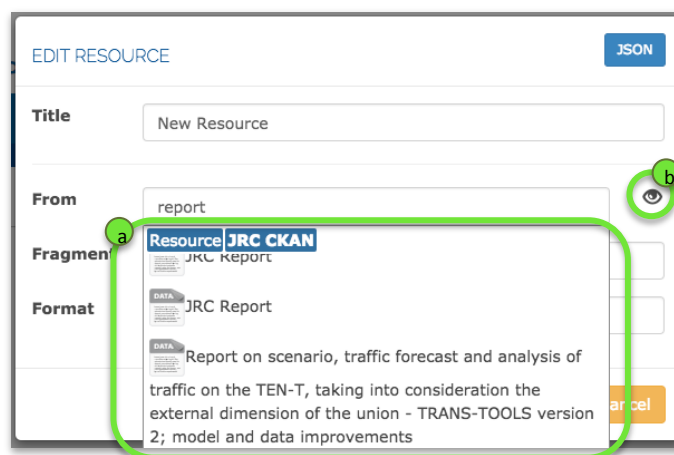b. Optionally, fill in the 'format' field with a mime type. If missing, the mime type will be guessed from the URL, the HTTP response headers or other criteria specific to each resource type. If selecting the resource from a catalogue result, this field may be set using the metadata of the catalogue.

c. the selected resource can be previewed by clicking on the 'eye' icon (Figure 8, b).



Figure 8. Selection fragment in preview

d. Once previewed, a **fragment can be selected** directly within the preview, by clicking on the plus sign that appears next to a selection (Figure 8).

The way to select a fragment depends on the type of resource being previewed: for HTML or PDF documents, it is done by selecting text with the pointer. For geospatial resources it is done by clicking the 'Draw box' icon in the toolbar (**Figure 9**, a). This allows the drawing of a box and brings the 'plus' icon (**Figure 9**, b). When clicking on it, it fills the fragment field with the corresponding fragment value (**Figure 9**, c).
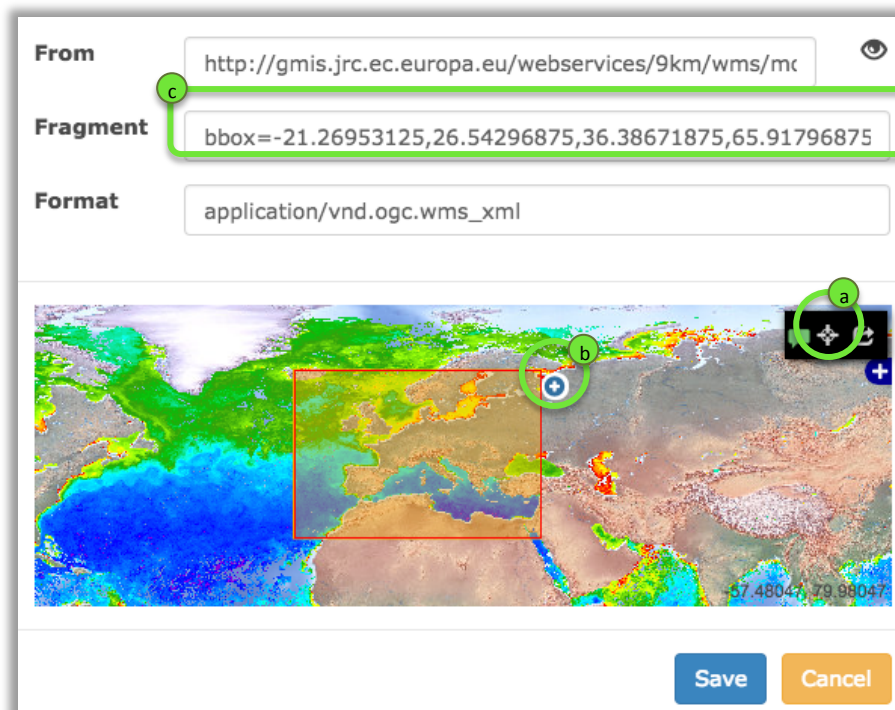


Figure 9. Geospatial fragment selection

### 2.3.3. Annotate a resource

While viewing a resource in the main interface, an annotation can be created with the following steps:

1. Select a fragment of the resource. For text-based documents such as HTML or PDF, this is done by selecting text with the pointer. Other types of resources have ad-hoc menu items to trigger selection.



**Figure 10. Annotate a text fragment**

2. Once a fragment is selected visually, a 'plus' signs appears next to the selection (Figure 10). Click it to bring the Annotation Edit dialog.

The Edit Annotation Dialog has the following (group of) fields (Figure 11).

   a. Title
   b. Link type
   c. Source resource and fragment thereof
   d. Target resource and fragment thereof

The Source fields are pre-filled with the fragment being annotated.

Source and target fields have the same search and preview mechanisms as in the Edit Resource dialog.

Link type is the optional semantic type of the annotation; typing a value will search through the configured vocabularies in the platform.



**Figure 11. Edit Annotation Dialog**

The 'switch' action (Figure 11, e) permutes the source and target resources, thereby inversing the direction of the annotation.

# 3. Technical Manual

This project involved both a modelling and standardization phase, and a prototyping phase. The modelling aimed at defining an abstract model for data mosaics and a set of vocabularies. The prototyping phase focused on the implementation of tools to author and consume data mosaics.

## 3.1. Design

### 3.1.1. Mosaic object model

A first attempt at a (Data) Mosaic model (Figure 12) consists of a collection of resource references, along with a collection of annotations that link those resources together. The resources constitute the central documents meant to be overlaid with annotations.

Both, resources or annotations are optional: a Mosaic can consist of annotations alone (being a pure set of stitched fragments with no background resource) or of resources alone (similar to a bookmarks list). As an example, a direct mapping of these requirements to a JSON representation can be:



**Figure 12. Mosaic Object Model**

```
{
    "id"    : "http://example.org/mosaic1",
    "type"  : "Mosaic",
    "title" : "Sample Mosaic",
    "resources": [
       {"id"    : "http://example.org/resource1"
        "type"  : "Resource",
        "title" : "DIRECTIVE 2007/2/EC",
        "url"   : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
        "mimeType" : "application/pdf"
       }
    ],
    "annotations": [
       {"id"    : "http://example.org/annotation1",
        "type"  : "Annotation",
        "title" : "Discovery Services",
        "from"  : "http://eur-
lex.europa.eu/.../?uri=CELEX:32007L0002#page=7&xywh=323,502,241,45",
        "to"    : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976#page=4"
}
```

where one PDF document (`...?uri=CELEX:32007L0002`) is annotated by one annotation that points at another fragment of another PDF document (`...?uri=CELEX:32009R0976`).

When it comes to the modelling of annotations, the Web Annotation Model (described in D1, section 3.2) offers a wider expressivity and range of properties. An example of use of the Web Annotation Model illustrating only the set of properties defined above (i.e. fromUrl, toUrl, text, type), and modelling only the annotation instance in the example above yields:

```
{
  "id": "http://example.org/annotation1",
  "type": "Annotation",
  "dc:title": " Discovery Services ",
  "body": {
    "source"    : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
    "selector"  : {
        "type"    : "FragmentSelector",
        "value"   : "page=7&xywh=323,502,241,45"
    }
  },
  "target": {
    "source"    : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976",
    "selector"  : {
        "type"    : "FragmentSelector",
        "value"   : "page=4"
    }
  }
}
```

Here, it should be noted here that properties absent from the Web Annotation Model can be substituted with properties from existing vocabularies, in this case "dc:title" from the DublinCore vocabulary.

However, the Web Annotation Model focuses solely on the description of annotations, while the mosaic model sketched above also includes the notions of Mosaic and Resource.

Two other constructs of the Web Annotation Model can be leveraged to express these notions: the Mosaic class can be seen as a subclass of the AnnotationCollection class, while the Resource class can be seen as an External Web Resource. Using the External WebResource construct, however, introduces 'fake' annotation instances that encapsulate resources, and implies the use of the 'purpose' property to discriminate these first-order resources from annotations (cf first item `http://example.org/resource1`).

```
{"id"   : "http://example.org/mosaic1",
 "type" : "AnnotationCollection",
 "title": "Sample Mosaic",
 "total": 2,
 "first": {
    "id"         : "http://example.org/mosaic1/page1",
    "type"       : "AnnotationPage",
    "startIndex": 0,
    "items": [
       {"id"    : "http://example.org/resource1",
        "type"  : "Annotation",
```

```
       "target": {
         "source"  : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
         "format"  : "application/pdf",
         "purpose" : "mainResource",
        }
       },
       {"id"    : "http://example.org/annotation1",
        "type"  : "Annotation",
        "body"  : {
          "source"  : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
          "selector": {
            "type"      : "FragmentSelector",
            "value"     : "page=7&xywh=323,502,241,45"
          }
        },
        "target": {
          "source"  : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976",
          "selector": {
            "type"      : "FragmentSelector",
            "value"     : "page=4"
          }
        }
      }
    ]
}
```

It must also be noted that the AnnotationCollection class enforces a paged structure, which introduces an unnecessary structural overhead.

As a conclusion, re-using the Web Annotation Model for the modelling of Mosaic and Resource classes introduces too much unnecessary overhead and impairs the readability of the model. It is therefore better to limit the reuse of the Web Annotation Model strictly to annotations, and keep a custom model for Mosaic and Resource, yielding:

```
{"id"    : "http://example.org/mosaic1",
 "type" : "Mosaic",
 "title": "Sample Mosaic",
 "resources": [
       {"id"    : "http://example.org/resource1",
        "type"  : "Resource",
        "source": "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
        "format": "application/pdf",
       }
    ],
 "annotations": [
       {"id"    : "http://example.org/annotation1",
        "type"  : "Annotation",
        "body"  : {
          "source"  : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
          "selector": {
            "type"      : "FragmentSelector",
            "value"     : "page=7&xywh=323,502,241,45"
          }
        },
        "target": {
```

```
          "source"  : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976",
          "selector": {
            "type"     : "FragmentSelector",
            "value"    : "page=4"
          }
        }
      }
    ]
}
```

### 3.1.1.2.    Decoupling mosaics from their visualisation

While one of the main goals of this project is to build efficient and user friendly visualisation tools to explore and consume mosaics, another major objective is to use mosaics to model knowledge as graphs of linked fragments. In that perspective, mosaics become structured and queriable datasets of their own, modelling a certain context, regardless of the mean of visualisation, if any.

Therefore, the mosaic model discussed here shall be decoupled from any visualisation configuration, leaving it to display tools to come up with the best visualisation solution.

However, some media types (e.g. CSV) can hardly have useful generic viewers, and require display instructions specific to the viewing context. It is therefore desired to have a placeholder in the Resource class model that will hold display configuration. Such a display configuration is required for some media types (e.g. what columns of a CSV to display, along which axis), but can be useful for others too (e.g. define a base map or a projection for geospatial data).

In the mosaic model described above, the optional `vizConfig` element provides a placeholder for such display hints. The content of such an element is to be defined in the implementation phase.

### 3.1.1.3.    Multilingualism

Multilingualism is an issue at least at two levels: inline annotation text, and multilingual annotated resources, resulting in multiple URIs with different respective fragments expressions.

If using the "Web Annotation Model", both these levels can be tackled with instances of the 'Choice' class, allowing to provide a choice of several entities for the body or the target of an annotation. This mechanism is not dedicated to multilingualism; it is rather up to the viewing application to somehow choose the best entity to display among all proposed.

Below is an example of both an annotated resource (the target) and an annotation text (the body) with multilingual values.

```
{
  "id"    : "http://example.org/annotation1",
```

```json
  "type" : "Annotation",
  "body" : {
    "type"  : "Choice",
    "items" : [
      { "value"   : "the View Service description",
        "language": "en"},
      { "value"   : "une description du Service de Consultation",
        "language": "fr"
      },
    ]
  },
  "target": {
    "type"  : "Choice",
    "items" : [
      {
        "source"  : "http://eur-lex.europa.eu/.../EN/.../?uri=CELEX:32007L0002",
        "selector": {
          "type"     : "FragmentSelector",
          "value"    : "page=7&xywh=323,502,241,45"
        },
        "language": "en"
      },
      {
        "source"  : "http://eur-lex.europa.eu/.../FR/.../?uri=CELEX:32007L0002",
        "selector": {
          "type"     : "FragmentSelector",
          "value"    : "page=7&xywh=323,462,241,60"
        },
        "language": "fr"
      },
    ]
  }
}
```

### 3.1.2. Linked Data & Semantic Web

#### 3.1.2.1. The need for triples

While the choice of a storage solution is an implementation choice that should not be set in stone, it is necessary to consider the linked data nature of the modelled problem to better frame architectural choices.

One of the purpose of this project is to express links between fragments of data, annotate these links and fragments using domain-specific vocabularies, and group these links into contextualized mosaics. The resulting data is therefore a graph that should be browsable, and queriable across mosaic instances. The implementation choices for the storage solution shall take into account this graph nature of the data, and provide support for the functionalities described in the next paragraphs.

#### 3.1.2.2. Usage of vocabularies

Apart from the vocabularies needed to describe the model itself, it is desired to be able to qualify annotations with terms from arbitrary vocabularies. In the perspective of a knowledge representation paradigm that focuses on linked media fragments and their

forming a meaningful graph, it is essential to be able to capture the semantics of those links, using any domain-specific vocabulary.

Such semantic tagging should be possible either to tag a resource with a IRI (e.g. a (fragment of) a document is classified by that IRI ), or to qualify the nature of an annotation (e.g. the two fragments of an annotation are bound by a relation that is classified by that IRI).
An example of a domain vocabulary used to classify resources and annotations will be described in the prototype section.

▶  *Classify a resource*

In the example below, the `purpose` statement is used together with a concept URI to classify a fragment of a document.

```
{
  "id"     : "http://example.org/annotation1",
  "type"   : "Annotation",
  "body"   : {
    "purpose" : "classifying",
    "source"  : "http://example.org/vocabulary/NetworkService"
  }
  ,
  "target": {
    "source"   : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976",
    "selector": {
      "type"       : "FragmentSelector",
      "value"      : "page=4"
    }
  }
}
```

▶  *Classify an annotation relation*

Below is an extension of the previous example, where a second annotation body is added. It contains a link to a second document fragment. As per the Web Annotation Model, the `type` attribute of the body or target is not meant to hold generic semantic typing, so a dedicated `linkType` property is used to qualify that second annotation body with a term from an external taxonomy.

```
{
  "id"    : "http://example.org/annotation1",
  "type" : "Annotation",
  "body" : [{
    "linkType" : "http://example.org/vocabulary/hasImplementationDirective",
    "source"    : "http://eur-lex.europa.eu/.../?uri=CELEX:32007L0002",
    "selector" : {
      "type"        : "FragmentSelector",
      "value"       : "page=7&xywh=323,502,241,45"
    }
  },
  {
    "purpose"  : "classifying",
```

```
    "source"   : "http://example.org/vocabulary/NetworkService"
  }]
  ,
  "target": {
    "source"   : "http://eur-lex.europa.eu/.../?uri=CELEX:32009R0976",
    "selector": {
      "type"     : "FragmentSelector",
      "value"    : "page=4"
    }
  }
}
```

### 3.1.2.3.    Querying

Storing mosaics in a triple store (as opposed to e.g. a relational database) allows flexible queries through the graph of resources, potentially leveraging vocabularies used in tagging the annotations. Below are some representative SPARQL queries that can be done on the Mosaics Service.
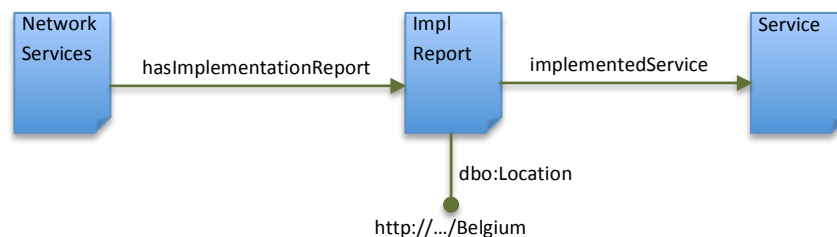
▶ *Find all fragments annotated with a certain relation*

```
SELECT ?target ?body
WHERE {
    ?annot oa:hasTarget ?target.
    ?annot oa:hasBody   ?body.
    ?body  moz:linkType <http://example.org/taxonomy/hasImplementationDirective>.
}
```

▶ *Find all services mentioned in implementation reports of Network services from Belgium*



```
SELECT ?serviceURL
WHERE {
    ?networkServices   oa:source   <http://.../networkServicesDirective.pdf>
    ?hasImplAnnotation oa:hasTarget ?networkServices
    ?hasImplAnnotation moz:linkType <http://.../hasImplementationReport>
    ?hasImplAnnotation oa:hasBody   ?implReport
    ?implReport        dbo:Location <http://.../Belgium>
    ?serviceAnnotation oa:hasTarget ?implReport
    ?serviceAnnotation moz:linkType <http://.../implementedService>
    ?serviceAnnotation oa:hasBody   ?serviceResource
    ?serviceResource   oa:source   ?serviceURL
}
```

### 3.1.3. Integrate external catalogues

To ease the insertion of existing documents and resources in a mosaic, the authoring tool should leverage data catalogues standards, by offering a mean to register instances of data catalogues and browse them for relevant resources to annotate.

Below is a non-exhaustive list of such relevant catalogue standards that have been implemented over the past years, in particular by public entities. Enabling the integration of such catalogues would leverage many years of standardization efforts.

Implementation details of the integration of these standards is left to the prototype scope.

#### 3.1.3.1.    DCAT and DCAT-AP

DCAT is not a catalogue API standard but a catalogue data vocabulary. It is used in particular, through its application profile DCAT-AP, for data catalogues in the European public sector.

#### 3.1.3.2.    CKAN

As a de facto standard for open data publication, deployed by many governments and administrative entities (including EU bodies) to publish their data, the CKAN platform is an access point to a large set of existing documents that are subject to be part of mosaics.

#### 3.1.3.3.    CSW and INSPIRE

The CSW standard is used to publish geospatial metadata. Being part of the European INSPIRE directive requirements, it is widely deployed and used to publish geospatial datasets.

## 3.2. Implementation

### 3.2.1.    Object model representation

The object model described in section 3.1.1 must be encoded in a representation format suited for the interaction with a browser-based application.

A natural encoding for such a messaging channel is JSON, allowing native parsing in a browser-based environment. On the other hand, as is reflected by the object model discussed above, the linked nature of the data is best modelled using triples. As a result, a JSON-LD based representation format is recommended.

The Web Annotation Model used as the basis for the mosaic model also uses JSON-LD as its preferred representation format. Its encoding conventions will be re-used for mosaics representations.
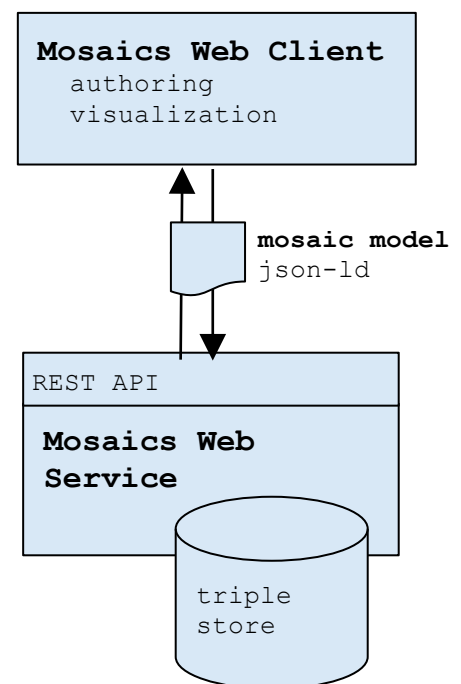


**Figure 13. Architecture**

### 3.2.2. Vocabularies

As described in section 3.1.2.2 (Usage of Vocabularies), it is necessary to define a domain vocabulary that will be used to classify and query annotations.

In the scope of this demonstrator, such a vocabulary should contain at least the following terms :

▶  *hasImplementationRule*

Qualifies an annotation between an element of the Inspire Directive and the relevant Implementation Rule document.
IRI : http://inspire.ec.europa.eu/mosaics/hasImplementationRule

▶  *hasTechnicalGuideline*

Qualifies an annotation between an element of the Inspire Directive and the relevant Technical Guidelines documents.
IRI : http://inspire.ec.europa.eu/mosaics/hasTechnicalGuideline

▶  *hasNationalTransposition*

Qualifies an annotation between an element of the Inspire Directive and its implementation in the legal text of a member state, at the national level.
IRI : http://inspire.ec.europa.eu/mosaics/hasNationalTransposition

▶  *hasSubNationalTransposition*

Qualifies an annotation between an element of the Inspire Directive and its implementation in the legal text of a subdivision of a member state.
IRI : http://inspire.ec.europa.eu/mosaics/hasSubNationalTransposition

▶  *hasImplementationReport*

Qualifies an annotation between an element of the Inspire Directive and a (fragment of a) report document describing its implementation by a member state.
IRI : http://inspire.ec.europa.eu/mosaics/hasImplementationReport

▶  *implementingResource*

Qualifies an annotation that links to a resource that is part of a member state implementation of the Inspire directive. E.g. an OGC service.
IRI : http://inspire.ec.europa.eu/mosaics/implementingResource

The controlled vocabulary used in the demonstrator can be found at
http://jrc.highlatitud.es/vocabularies/InspireLinkTypes.jsonld

### 3.2.3. Server implementation

The Mosaics Web Service is implemented in Java, using an open-source stack made of several components, in particular:
- Grizzly: a generic web service framework.
- Jersey: a reference implementation of the JAX-RS specification, used to describe and expose REST endpoints.

- Sesame: an API and storage implementation for a triple store.
- json-ld-java: a serialization library to map triples to/from JSON representations.

### 3.2.4. Web application

The Mosaics Web Client application is a HTML5 application based on Angular 1 javascript framework.

The codebase is modularized into:
- A Media fragments viewing library (`fragmentsviz`), gathering the functionalities needed to parse fragments and display them in appropriate viewers for each possible mime type.
- A Mosaics library (`ng-mosaics`), used to interact with the Mosaics Web Service REST endpoints.
- The actual Mosaics Web Client (`mosaics-jrc-app`), that orchestrates above-mentioned components to offer a comprehensive authoring/visualization user experience.

## 3.3. Technical limitations

### 3.3.1. Browser security constraints

As per the security rules known as the Same-Origin policy, browsers do not allow a web application such as the Mosaics Web Client to import and manipulate resources from other domains. Short of being able to configure access control headers in all the servers that host the remote resources, the only solution is to use a proxy server that will fake the same origin of the data.

While being a common practice used in many mainstream applications, the use of such a proxy is a performance bottleneck, forcing all the data to go through a single server. From a financial point of view, it is also a bigger cost, as hosting service providers will bill resources such as bandwidth.

A better option would be a native plugin for the browser that would disable the Same-Origin policy for a trusted application such as the Mosaics Web Client.

### 3.3.2. Storage performance

The Mosaics Service is implemented by doing on-the-fly retrieval and encoding of triples into JSON-LD mosaics. A more efficient solution performance-wise is to store JSONLD representations along the triple store and return those when querying mosaics. Cached representations would be re-computed on mosaic updates, and can be indexed using a full-text search engine for fast search access. The triple store would maintain its SPARQL querying mechanism for complex search across the graph of annotations.

## 3.4. Installation instructions & source code

Instructions to build and run the complete server can be found on the homepage of the Git project, at https://github.com/hilats/mosaics-jrc.

The demonstrator built from that source code can be found at http://jrc.highlatitud.es