Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# Integrating Scholarly Publications and Research Data - Preparing for Open Science, a Case Study from High-Energy Physics with Special Emphasis on (Meta)data Models

Piotr Praczyk[1,2]     Javier Nogueras-Iso[2]     Sunje Dallmeier-Tiessen[1]     Mike Whalley[3]
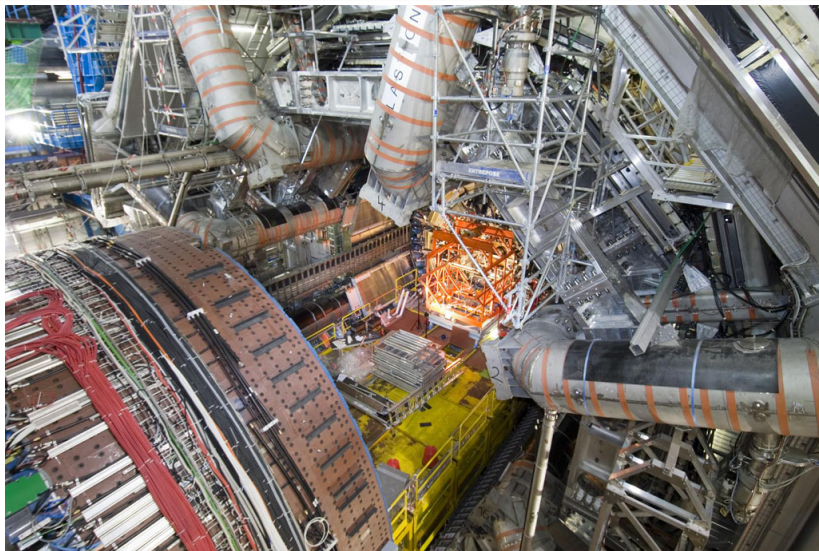
[1]CERN, Geneva, Switzerland
[2]Computer Science and Systems Engineering Dept.,Universidad de Zaragoza, Spain
[3]IPPP, Durham University, UK

November 28, 2012

Introduction
Architecture of the integration
Benefits derieved from the integration
Conclusions and future work

**Introduction**
Architecture of the integration
Benefits derieved from the integration
Conclusions and future work

# High Energy Physics

**Introduction**
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# INSPIRE



- The digital library of High Energy Physics (HEP) - Successor of SPIRES
- Over $10^6$ records
- Harvesting from ArXiv.org, manual curation
- Added services: Reference counting, Author disambiguation, figures, ...

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# HepData

- Data behind around 7000 publications - over 50'000 datasets. (Data behind tables, plots and additional datasets)

- Manual curation

- Additional data submitted by authors of the publications

- Developed and maintained at the Durham University



**The Durham HepData Project**

Durham University

REACTION DATABASE • DATA REVIEWS • PARTON DISTRIBUTION FUNCTION SERVER • OTHER HEP RESOURCES

**Reaction Database Standard Search Interface**

Database of Numerical HEP scattering cross sections

Enter query:

[        ] Search

examples: re gamma gamma, re p p --> p p and obs sig, exp cern

Search Help — Output Help — Form Search — Browse Keywords —
**Latest LHC DATA**

**To search the database:** Enter your query command comprising keyword-value pairs joined with Boolean ANDs. A null entry will retrieve all records.
**The basic keywords are:**
**reac** - the reaction (eg. p p --> charged x) also **beam** and **fsp**.
**obs** - the observable (eg. SIG, DSIG/DX, DN/DPT).
**sqrts** - lower bound of the centre-of-mass energy in GeV.
**exp** - the experiment/laboratory name (eg. ZEUS, CERN, LHC).
**date** - the year of the publication/preprint.
**auth** - the first author name on the paper.
**ref** - the publication/preprint reference.
Use % as the right or left truncation character to search for values beginning or ending with the value. All searches are
**case-insensitive.** More details in the Search Help

**Quick link to HepData data reviews**

- NEW Quarkonia data in Hadronic interactions
- Structure functions in DIS
- Single photon production in hadronic interactions
- Two-photon reactions leading to hadron final states
- Drell-Yan cross-sections
- Inclusive particle production data in e+e- interactions
- Hadronic total cross-sections (R) in e+e- interactions
- Low-energy neutrino cross-sections
- Event shapes in lepton-lepton and lepton-nucleon interactions

**Predefined event shape / jet searches**

- Event shapes (thrust, etc...)
- Event shapes in e+e- collisions
- Event shapes in non-e+e- collisions
- Jet production (in any process)
- Jet production in e+e- collisions
- Jet production in non-e+e- collisions

About HepData — Submitting your data to HepData

HepData also maintains the UK mirror of the **PDG**                    **Contact Us**

Introduction
Architecture of the integration
Benefits derieved from the integration
Conclusions and future work
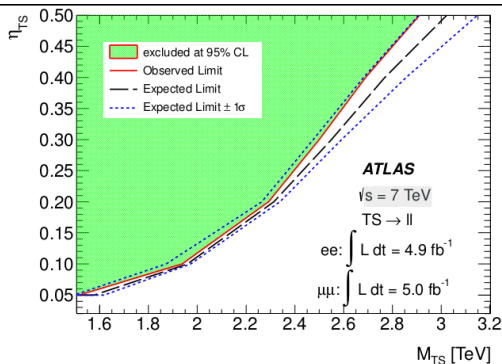
# Fragment of a sample publication page



**Figure 6.** Exclusion regions in the plane of $\eta_{TS}$ versus Torsion mass for the combination of dielectron and dimuon channels. The region above the curve is excluded at 95% CL.

## 13 Limits on Torsion models

The Torsion heavy state (TS) can be treated as a fundamental propagating field character-

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# HepData entry describing a single figure

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

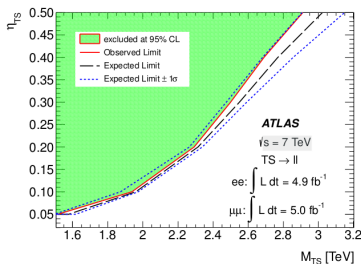# Figure inside a publication and the corresponding HepData entry



**Figure 6.** Exclusion regions in the plane of $\eta_{TS}$ versus Torsion mass for the combination of dielectron and dimuon channels. The region above the curve is excluded at 95% CL.

Introduction
Architecture of the integration
Benefits deriewed from the integration
Conclusions and future work

## Why should we integrate?

- INSPIRE contains publications and high-quality meta-data, HepData contains data.
- HepData stores only metadata of datasets and very minimal description of publications.
- Assignment of Digital Object Identifiers.

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# A bibliographic record before...

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

## ... and after

Introduction
**Architecture of the integration**
Benefits derived from the integration
Conclusions and future work

## Interaction between systems before the integration

Introduction
**Architecture of the integration**
Benefits derived from the integration
Conclusions and future work

## Interaction between systems before the integration



## ... and after the integration

Introduction
**Architecture of the integration**
Benefits derived from the integration
Conclusions and future work

## Harvesting

Assumption: We did not want to make deep modifications to HepData

- HepData does not maintain dates nor times of the last modification of a dataset
- Extensions of HepData:
  - HepData exports all the INSPIRE identifiers of records for which datasets are stored
  - HepData allows to access datasets by addressing them with INSPIRE id and number of a dataset

Introduction
Architecture of the integration
Benefits deriewed from the integration
Conclusions and future work

## Harvesting (2)

- INSPIRE reads the entire list of publication IDs stored in HD and retrieves all datasets related to those publications
- Datasets within a single publications are matched with existing INSPIRE records and necessary updates are applied

Introduction
**Architecture of the integration**
Benefits derived from the integration
Conclusions and future work

## Storing datasets in INSPIRE



■ Data qualifiers (MARC: 6531)
■ Data (Attached data file, not stored in MARC - not metadata)
■ Column headers and titles (MARC: 910)
■ Part of a reference to a paper (MARC: 786)

Introduction
**Architecture of the integration**
Benefits derived from the integration
Conclusions and future work

Table 2 ( T 2 ) [HIDE DATA] or as: plain text, AIDA, PyROOT, YODA, ROOT, mpl or jhepwork

| | RE : P P --> P P | P P --> P P | PBAR P --> PBAR P | PBAR P --> PBAR P |
|---|---|---|---|---|
| | SQRT(S) : 31.0-62.0 GeV | | | |
| SQRT(S) IN GEV | D(SIG)/D(T) (AT T=0) IN MB/GEV**2 | SLOPE IN GEV**-2 | D(SIG)/D(T) (AT T=0) IN MB/GEV**2 | SLOPE IN GEV**-2 |
| | | | | [HIDE DATA] |
| 31 | 93.0 ± 5.5 | 11.70 ± 0.62 | 90.4 ± 5.1 | 11.37 ± 0.60 |
| 31 | 74.0 ± 3.6 | 10.92 ± 0.15 | 75.6 ± 4.6 | 11.16 ± 0.20 |
| 53 | 72.5 ± 2.2 | 11.06 ± 0.11 | 78.0 ± 3.2 | 11.50 ± 0.15 |
| 62 | 66.4 ± 1.7 | 10.71 ± 0.08 | 72.3 ± 3.0 | 11.12 ± 0.15 |
| | Plot Select Plot | Plot Select Plot | Plot Select Plot | Plot Select Plot |

```
001__ 1157867
245__ $$9HEPDATA$$aAdditional data from: A MEASUREMENT OF anti-p p AND p p ELASTIC (...)
336__ $$tDATASET
520__ $$9HEPDATA
6531_ $$c4$$c3$$c2$$c1$$kSQRT(S)$$v31.0-62.0 GeV
710__ $$gAMES-BOLOGNA-CERN-DORTMUND-HEIDELBERG-WARSAW COLLABORATION
786__ $$hT 2.$$q2$$rCERN-EP/84-105$$w204422
8564_ $$uhttp://inspirehep.net/record/1157867/files/Data.txt$$ydata extracted from the table
910__ $$dSQRT(S) IN GEV$$$n0
910__ $$dD(SIG)/D(T) (AT T=0) IN MB/GEV**2$$n1$$tRE : P P --&gt; P P
910__ $$dSLOPE IN GEV**-2$$n2$$tP P --&gt; P P
910__ $$dD(SIG)/D(T) (AT T=0) IN MB/GEV**2$$n3$$tPBAR P --&gt; PBAR P
910__ $$dSLOPE IN GEV**-2$$n4$$tPBAR P --&gt; PBAR P
911__ $$x1$$y4
980__ $$aDATA
```

Introduction
Architecture of the integration
**Benefits derieved from the integration**
Conclusions and future work

## Improvements in the display



- Data qualifiers parsed and displayed in a more intuitive way
- Directly integrated plotting facility
- More compact view

Introduction
Architecture of the integration
Benefits derived from the integration
Conclusions and future work

# Plotting of data integrated within the dataset view

Introduction
Architecture of the integration
**Benefits derived from the integration**
Conclusions and future work

# HepData using INSPIRE search engine

- INSPIRE store meta-data about publications
- HepData stores only meta-data describing separate datasets
- Using publication-level attributes indexed in INSPIRE allows to increase the accuracy of HepData search

Introduction
Architecture of the integration
Benefits derieved from the integration
**Conclusions and future work**

## Conclusions

The integration betweem HepData and INSPIRE extends the
possibilities of users of both systems and improves their experience.
However, some problems have been identified:

- Changes in the display format of HepData can have
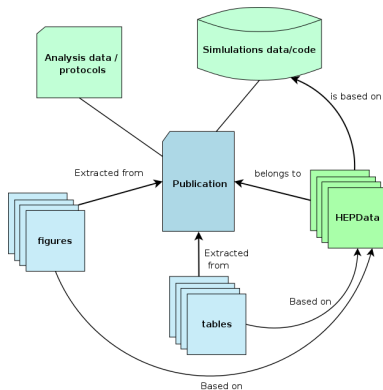  non-predictible effects on the INSPIRE dataset
- We need to harvest the entire database of HepData every time
- HepData is expected to provide the correct data, Inspire is
  expected to archive the publication data

Introduction
Architecture of the integration
Benefits derived from the integration
**Conclusions and future work**

## Harvesting 2.0

- Data exported in fixed XML format
- All modification in HepData marked with timestamps
- Retrieving only publications for which at least one dataset has been updated
- Datasets should never be removed in INSPIRE.
- Small dataset changes (mistakes in transcription) : Update the record
- Severe updates (Data has been completely reuploaded) : Create a new record linking to the old one

Introduction
Architecture of the integration
Benefits derieved from the integration
**Conclusions and future work**

## Open Linked Data

- Linking datasets with figures (and other artefacts) in an automatic way
- Exposing data behind publications in RDF
- Assignment of Digital Object Identifiers

Introduction
Architecture of the integration
Benefits derieved from the integration
**Conclusions and future work**

# Questions?

piotr.praczyk@cern.ch