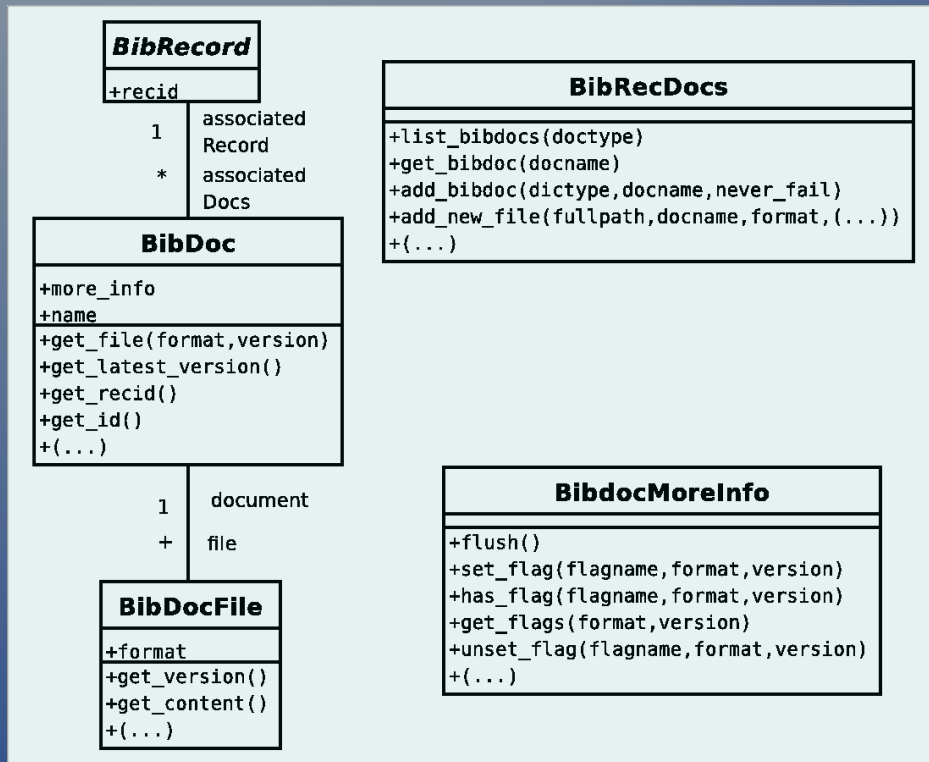


BibDocfile extensions

Piotr Praczyk piotr.praczyk@cern.ch

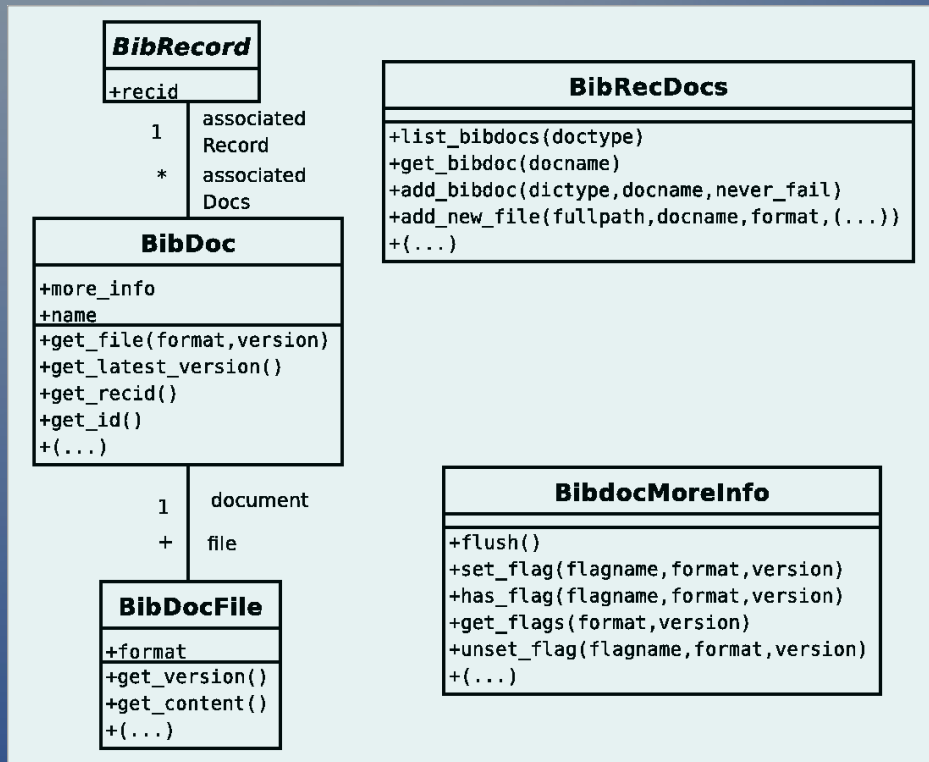
Somewhere between Meyrin and Saint Genis-Pouilly, 14.11.2012

Non-bibliographical data in Invenio



- Documents represented as BibDoc instances
- Document supports versions and different formats
- Internal data stored in a BibdocMoreInfo instance

Non-bibliographical data in Invenio

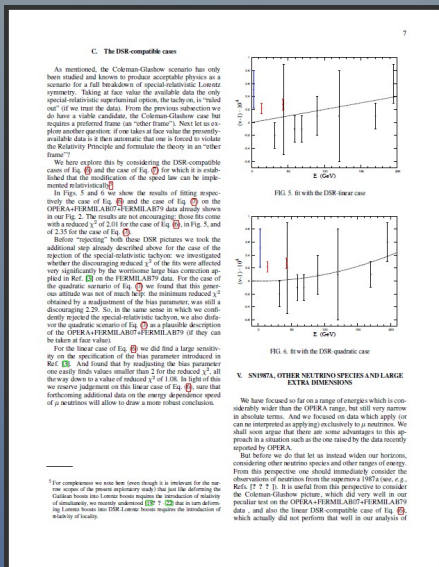


- Internal meta-data stored in a MoreInfo instance
- Link between a MARC record and the document
 - Every document must belong to exactly one record

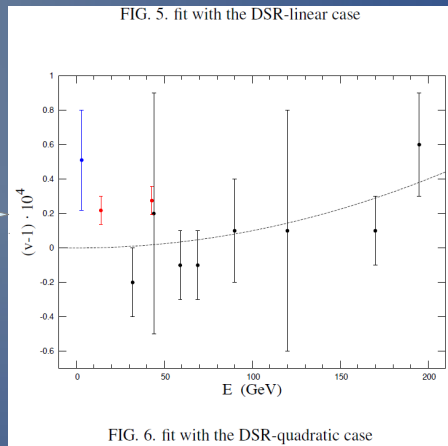
Initial assumptions

- We need to make documents independent of the records
- We need a more elastic storage for document-dependent meta-data

Example: Figures



Extracted from



Extracted from

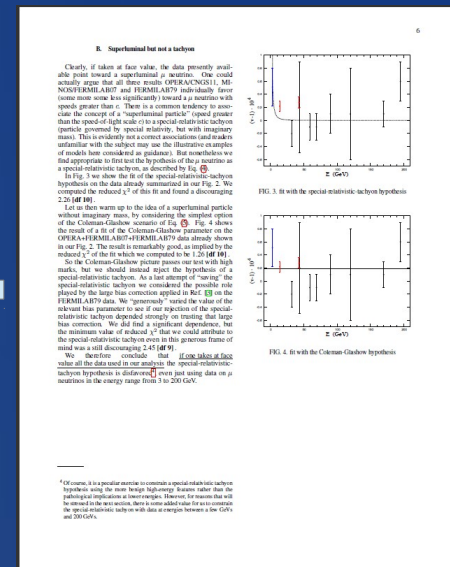


Figure 1

Describes the same data

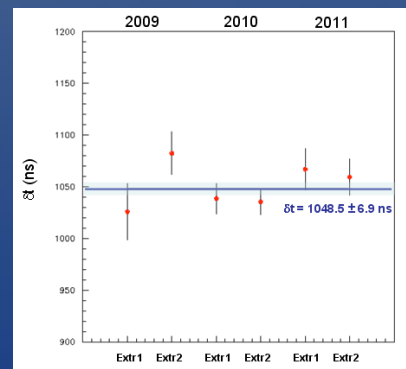


Figure 2 (extracted from different publication)

BibRecord

- * associated Records
- * associated Objects

BibDoc

+more_info: MoreInfo
+type
+get_version(version_number)
+list_versions()
+get_latest_version()
+get_id()

BibDocFile

+get_bibdocid()
+get_format()
+get_subformat()
+get_superformat()
+get_size()
+get_version()
+get_recid()
+get_content()

BibRelation

+more_info: MoreInfo
+type

MoreInfo

+serialize()
+set(namespace,key,value)
+get(namespace,key)

BibRecDocs

*The same as
BibRecDocs in
the old model*

BibFulltextDoc

+get_text()
+extract_text()

BibFigure

+get_caption()
+get_fulltext()

BibDataObject

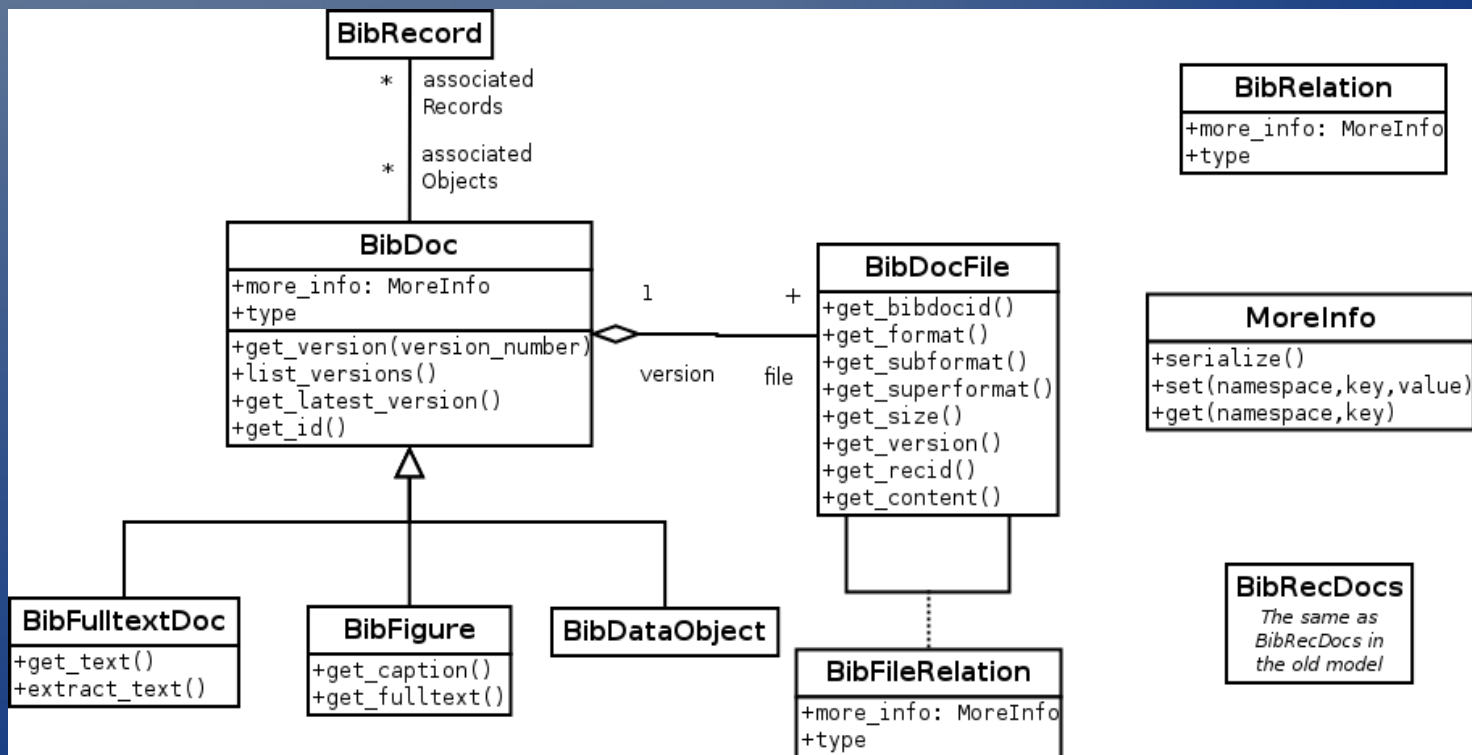
BibFileRelation

+more_info: MoreInfo
+type

1 +
version file

Data model and figures storage

Data model	Figures
BibDoc	Figure
BibDocFile	Particular encoding of a figure
BibRelation	Relation between figure and original document
BibRelation	Relation between two figures



Making documents more standalone

- Docnames cannot be identified by their names outside of the scope of a record.
 - `get_name()` ----> `get_storagename()`
- ID numbers made visible to the external world (in the FFT fields and so on)

BibRelation – link between entities

- Allows to describe dependencies and connections between different entities of the data model
- Allows specifying an arbitrary type of the relation (for example „is extracted from“, „is the same as“ etc...)

MoreInfo: custom meta-data container

Namespace → key → value

- Can be attached to any entity (BibObject, BibVersion, BibFile, BibRelation)
- Persistently stores a generic dictionaries (every module has their own identified by the namespace)
- Stored with smaller granularity

BibUpload extensions

- BibUpload = local replacement of the database transactions (assures the data consistency)
- Currently documents uploaded to Invenio using artificial FFT tag.
- Extended syntax of FFT, new tags for the new types of objects

Uploading data in new format

- New artificial MARC XML fields:
 - BRT (Uploading and modifying relations between documents)
 - MIT (Uploading MoreInfo fields)
- BDR (Attaching existing objects to records)

Uploading MoreInfo

- Externally (MIT field) or internally (from within FFT/BRT)
- Values encoded in serialised Python objects

```
{  
  "namespace": {  
    "key": "value",  
    "key2": "value2"  
  }  
}
```
- Semantics decoupled from BibUpload modes (insert/replace/correct/...)

Temporary identifiers

- Internal Invenio identifiers are assigned during the execution of BibUpload.
- We need to be able to upload relations between BibDocs whose IDs might be not know yet.

SOLUTION:

- Temporary identifier = identifier unique within the input MARC XML file

Example of TMP ID usage:

```
<collection xmlns="http://www.loc.gov/MARC21/slim">
  <record>
    <datafield tag="100" ind1=" " ind2=" ">
      <subfield code="a">This is a record of the publication</subfield>
    </datafield>
    <datafield tag="FFT" ind1=" " ind2=" ">
      <subfield code="a">http://somedomain.com/document.pdf</subfield>
      <subfield code="t">Main</subfield>
      <subfield code="n">docname</subfield>
      <subfield code="i">TMP:id_identifier1</subfield>
      <subfield code="v">TMP:ver_identifier1</subfield>
    </datafield>
  </record>
```

```
<collection xmlns="http://www.loc.gov/MARC21/slim">
```

```
<record>
  <datafield tag="100" ind1=" " ind2=" ">
    <subfield code="a">This is a record of the publication</subfield>
  </datafield>
  <datafield tag="FFT" ind1=" " ind2=" ">
    <subfield code="a">http://somedomain.com/document.pdf</subfield>
    <subfield code="t">Main</subfield>
    <subfield code="n">docname</subfield>
    <subfield code="i">TMP:id_identifier1</subfield>
    <subfield code="v">TMP:ver_identifier1</subfield>
  </datafield>
</record>
```

```
<record>
  <datafield tag="100" ind1=" " ind2=" ">
    <subfield code="a">This is a record of a dataset extracted from the publication</subfield>
  </datafield>

  <datafield tag="FFT" ind1=" " ind2=" ">
    <subfield code="a">http://sample.com/dataset.data</subfield>
    <subfield code="t">Main</subfield>
    <subfield code="n">docname2</subfield>
    <subfield code="i">TMP:id_identifier2</subfield>
    <subfield code="v">TMP:ver_identifier2</subfield>
  </datafield>

  <datafield tag="BRT" ind1=" " ind2=" ">
    <subfield code="i">TMP:id_identifier1</subfield>
    <subfield code="v">TMP:ver_identifier1</subfield>
    <subfield code="j">TMP:id_identifier2</subfield>
    <subfield code="w">TMP:ver_identifier2</subfield>

    <subfield code="t">is_extracted_from</subfield>
  </datafield>
</record>
```


Questions?