# Invenio:
# a modern digital library system

**Samuele Kaplun <Samuele.Kaplun@cern.ch>**
*(on behalf of the Invenio Development Team)*

## <http://invenio-software.org/>

- CERN adopted Open Access since almost **60 years**

- **Explored and invented** the tools which were needed to accomplish this (starting with *paper dissemination* in 1954)

- **Invenio** (formerly known as CDSware) answers the need of having a large scale digital library (since 2002)

From the CERN Convention:
*Paris, 1st July 1953*

[…] **the results of its experimental and theoretical work shall be published or otherwise made generally available**. [...]

# History II

| 1993 | 1996 | 2000 | 2002 | 2006 |
|---|---|---|---|---|
| Preprint Server | WebLib | CERN Document Server | CDSware | Invenio |

- **2009** – Invenio takes part in EU projects:

  - **OpenAIRE** – to implement the Orphan Record Repository

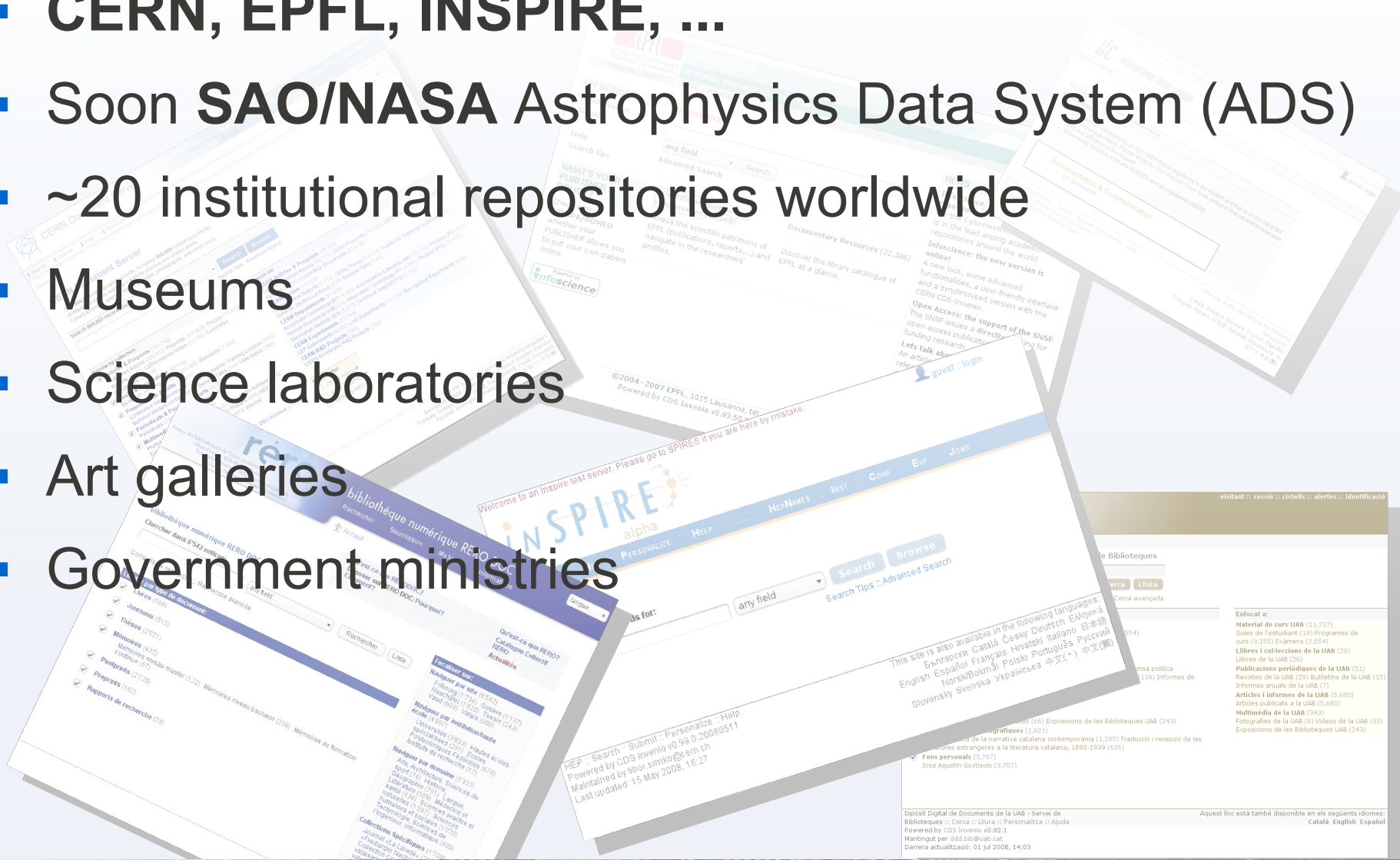  - **D4Science II** – through collaboration with INSPIRE

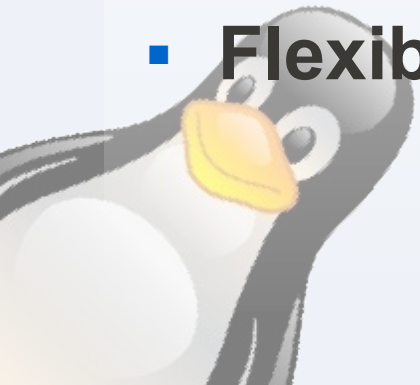- **2010** – target release date of **Invenio 1.0**

## \<http://invenio-software.org/\>

# Usage

- **CERN, EPFL, INSPIRE, ...**

- Soon **SAO/NASA** Astrophysics Data System (ADS)

- ~20 institutional repositories worldwide

- Museums

- Science laboratories

- Art galleries

- Government ministries

# Architecture

- **Open source** GNU General Public License project

- **Python** (and C and Lisp and Javascript), **MySQL** and **Apache**

- **Modular architecture**

- Implements **open standards** (MARCXML, MARC21, OAI-PMH, OpenURL...)

- **Medium to big** size record repositories (~5M)
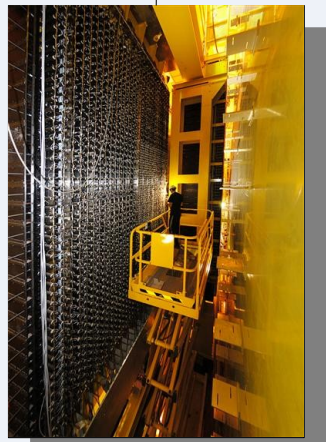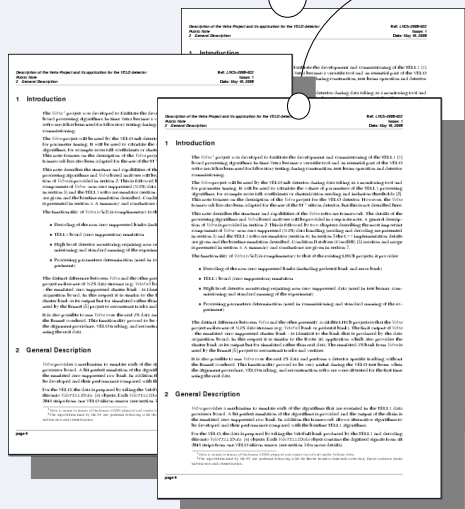
- **Flexible** in every layer

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

marcxml

Multiple revisions

Multiple formats

# Ingestion modules

- **OAI-PMH Harvesting** (DC and MARCXML out of the box)

- **Mail based** document deposition

- **Web based** flexible framework to design submission web forms and approval workflows

  - Several workflow building blocks, including fulltext watermarking, plots extractions, ...

- **SWORD** server side (development ongoing)

# Background processing modules I

- **Collections** (based on metadata)

- Pre-caching of common **formats**

- Pre-caching of common **splash-pages**

- Variety of **ranking methods**

- Extremely fast **indexing data structure** (to implement IR boolean model)

- Any metadata field can be indexed

- Special treatment for author and journal names, dates and fulltext documents

- Garbage collector

- Database dumper

- Harvester and OAI updater

- Task scheduler supporting priorities

- Automated citation extractor

- Taxonomy based keyword classifier

- Automatic author disambiguation

- Any Invenio heavy computational task implemented as a *daemon*

- **Bibliographic Task Scheduler** that takes care of executing tasks at proper times, respecting dependencies and priorities

- Plugin framework to quickly write ***tasklets*** to add additional recurrent procedures in the day-to-day repository life

- Thanks to the partnership with D4Science II, cloud and grid computating are being explored to perform heavy computational tasks

- **Words**, **pairs** of words and **phrase** indexes

- Using forward and reverse indexes

  - By using **extremely fast bit set** in the form of a C extension to Python

- **Stemming**, stop words, TeX formulas, TeX/HTML **stripping**

- Investigation on Solr and Lucene integration

- Planned implementation of generalized *derived logical fields* framework i.e. indexable metadata fields generated on the fly from any piece of information related to a record, via plug-ins

# Dissemination and publishing modules I

- Fully configurable **regular and virtual collection tree** (defined on the metadata)

- Multiple pluggable and flexible **output formats** (e.g. MODS, BibTeX, EndNote, Excel, EndNote, DC, NLM, RefWorks...)

- **Fast response** to search queries

- Search available on **any metadata field, fulltext, citations, plot captions...**

- Personalized **RSS feeds**

- **OpenURL** resolver

- **OAI-PMH** 2.0 exporting (MARC and OAIDC out of the box)

- SWORD pushing (currently to arXiv)

- Integration with **LibX** Firefox/IE toolbar

- Integration with **Cooliris**

- Integration with **Zotero**

- **OpenSearch** support

- **Journal** module (web based journal generated from records)
  - An Open Access publishing platform



Submit → Review → Publish → Announce → Submit

# Dissemination and publishing modules IV



ATLANTIS INSTITUTE OF FICTIVE SCIENCE

Search | Submit | Personalize ▾ | Help | **Administration** ▾

Home > Admin Area > BibCirculation Admin

## BibCirculation Admin

Loan | Return | Request | Borrowers | Items | Lists | Libraries | Vendors | Acquisitions | ILL | Hel

WELCOME TO CDS INVENIO BIBCIRCULATION ADMIN

- Last loans
- Overdue loans
- Items on shelf with holds
- Items on loan with holds
- Overdue loans with holds
- Ordered books

admin :: logout

- **Circulation** module (to use Invenio to manage items and borrowers)

- CERN Library catalogue integrated into Invenio

- **Library-specific features** implemented (loans/holdings/reservations/warning/renewals...)

- Ideal solution for "*library meets repository*"

- **Baskets/bookshelves** (to organize personal collections of interesting records and to share them)

- **Commenting**, **reviewing**

- Web **messages**

- Automatic **alerts** for new results

- **Tagging** (soon to be available)

# Curation modules I

- **MARC Editor**
  - Ajax record editor, with undo/redo, holding pen, knowledge base lookup, and more
  - Multi-record editor to modify in batch several records
  - Record merger to visualize and merge differences
- **Knowledge bases** (authority files, vocabularies, dynamically generated)
  - Integration with D-NET vocabularies (in the framework of the OpenAIRE collaboration)

# Curation modules II

- **Flexible configuration** based on web interfaces for:
  - Logical fields, indexes, ranking algorithms, collection structures, deposition interfaces, collection appearance and hierachy, OAI-PMH harvesting and exporting, formats, ...

Configure BibFormat

Configure BibIndex

Configure BibRank

Configure Bibknowledge

Configure OAI Harvest

Configure OAI Repository

Configure WebAccess

Configure WebComment

Configure WebJournal

Configure WebSearch

Configure WebSubmit

Run Batch Uploader

Run BibCirculation

Run Document File Manager

Run Multi-Record Editor

Run Record Editor

Run Record Merger

Statistics

- Automatic metadata **inconsistency checker** (with optional automatic correction)

- Automatic filter of incoming records against current repository to spot e.g. **duplicates**

- **Batch uploading** of fulltext files

- Web interface to **manage full-text files** of a given record

- Author disambiguation tool project

- **Plug-in based** authentication framework:

  - To exploit already existing credentials data-bases

  - **LDAP** support

  - **Shibboleth** support to implement Single Sign-On.

- Local and external **groups** (useful for authorization, exchanging messages, defining shared baskets)

- *Role-Based Access Control* (RBAC) authorization model

- Role defined:
  - Explicitly by attaching users to roles
  - Implicitly via *Firewall-like Role Definitions*
- **Any user detail** can be used to define a role (i.e. group membership, IP address, regular expression on the institution name…)

- Extensible **ranking** framework
- Support out of the box for:
  - **Word similarity** (implementing a flavour of the classical IR vectorial model)
  - advanced **citation network analysis** (co-cited with, time-decayed rankings)
  - **Record access** and **download**

# Statistics

- Extensible **event recording** framework

- **Access and download** statistics

- **Knowledge Exchange guidelines** to facilitate the exchange of usage statistics (in the framework of the OpenAIRE collaboration)

# Digital objects support I

- **Multiple documents** can be attached to a record

- **Multi-format** and **multi-revision** support

- **Icon** creation

- PDF **Stamping**

- **PDF/A** conversion

- Text extraction (including **OCR**)

- **Scanned PDF to PDF** with OCRed text

- HTTP streaming supporting **range requests** (e.g. to stream multimedia material, or linearized PDFs)

- Automated **checksum** verification

- **Batch uploading** of documents

- Multiple **protection and authorization** levels

- Automatic and plug-in based **meta-data extraction** (e.g. XMP, EXIV, …)

- Plugin-based **format conversion framework** (development ongoing)

- Invenio interface exists in 26 languages

# Development support

- *Hacking guides* (i.e. recipes on how to extend and tune Invenio even further)

- Unit, regression and web **test suites**

- **Unhandled exceptions** sent via email to admin

  - **SMS sending** for emergencies (e.g. when bibliographic task manager stops after an error)

- Inherent **code quality** checking

# Conclusions

- **Extremely flexible** integrated digital library and digital repository software

- Suitable for **medium to large** repositories (up to 10M records)

- **Broad range of features**

- **Highly configurable** to cover any large institution needs

# Links

- Invenio home page: 

  - <http://invenio-software.org/>

- Invenio demo site (running v0.99.1):

  - <http://invenio-demo.cern.ch/>

  - <http://invenio-software.org/wiki/General/Demo>

- OpenAIRE:

  - <http://www.openaire.eu/>

- INSPIRE:

  - <http://inspirebeta.net/>