# Zenodo/Invenio for EIC: Yellow Report and beyond

Maxim Potekhin
(BNL, NPPS)
04/29/2020

# Terminology

- **Zenodo** is an open science data repository at CERN
  - In a nutshell, storage+metadata
  - Any data within the set limits
- **Invenio** is a toolkit used to in a number of CERN systems *including* Zenodo
  - A complex and capable framework.
  - Framework, not a system. *An application is needed to make use of its functionality*.
  - *cf. Zenodo is an Invenio-based application.*
- **Invenio RDM** ("research data management") is a new product aiming to achieve
  - Portability (currently installing and configuring Invenio requires a high level of expertise)
  - Configurability i.e. eliminating the need for a custom app - a turnkey solution
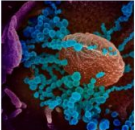  - ETA: late 2020

https://zenodo.org/ - named after Ζηνόδοτος, inventor of metadata in 280 BC

# Zenodo "in a nutshell"

- General purpose digital repository
- Version control
- Data (storage space) + Metadata (DB)
- Extensive query capabilities
  - Full-text search is in the works
- DOI management (**doi.org** integration)
- ORCID-aware
- Gateway to other repositories
- GitHub integration (citeable code)
- Currently a service instance at CERN, being transformed into a more portable system under the "Invenio RDM" brand

Zenodo in a nutshell

- **Research. Shared.** — all research outputs from across all fields of research are welcome! Sciences and Humanities, really!
- **Citeable. Discoverable.** — uploads gets a Digital Object Identifier (DOI) to make them easily and uniquely citeable.
- **Communities** — create and curate your own community for a workshop, project, department, journal, into which you can accept or reject uploads. Your own complete digital repository!
- **Funding** — identify grants, integrated in reporting lines for research funded by the European Commission via OpenAIRE.
- **Flexible licensing** — because not everything is under Creative Commons.
- **Safe** — your research output is stored safely for the future in the same cloud infrastructure as CERN's own LHC research data.

# Zenodo: durability

Safe

**— more than just a drop box!**

Your research output is stored safely for the future in same cloud infrastructure as research data from CERN's Large Hadron Collider and using CERN's battle-tested repository software Invenio, which is used by some of the world's largest repositories such as INSPIRE HEP and CERN Document Server.

# DOI, keywords, conference-awareness

# Motivations

- Managing documents and other materials is a universal necessity in the field
  - Consider the needs of the Yellow Report working groups (papers, presentations, tables etc)
  - Not a replacement of the Wiki (which is not a document handling system in the first place)
- Not too many products exist in that area
  - DocDB is used at FNAL, BNL and a few other places, it's an aging product, no clear API
  - CERN CDS is not portable (NB shares the Invenio back-end with Zenodo)
- In EICUG here is currently not a single accepted solution or a policy
- The new EIC Software website is not designed as a general purpose document store (scalability, lack of proper metadata etc)
- Zenodo is an obvious contender

# Invenio at BNL

- BNL SDCC is testing installation of **Invenio RDM**, BNL being an official partner as a test site
    - Install is currently pretty hard and not ready for prime time
    - In fairness, this is a complex product with many moving parts
- sPHENIX moved to use a custom app based on pre-RDM Invenio to use in lieu of DocDB
- In more recent news, Zenodo instance(s) have been created
    - TBD in an upcoming meeting

# Recent Zenodo activities

- Approved for PHENIX Data and Analysis Preservation (DAP)
- "Zenodo Communities" - see next slide - functional testing started
  - A "PHENIX Collaboration" community created, started populating it with materials
- Communication with the developers, looking for guidance regarding
  - Possible future data migration from Zenodo to Invenio RDM
  - Feature requests for community management
  - Storage allocation and use pattern discussion
- GitHub integration - "nice to have" but not core - initial testing done
  - Additional cloud replica of your GitHub release tagged with arbitrary metadata (discoverability)
  - Citeable via DOI

# Zenodo Community (another way to tag material)

- A way to organize material, and to consistently attribute materials to a collaboration/project/experiment - keeping a consistent brand
- An improvement in visibility/discoverability/PR
  - An addition to the already existing metadata query aids in discovery of materials
- Anyone can upload a material to the community which is subject to **curation**
  - The curator gets notified and inspects the submission
    - If accepted, it becomes posted under the community umbrella
    - If rejected, it still remains on Zenodo site but is not officially owned/acknowledged by the community, this is an accordance to the "open access" platform
  - *There is currently one curator per community and there is no easy way to transfer this duty to a different account* (something few people expected) but a fix is on the way according to the lead developer and other team members. Unofficial ETA is late 2020.

# A Community Example

# Advanced search capabilities

By default all searches are sorted according to an internal ranking algorithm that scores each match against your query. In both the user interface and REST API, it's possible to sort the results by:

- Most recent
- Publication date
- Title
- Conference session
- Journal
- Version

## Regular expressions

Regular expressions are a powerful pattern matching language that allow to search for specific patterns in a field. For instance if we wanted to find all records with a DOI-prefix 10.5281 we could use a regular expression search:

**Example:** `doi:/10\.5281\/.+/`

Careful, the regular expression must match the *entire* field value. See the regular expression syntax for further details.

## Missing values

It is possible to search for records that either are missing a value or have a value in a specific field using the `_exists_` and `_missing_` field names.

**Example:** `_missing_:notes` (all records without notes)

**Example:** `_exists_:notes` (all records with notes)

## Advanced concepts

### Boosting

You can use the boost operator `^` when one term is more relevant than another. For instance, you can search for all records with the phrase *open science* in either *title* or *description* field, but rank records with the phrase in the *title* field higher:

**Example:** `title:"open science"^5 description:"open science"`

### Fuzziness

You can search for terms similar to but not exactly like your search term using the fuzzy operator `~`.

**Example:** `oepn~`

Results will match records with terms similar to `oepn` which would e.g. also match `open`.

### Proximity searches

A phrase search like `"open science"` by default expect all terms in exactly the same order, and thus for instance would not match a record containing the phrase *"open access and science"*. A proximity search allows that the terms are not in the exact order and may include other terms inbetween. The degree of flexibility is specified by an integer afterwards:

**Example:** `"open science"~5`

### Wildcards

You can use wildcards in search terms to replace a single character (using `?` operator) or zero or more characters (using `*` operator).

**Example:** `ope? scien*`

Wildcard searches can be slow and should normally be avoided if possible.

## Fields reference

The table below lists the data type of each field. Below is a quick description of what each data type means and what is possible.

- **string**: Field does not require exact match (example field: `title` ).

12

# Policy issues

- Zenodo defines itself as an open science platform i.e. for the most part public
- It does have access tiers: private, restricted and public
  - "Restricted" means that a request for access is forwarded to the owner
  - Not designed to handle "roles" for large groups of people
- Consider the fact that DocDB instances are often protected
  - In reality I would say 95%+ of materials don't need to be protected

# GitHub/Zenodo mechanics

- A snapshot of a GitHub repository can be included in Zenodo organically+DOI
  - Integration/app link is in place: prepares and preserves tarballs of your releases
  - Makes your code easy to find (using the metadata) and to reference by a unique ID
  - Nice GUI
  - DOI reference to the code
- Easy to use
  - I tested this functionality and it was quite simple
  - DOIs take some time O(10min) to propagate to the DOI.org system

# Zenodo - GitHub panel - repo selection

# Zenodo - GitHub panel - published release

# Zenodo - GitHub panel - published release browser

# DOIs are an increasingly popular way to reference software

Persistent, durable link to archived software, can be nicely embedded in any page.

# GitHub/Zenodo integration benefits

- Not a core functionality by a long shot, however...
- ...provides a uniform way to reference digital products using DOI
- ...metadata is a good thing to have - better discoverability!
- ...can leverage the Zenodo "community" feature to organize materials and increase visibility
  - Cf. simulated data and the code used to produce it can be kept under the same umbrella
- Longer term - Data and Analysis Preservation
- In general, an "EIC Software" community on Zenodo may be a useful thing to have (papers, conference presentations etc)

# Zenodo - final notes

- A drop-in replacement for DocDB and as such can fulfill immediate needs of the EIC community
- ...but with more functionality, future proof and interesting new features
  - Can store almost any type of data and this may be helpful for the EIC YR
  - Obviously not a production resource i.e. we won't store datasets on Zenodo
  - However consider a versioned set of histograms as an example
  - Durable permalinks
  - API
- Search capabilities are significant and they will further improve
- An "EIC Software" community on Zenodo may be a useful thing to have (papers, conference presentations etc, good PR and a solid way to reference materials)
- A software group advice to the community?