# A Storage Model for Supporting Figures and
# Other Artefacts in Scientic Libraries: the Case Study of Invenio

Piotr Praczyk *(1,2)*, Javier Nogueras-Iso *(2)*, Samuele Kaplun *(1)*, Tibor Simko *(1)*

*(1)* CERN, Geneva, Switzerland
*(2)* Universidad de Zaragoza, Zaragoza, Spain

Berlin, 29.09.2011

# Outline

- The Invenio and Inspire projects

- Old data model

- New use cases of storin data

  - Figures

  - Data preservation

- New data model

- Uploading data into the repository

# Invenio



- Created to be a basis for CERN Document Server

- Meta-data represented in MARC

# Application layer

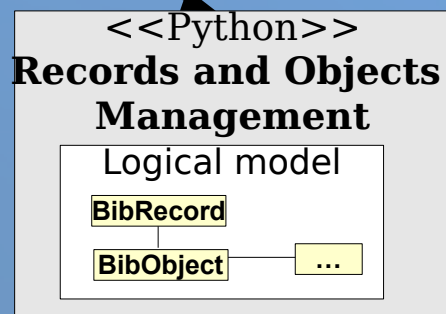**Main** file(s):

| | 📄 0101431 |
|---|---|
| version 1 | 0101431.pdf [782.53 KB] *28 Jun 2011, 14:21* |
| | 0101431.ps.gz [1.91 KB] *28 Jun 2011, 14:21* |

**Additional** file(s):

| | 📄 0101431.fig1 |
|---|---|
| version 1 | 0101431.fig1.ps.gz [2.31 KB] *28 Jun 2011, 14:21* |

| | 📄 0101431.fig2 |
|---|---|
| version 1 | 0101431.fig2.ps.gz [231.35 KB] *28 Jun 2011, 14:21* |

<<Python + JavaScript>>
**Web interface**

Query | Presentation

<<Python + JavaScript>>
**Web interface (administration)**

Users | Manual figures extraction

Records

<<Python>>
**Command Line Interface**

# Middleware layer

<<Python + C>>
**Search Engine**

<<Python>>
**Records and Objects Management**

Logical model

**BibRecord**

**BibObject** — ...

<<Python>>
**BibSched**

<<Python>>
**BibUpload**

<<Python>>
**BibIndex**

<<Python>>
**...**

# Storage layer

<<MySQL>>
**Search indexes**

<<MySQL>>
**Metadata (MARC21)**

<<Andrew File System>>
**Content of objects**

# SPIRES



- Database of preprints started in 70s
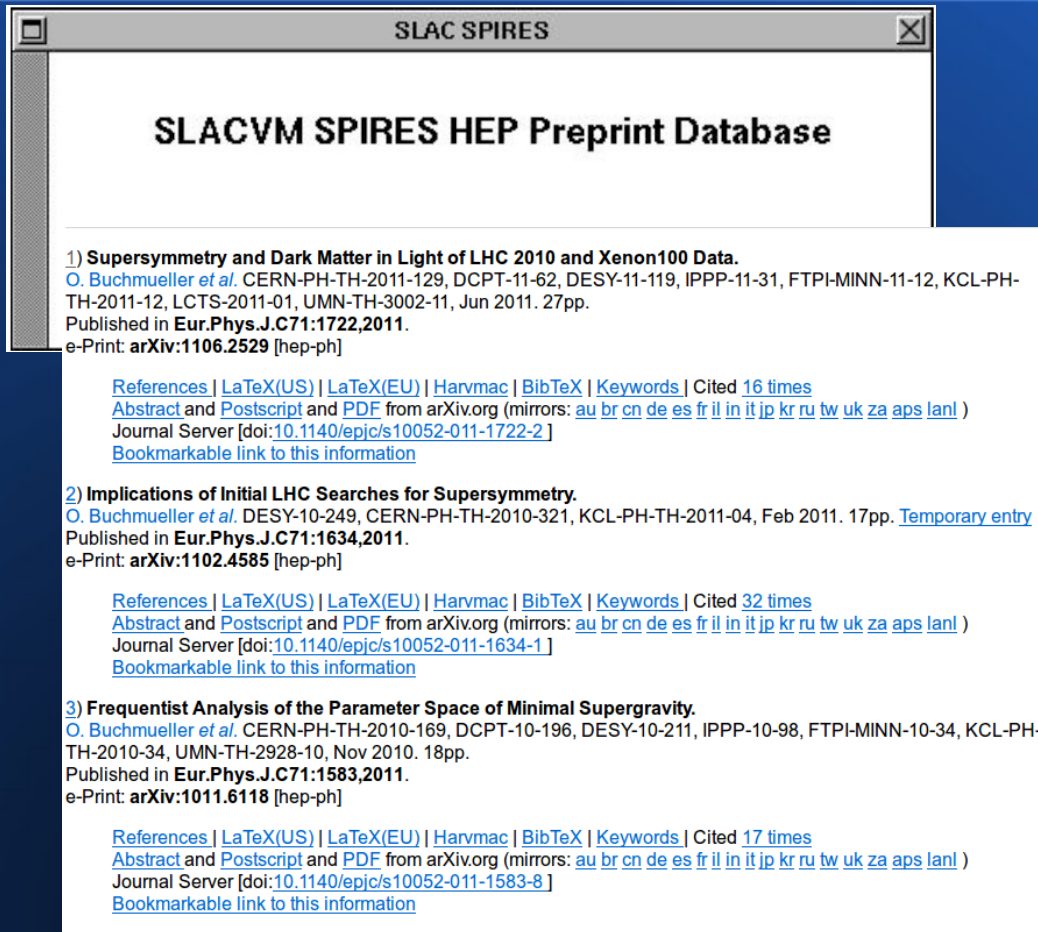
- In 90s the first WEB page in USA

- Very difficult to maintain, extremely slow

# Invenio + SPIRES = INSPIRE



- Large community of users

- Multiple sources of data (SPIRES, arXiv, diect submissions, publishers)

- Nearly 1000000 records

# Invenio/INSPIRE

- Invenio – digital library software developed at CERN to manage the repository of documents created in the institution

- SPIRES – The digital library of preprints created at SLAC.

- **In**venio + **SPIRE**S = INSPIRE

# 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

**BibRecord**
+recid

1   associated Record
*   associated Docs

**BibDoc**
+more_info
+name
+get_file(format,version)
+get_latest_version()
+get_recid()
+get_id()
+(...)

1   document
+   file

**BibDocFile**
+format
+get_version()
+get_content()
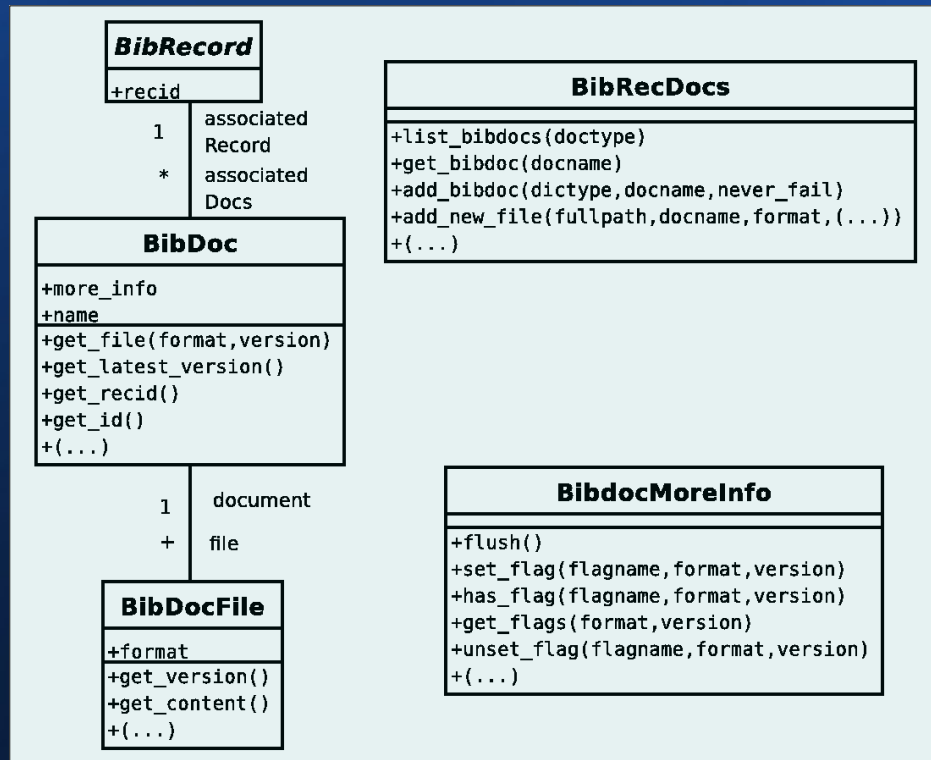+(...)

**BibRecDocs**
+list_bibdocs(doctype)
+get_bibdoc(docname)
+add_bibdoc(dictype,docname,never_fail)
+add_new_file(fullpath,docname,format,(...))
+(...)

**BibdocMoreInfo**
+flush()
+set_flag(flagname,format,version)
+has_flag(flagname,format,version)
+get_flags(format,version)
+unset_flag(flagname,format,version)
+(...)

- Internal meta-data stored in a MoreInfo instance

- Link between a MARC record and the document

  – Every document must belong to exactly one record

# Figures from scientific publications



Automatic extraction

Manual extraction

FIG. 5. fit with the DSR-linear case

FIG. 6. fit with the DSR-quadratic case

extraction

- Figures should exist independently from articles

- Some type of meta data should not be presented to users directly

# Figures from scientific publications



Extracted from

Extracted from

FIG. 5. fit with the DSR-linear case

FIG. 6. fit with the DSR-quadratic case

Figure 1

Describes the same data

Figure 2 (extracted from different publication)

# Data-model requirements

• Different types of relations between figures
(illustrates the same data, is subfigure of...)

• Relation of being extracted from a document

• Meta-data of figures, different versions of figure, relations between figures, links between figure and document it is extracted from

• Storage of more complicated data-types

# Examples of meta-data associated with different entities

| Figure (the most general meta-data) | Figure version (Appearance related meta-data) | Figure storage format | Relation between figure and document | Relation between two figures |
|---|---|---|---|---|
| Type of figure | Semantics of a figure | Access permissions | Position of a figure within the original document | Type of relation |
| Quantities presented on axis (in the case of plots) | Data extracted from figure | | Caption of a figure within a particular document | <additional type-dependent fields> |
| Units and scales of axis (in the case of plots) | | | References from within text | |
| ... | | | Figure identifier within a document | |

# HEP data preservation

- Storing raw data

- Storing intermediate analysis

- Storing additional documentation

- Assigning Digital Object Identifiers

# New architecture for storing non-bibliographical objects

# BibObject

- Abstract representation of a document (not

- Document-type specific functionalities are implemented by subclasses (defined by modules of Invenio and loaded dynamically)

- Identified by a globally-unique identifier and by a name unique in the scope of a bibliographical record

# BibVersion, BibFile

- Represent increasing specialisation of a document

- BibVersion represents a particular revision of an object (corresponding for example to correction of mistakes)

- BibFile describes a particular encoding of a version of an object (encapsulates the real file, remembers the format of a file)

# 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)

# Data model and figures storage

| Data model | Figures |
| --- | --- |
| BibObject | Figure |
| BibObjectVersion | Figure version |
| BibFile | Particular encoding of a figure |
| BibRelation | Relation between figure and original document |
| BibRelation | Relation between two figures |

# Comparison between data models

|  | Old model | New model |
|---|---|---|
| Attaching document to a bibliographic record | YES | YES |
| Attaching the same document to many records | NO | YES |
| Storing custom data keys in MoreInfo dictionaries | NO | YES |
| Creating documents not attached to any records | NO | YES |

# Uploading data to Invenio

```xml
<?xml version="1.0" encoding="UTF-8"?>
<collection xmlns="http://www.loc.gov/MARC21/slim">
  <record>
    <controlfield tag="001">929725</controlfield>
    <datafield tag="970" ind1=" " ind2=" ">
      <subfield code="a">SPIRES-9208755</subfield>
    </datafield>
    <datafield tag="100" ind1=" " ind2=" ">
      <subfield code="a">Artymowski, Michal</subfield>

    </datafield>
    <datafield tag="700" ind1=" " ind2=" ">
      <subfield code="a">Lalak, Zygmunt</subfield>
    </datafield>
    <datafield tag="856" ind1="4" ind2=" ">
      <subfield code="u">http://inspirebeta.net/record/929725/files/arX
iv:1109.5901.pdf</subfield>
    </datafield>
    <datafield tag="856" ind1="4" ind2=" ">

      <subfield code="u">http://inspirebeta.net/record/929725/files/BBH
hnkanon.png</subfield>
      <subfield code="y">00003 The left panel shows the evolution of sc
ale factors for the domination of a massless vector field with non cano
nical kinetic term and $f\propto a^{-4}$. The right panel presents the
evolution of Hubble parameters in the same model. One can see, that aft
er $w\sim t\sqrt{\rho_I}\sim 5$ the Universe becomes isotropic and ente
rs the era of exponential expansion.</subfield>
    </datafield>
    <datafield tag="856" ind1="4" ind2=" ">
      <subfield code="u">http://inspirebeta.net/record/929725/files/BBa
bnkanon.png</subfield>
```
`:--- sample.xml    Top L7    (nXML Valid)-------------------------------`

- New record is encoded in MARC XML
- BibUpload is executed adding uploading task to the BibSched queue
- BibSched uploads data to the main database

# FFT = Fulltext File Transfer

```xml
<?xml version="1.0" encoding="UTF-8"?>
<collection xmlns="http://www.loc.gov/MARC21/slim">
  <record>
    <datafield tag="FFT" ind1=" " ind2=" ">
      <subfield code="a">http://invenio-software.org/download/invenio-demo-site-files/0106015_01.jpg</subfield>
      <subfield code="r">restricted_picture</subfield>
    </datafield>
    <datafield tag="FFT" ind1=" " ind2=" ">
      <subfield code="a">http://invenio-software.org/download/invenio-demo-site-files/0106015_01.gif</subfield>
      <subfield code="f">.gif;icon</subfield>
      <subfield code="r">restricted_picture</subfield>
    </datafield>
  </record>
</collection>
```

| subfield | explanation |
|---|---|
| a | URL of the file to upload |
| t | Function of the document within the record |
| f | Format of the file |
| | ... |

- Artificial, interpreted and removed during the BibUpload phase

- One entry represents one file

- Enforces documents to be attached to a record

# Uploading data in new format

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

- Open for extension supporting METS

# Uploading MoreInfo

- Externally (MIT field) or internally (from within FFT/BRT)

- Values encoded in Json or serialised Python objects

- Semantics completely decoupled from BibUpload modes (insert/replace/correct/...)

# Conclusions & outlook

- The proposed a model is a flexible approach that facilitates the support of custom objects (figures, data files, software, …)

  - Based on it, new applications for searching and accessing digital objects can be developed

- Issues to address in the future

  - Integration of extended objects within the INVENIO platform

    - Search and display information about custom objects

  - Need of assigning Digital Object Identiers (DOI) to stored data objects

    - To store the persistent state of a data object (management of versions)

# Thank you !

http://invenio-software.org/
http://www.projecthepinspire.net/

piotr.praczyk@cern.ch