CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

# CDS-Invenio for librarians

Jean-Yves Le Meur - CERN Document Server Project Leader

September 22, 2009

**CDS-Invenio for librarians**
Record definition
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

**What is a record ?**
Internal representation of records: MARCXML
The "life" of a record

- Who
- When
- Where
- What
- With whom
- ...

MARCXML

CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
**Internal representation of records: MARCXML**
The "life" of a record

- MARC Standard exists since the 60s
- established by the Library of Congress, it is used by most library systems to describe bibliographic information of records
- very large number of fields, and freedom to "invent" new ones
-
```xml
<record>
  <datafield tag="041" ind1=" " ind2=" ">
    <subfield code="a">eng</subfield>
  </datafield>
  <datafield tag="088" ind1=" " ind2=" ">
    <subfield code="a">PRE-25553</subfield>
  </datafield>
  <datafield tag="088" ind1=" " ind2=" ">
    <subfield code="a">RL-82-024</subfield>
  </datafield>
  <datafield tag="100" ind1=" " ind2=" ">
    <subfield code="a">Ellis, J</subfield>
    <subfield code="u">University of Oxford</subfield>
  </datafield>
  <datafield tag="245" ind1=" " ind2=" ">
    <subfield code="a">Grand unification with large supersymmetry breaking</subfield>
  </datafield>
  <datafield tag="260" ind1=" " ind2=" ">
    <subfield code="c">Mar 1982</subfield>
  </datafield>
  <datafield tag="300" ind1=" " ind2=" ">
    <subfield code="a">18 p</subfield>
  </datafield>
```
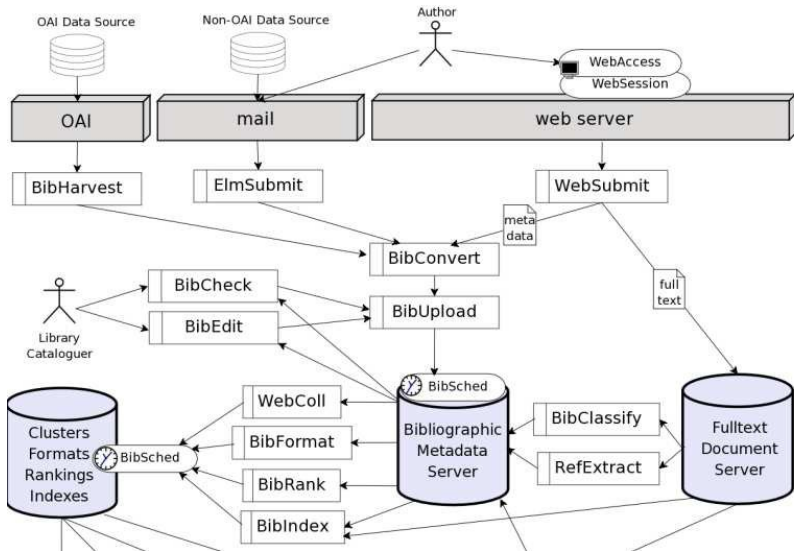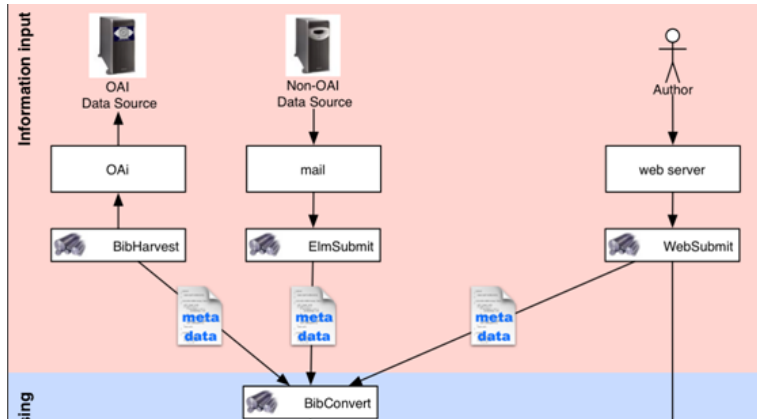
CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
**Internal representation of records: MARCXML**
The "life" of a record

| Information | References | Discussion | Usage statistics | Fulltext |

## Internal Note

| | |
|---|---|
| Report number | CERN-IT-Note-2007-025 |
| Title | **Authentication/authorization issues and fulltext document migration for the CERN Document Server** |
| Author(s) | Kaplun, S |
| Corporate author(s) | CERN. Geneva. IT Department |
| Collaboration | UDS |
| Imprint | 16 Oct 2007 - 91 p |
| Subject category | Computing and Computers |
| Free keywords | CERN ; CDS Invenio   Authentication ; Authorization ; Single Sign-On ; FireRole ; IntBitSet ; S2D ; Role Based Access Control |
| Abstract | This thesis describes a master degree project, ending studies at Università degli Studi di Milano Bicocca of Computer Science, Milano This work has been realized at the European Organization of Nuclear Research (CERN), in Geneva. The aim of the project was to enhance CDS Invenio, a digital library software developed by CERN, in the authentication/authorization area, to develop an automatic migration tool for moving documents from the legacy architecture and to develop an extension to Python in C for solving indexing time issues. |

Email contact: Samuele.Kaplun@cern.ch

Record created 2007-10-18, last modified 2007-10-20                    Similar records
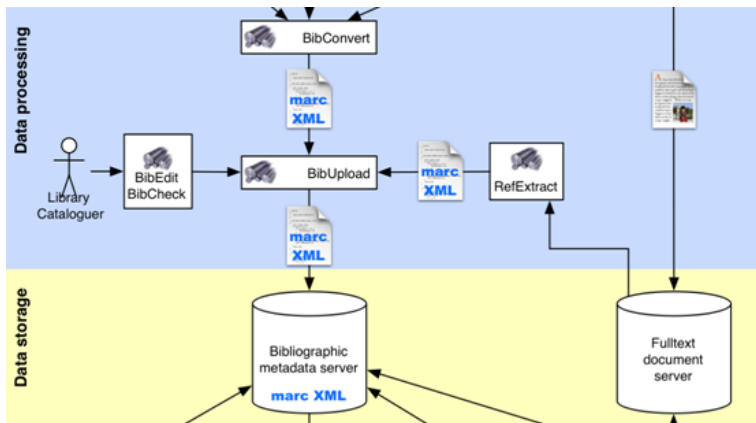
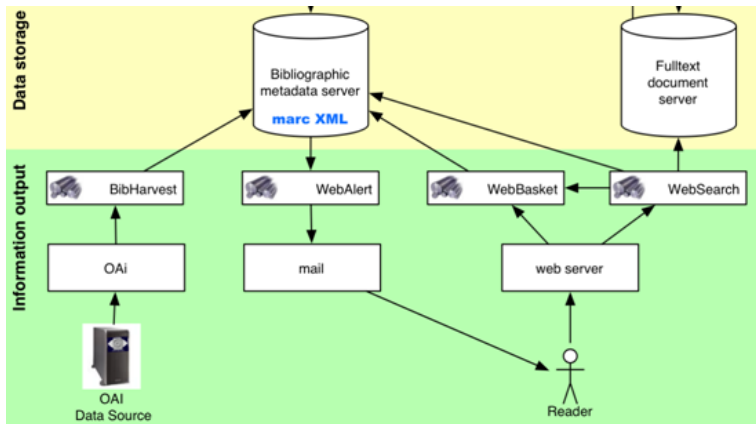CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
Internal representation of records: MARCXML
**The "life" of a record**

CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
Internal representation of records: MARCXML
**The "life" of a record**

# Information Input

CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
Internal representation of records: MARCXML
**The "life" of a record**

# Data Process and Storage

CDS-Invenio for librarians
**Record definition**
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

What is a record ?
Internal representation of records: MARCXML
**The "life" of a record**

# Information Output

CDS-Invenio for librarians
Record definition
**Document Submission by Authors**
Batch Acquisition of records
Editing tools for librarians
Conclusion

**Web Submission**
Submissions at CERN
Email Submission

# Web submissions (I)

- ▶ Each collection can have its own submission policy
  - ▶ Direct submission
  - ▶ Submission with monitoring
  - ▶ Submission with simple approval
  - ▶ Submission with peer review/refereeing and editorial board
- ▶ Each collection can have its own record definition
  - ▶ Metadata fields (mandatory, optional, controlled at input time)
  - ▶ Full text formats
  - ▶ Revised versions

CDS-Invenio for librarians
Record definition
**Document Submission by Authors**
Batch Acquisition of records
Editing tools for librarians
Conclusion

**Web Submission**
Submissions at CERN
Email Submission

# Web submissions (II)

- ▶ Each submission has its own process management
    - ▶ It can be configured with an HTML administration interface
    - ▶ To define submission screens
    - ▶ To define actions to be applied when the document is transferred
    - ▶ Examples:
        - ▶ When finishing the submission of a videotape by video service, the label is created (PDF) to be stick on the tape
        - ▶ When a note is submitted by a collaboration, members of the collaboration are immediatly notified by email for comments

CDS-Invenio for librarians
Record definition
**Document Submission by Authors**
Batch Acquisition of records
Editing tools for librarians
Conclusion

Web Submission
**Submissions at CERN**
Email Submission

## Submit

**Document types available for submission:**

Please select the type of document you want to submit:
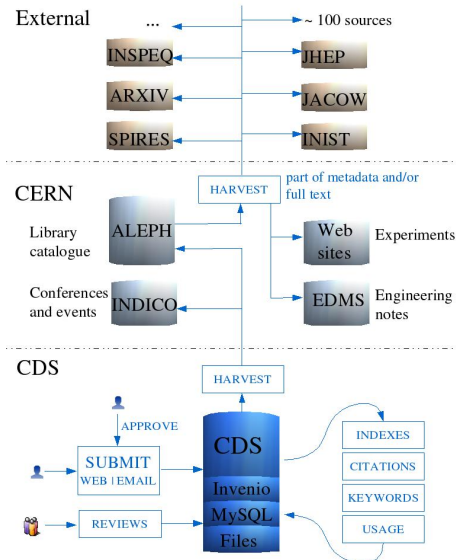
- **Preprints, Notes, Articles...**
  - CERN Preprints Automatic Numbering
  - CERN Preprints
  - CERN Thesis
  - CERN Open Documents
  - Departmental Internal Notes
  - TOTEM Notes
  - imXgam Internal Notes
  - ANTARES Plots
  - EGEE Publications and Technical Reports

- **Experiments Committees**
  - LHCC
    - LHCC Documents (Proposals etc)
    - LHCC Meeting Documents (protected)
    - LHCC Technical Design Reports
  - INTC
    - INTC Documents (Proposals etc)
    - INTC Meeting Documents (protected)
  - SPSC
    - SPSC Documents (Proposals etc)
    - SPSC Meeting Documents (protected)

CDS-Invenio for librarians
Record definition
**Document Submission by Authors**
Batch Acquisition of records
Editing tools for librarians
Conclusion

Web Submission
Submissions at CERN
**Email Submission**

# Email Submission

- ▶ Users can send their document as attachement by email
- ▶ The destination must be a dedicated Invenio email address (eg: submit@invenio.cern.ch)
- ▶ Subject and Body of the email contain the bibliographic information
- ▶ It must be well formatted to be accepted by the server
- ▶ The server returns aknowledgment of the succesful submission
- ▶ Almost not used at CERN, but it could become the most efficient system for end-users

CDS-Invenio for librarians
Record definition
Document Submission by Authors
**Batch Acquisition of records**
Editing tools for librarians
Conclusion

Acquisition from other systems
Harvesting Open Access Compatible sites
Converting Records
Harvesting and Converting Fulltexts

CDS-Invenio for librarians
Record definition
Document Submission by Authors
**Batch Acquisition of records**
Editing tools for librarians
Conclusion

Acquisition from other systems
**Harvesting Open Access Compatible sites**
Converting Records
Harvesting and Converting Fulltexts

## Open Access Harvesting

▶ Collaborative framework with Data and Service Providers

▶ Information interoperability

▶ Value added processing (metadata brokering)

▶ OAI-PMH protocol: Data Provider $< -- >$ Service Provider

▶ Invenio has a Web Interface to manage Harvesting and Dissemination via OAI-PMH protocol

CDS-Invenio for librarians
Record definition
Document Submission by Authors
**Batch Acquisition of records**
Editing tools for librarians
Conclusion

Acquisition from other systems
**Harvesting Open Access Compatible sites**
Converting Records
Harvesting and Converting Fulltexts

# BibHarvest Admin Interface

Overview of sources  [?]          [add new OAI source]
  **1 OAI sources currently present in the database**

| name | baseURL | metadataprefix | frequency | bibconvertfile | postprocess | actions |
|------|---------|----------------|-----------|----------------|-------------|---------|
| cds | http://cdsweb.cern.ch/oai2d | marcxml | daily | oaimarc2marcxml.xsl | h-c-u | edit / delete |

Harvesting status  [?]

  **Next oaiharvest task**
       - scheduled time: 2008-07-04 15:09:12
       - current status: WAITING

| name | last update |
|------|-------------|
| cds | 2008-07-04 10:08:52 |

CDS-Invenio for librarians
Record definition
Document Submission by Authors
**Batch Acquisition of records**
Editing tools for librarians
Conclusion

Acquisition from other systems
Harvesting Open Access Compatible sites
**Converting Records**
Harvesting and Converting Fulltexts

## Converting Records

▶ It can be performed using ad hoc scripts (e.g: using yaz), to deal with specific formats

▶ Lot of converters from *anything* to MARCXML already exist (from LoC site)

▶ Invenio proposes two internal ways of converting incoming records:
  ▶ XSLT: most standard tool to transform XML files
  ▶ BibConvert descriptive tagging

▶ BibConvert
  ▶ You first describe your source file with tags
  ▶ You then describe your target file with the same tags
  ▶ The converter moves the information within the new scheme

CDS-Invenio for librarians
Record definition
Document Submission by Authors
**Batch Acquisition of records**
Editing tools for librarians
Conclusion

Acquisition from other systems
Harvesting Open Access Compatible sites
Converting Records
**Harvesting and Converting Fulltexts**

► The FFT Protocol: Fulltext File Transfer
  ► Additional "FFT" tag in MARC records when submitted to BibUpload: $< datafieldtag = "FFT" ind1 = "" ind2 = "" >< subfieldcode = "a" > http : //cdsware.cern.ch/download/invenio − demo − site − files/0101431.ps.gz < /subfield >< /datafield >$

► Fulltext files can be uploaded at the same time as the metadata

► If files are scanned files (bitmap), OCR process can be started automatically
  ► to recognize text and index it
  ► to let users copy/paste text inside the files

► Files are managed with a dedicated module: BibDocFile

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
**Editing tools for librarians**
Conclusion

**Cataloguing tools overview**
Bibliographic Information Edition
Record maintenance

- ▶ goal: reproducing traditional library systems cataloguer-level functionality
- ▶ record editing interface
- ▶ record checking tools
- ▶ record maintenance tools
- ▶ record inputting workflow
- ▶ record harvesting workflow
- ▶ knowledge bases

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

Cataloguing tools overview
Bibliographic Information Edition
Record maintenance

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
Conclusion

Cataloguing tools overview
Bibliographic Information Edition
Record maintenance

## Editing Tool for cataloguers

- ▶ Direct access to any record (authorized)
- ▶ Access control on who can edit what
- ▶ Web interface with Javascript (JQuery)
- ▶ Designed for distributed cataloguing
- ▶ Record locking mecanism to avoid parallel modifications

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
**Editing tools for librarians**
Conclusion

Cataloguing tools overview
**Bibliographic Information Edition**
Record maintenance

## Record #1

**Your changes are TEMPORARY. To save this record, please click on submit.**

Record #1

Action: Cancel Record: Add Field | Delete Display Verbose MARC

| 037___ | s$a | CERN-EX-0106015 | |
|---|---|---|---|
| 100___ | s$a | Photolab | |
| 245___ | s$a | ALEPH experiment: Candidate of Higgs Boson production | |
| 246_1 | s$a | Expérience ALEPH: Candidate de la production d'un boson Higgs | |
| 260___ | $$c | 14 06 2000 | |
| 340___ | s$a | FILM | |
| 520___ | s$a | Candidate for the associated production of the Higgs boson and Z boson. Both, the Higgs and Z boson decay into 2 jets each. The green and the yellow jets belong to the Higgs boson. They represent the fragmentation of a bottom and anti-bottom quark. The red and the blue jets stem from the decay of the Z boson into a quark anti-quark pair. Left: View of the event along the beam axis. Bottom right: Zoom around the interaction point at the centre showing details of the fragmentation of the bottom and anti-bottom quarks. As expected for b quarks, in each jet the decay of a long-lived B meson is visible. Top right: "World map" showing the spatial distribution of the jets in the event. | |
| 595___ | s$a | Press | |
| 65017 | s$2 | SzGeCERN | |
| | s$a | Experiments and Tracks | |
| 6531_ | s$a | LEP | |
| 8560_ | $$s | nel.cader@cern.ch | |

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
**Editing tools for librarians**
Conclusion

Cataloguing tools overview
**Bibliographic Information Edition**
Record maintenance

# Record Editor: Record #6 (view history)

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
**Editing tools for librarians**
Conclusion

Cataloguing tools overview
Bibliographic Information Edition
**Record maintenance**

- ▶ BibCheck module: allows to run quality checking on the metadata (command line)
- ▶ BibKnowledge module: allows to maintain Authority files and apply them to set of records
- ▶ BibExport module: allow to run batch export of records in many formats, for manual corrections and re-upload
- ▶ Holding pen: this is the place where the pending treatments to perform to records are listed
- ▶ BibMerge module: allows to merge identical records, preserving the information of the two records

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
**Editing tools for librarians**
Conclusion

Cataloguing tools overview
Bibliographic Information Edition
**Record maintenance**

# Record Merger

Record1: 2
Record2: 148

Compare
Method: Revisions
Submit  Cancel

CDS-Invenio for librarians
Record definition
Document Submission by Authors
Batch Acquisition of records
Editing tools for librarians
**Conclusion**

## Conclusion

- ▶ In Invenio, a record is described in XMLMARC, and it can be associated with multiple files

- ▶ A record moves in the system through multiple modules, where it can be enriched with more information

- ▶ Submission of documents by (or in the name of) authors can be configured in multiple ways

- ▶ Records can also be imported from other systems in many different ways, including fulltexts

- ▶ Librarians have powerful tools to edit and improve existing records