

Using the Zenodo document and data repository

History:

- developed under **OpenAIRE** programme (EU initiative (FP8) to support open access)
- operated by CERN
- launched in 2013 as a document repository; since 2015 also for datasets



* Zenodotus (gr.: Ζηνόδοτος): first head librarian of the Library of Alexandria

Technical Information:



- primarily hosted in the CERN data centre infrastructure (OpenStack based)
- files stored on CERN EOS system
- metadata in PostgreSQL database
- uses a modified version of the **Invenio software**
 - originally developed by CERN; now an open source (MIT licence) project with multiple contributing institutions / persons
 - „modules“ for generic repositories, integrated library systems & RDM repositories (InvenioRDM = basically Zenodo)
 - also invenio-based: Inspire-hep, B2share, Caltech Library Catalog, ...



Features:

- free-of-charge
- size limit on single deposit: 50 GB
- deposits can be assigned to „communities“
- DOI service (including concept DOI and version DOI) + versioning of files

Ruya: Memory-Aware Iterative Optimization of Cluster Configurations for Big Data Processing

Jonathan Will, Lauritz Thamsen, Jonathan Bader, Dominik Scheinert, Odej Kao

Selecting appropriate computational resources for data processing jobs on large clusters is difficult, even for expert users like data engineers. Inadequate choices can result in vastly increased costs, without significantly improving performance. One crucial aspect of selecting an efficient resource configuration is avoiding memory bottlenecks. By knowing the required memory of a job in advance, the search space for an optimal resource configuration can be greatly reduced.

Therefore, we present Ruya, a method for memory-aware optimization of data processing cluster configurations based on iteratively exploring a narrowed-down search space. First, we perform job profiling runs with small samples of the dataset on just a single machine to model the job's memory usage patterns. Second, we prioritize cluster configurations with a suitable amount of total memory and within this reduced search space, we iteratively search for the best cluster configuration with Bayesian optimization. This search process stops once it converges on a configuration that is believed to be optimal for the given job. In our evaluation on a dataset with 1031 Spark and Hadoop jobs, we see a reduction of search iterations to find an optimal configuration by around half, compared to the baseline.

Comments: 9 pages, 5 Figures, 3 Tables; IEEE BigData 2022. arXiv admin note: substantial text overlap with [arXiv:2206.13852](https://arxiv.org/abs/2206.13852)

Subjects: **Distributed, Parallel, and Cluster Computing (cs.DC)**

ACM classes: C.2.4; I.2.8; I.2.6

Cite as: [arXiv:2211.04240](https://arxiv.org/abs/2211.04240) [cs.DC]

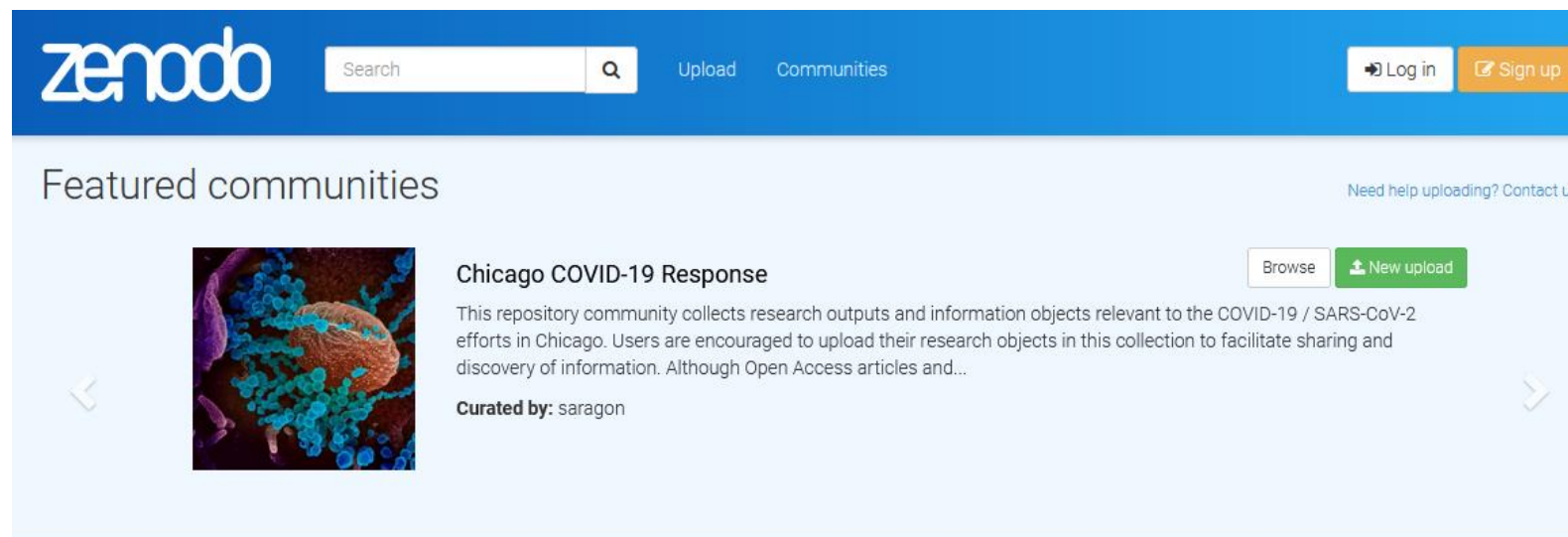
(or [arXiv:2211.04240v1](https://arxiv.org/abs/2211.04240v1) [cs.DC] for this version)

<https://doi.org/10.48550/arXiv.2211.04240> 

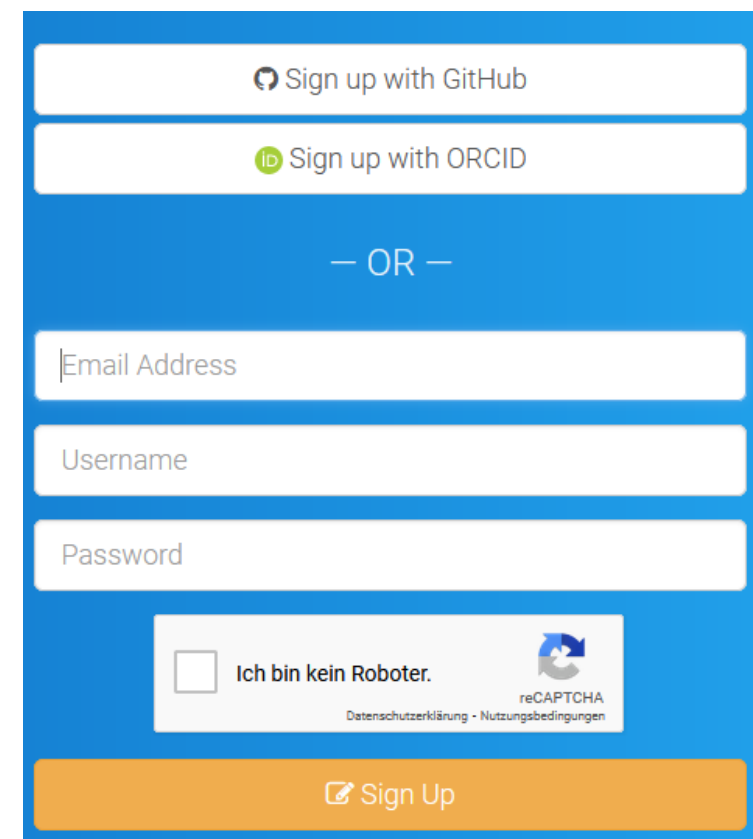


Features:

- free-of-charge
- size limit on single deposit: 50 GB
- deposits can be assigned to „communities“
- DOI service (including concept DOI and version DOI) + versioning of files
- non-static metadata
- Different access levels: open access, restricted, closed (metadata is always CC0 → FAIR compliant)
- GitHub integration
- metadata deposited in OpenAIRE (mandatory for EU-funded projects)
- no mandatory file formats / types
- user is responsible for upload (copyright, data protection ...)



The screenshot shows the Zenodo homepage. At the top, there is a blue navigation bar with the Zenodo logo on the left, a search bar in the center, and 'Upload' and 'Communities' links on the right. Further right are 'Log in' and 'Sign up' buttons. Below the navigation bar, the main content area is titled 'Featured communities'. The first featured community is 'Chicago COVID-19 Response', which includes a thumbnail image of a virus particle, a description of the repository's purpose, and a 'Curated by: saragon' note. There are 'Browse' and 'New upload' buttons next to the community name. A link for 'Need help uploading? Contact us' is also present.



The screenshot shows the Zenodo sign-up form. It features a blue background with white input fields. The form starts with two social login options: 'Sign up with GitHub' and 'Sign up with ORCID'. Below these is a separator '— OR —'. The main form consists of three input fields: 'Email Address', 'Username', and 'Password'. At the bottom of the form is a reCAPTCHA widget with the text 'Ich bin kein Roboter.' and a 'reCAPTCHA' logo. A large orange 'Sign Up' button is positioned at the very bottom of the form.

Recent uploads

July 2, 2018 (v2022-11-03) Dataset Open Access

View

Gene Ontology Data Archive

Carbon, Seth; Mungall, Chris

Archival bundle of GO data release.

Uploaded on November 7, 2022

49 more version(s) exist for this record

November 6, 2022 (v139) Dataset Open Access

View

A large-scale COVID-19 Twitter chatter dataset for open scientific research - an international collaboration

Need help?

Contact us

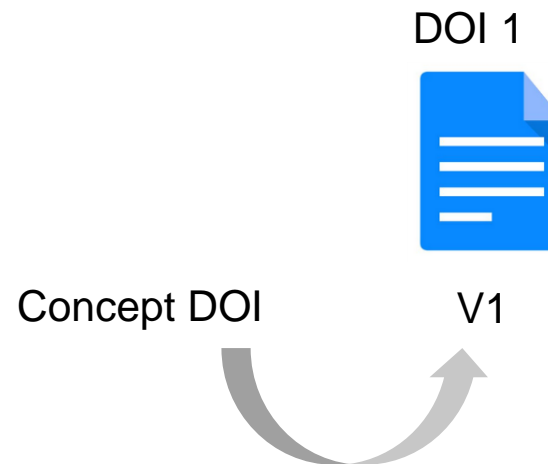
Zenodo prioritizes all requested related to the COVID-19 outbreak.

We can help with:

- Uploading your research data, software, preprints, etc.
- One-on-one with Zenodo supporters.
- Quota increases beyond our default policy.
- Scripts for automated uploading of larger datasets.

-
- To sign up:
 - basically any email address → works with a shared mail box etc (useful for managing a community)
 - „full“ functionality after 14 days

- Metadata and data can be uploaded / entered in „staging mode“ → can be freely changed and saved
- Uploads are tied to the account of the uploader (= only account that can edit / correct metadata)
- After publishing:
 - metadata can still be changed and corrected (e.g. add ORCID, fix a typo in the description, add a new publication where the dataset / software/ ... was used, etc)
 - data is permanently fixed → changing data creates a **new version** (changes are saved incrementally) with a new DOI
 - Zenodo provides a so-called concept DOI



Delete
Save
Publish

New upload

restrictions: (1) Upload minimum one file or file in required fields (marked with a red asterisk). (2) Press 'Save' to save your upload for editing later. (3) When ready, press 'Publish' to finalize and make your upload public.

Files
Choose files
Start upload

Drag and drop files here

— or —

[Choose files](#)

(minimum 1 file required, max 50 MB per dataset - contact us for larger datasets)
If you're experiencing issues with uploading larger files, read our FAQ section on file upload issues.

Communities
recommended

Specify communities which you wish your upload to appear in. The owner of the community will be notified, and can either accept or reject your request. Please make sure your record complies with the content policy of the communities you add; reported abuse will be followed by account inactivation.

FAU
Friedrich-Alexander-Universität
Erlangen-Nürnberg - FAU

FRASCAL
Fracture across Scales - GRK 2423
FRASCAL

Upload type
required

Publication
Poster
Presentation
Dataset
Image
Video/Audio
Software
Lesson
Physical Object
Workflow
Other

Publication type
Journal article

Basic information
required

Digital Object Identifier Optional. Did your publisher already assign a DOI to your upload? If not, leave the field empty and we will register a new DOI for you. A DOI allows others to easily and unambiguously cite your upload. Please note that it is NOT possible to edit a Zenodo DOI once it has been registered by us, while it is always possible to edit a custom DOI.

[Reset DOI](#)

Publication date Required. Format: YYYY-MM-DD. In case your upload was already published elsewhere, please use the date of first publication.

Title Required.

Authors Optional.

[Add another author](#)

Related/alternate identifiers

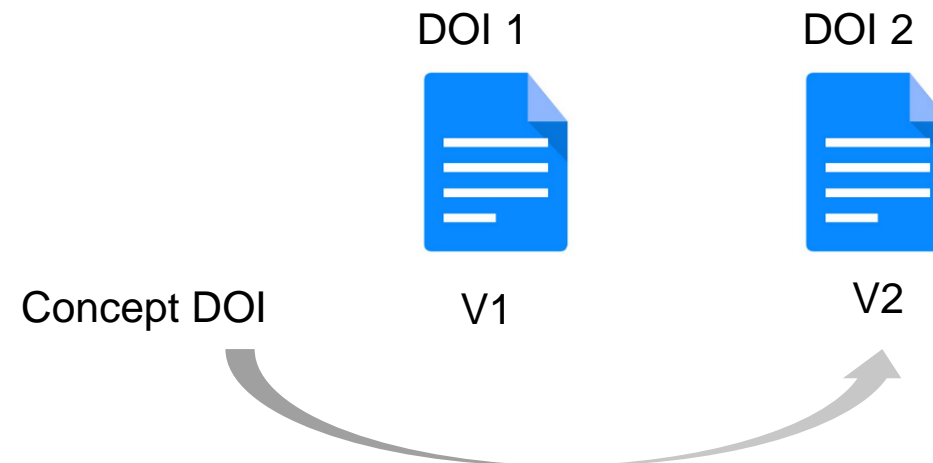
recommended

Specify identifiers of related publications and datasets. Supported identifiers include: DOI, Handle, ARK, PURL, ISSN, ISBN, PubMed ID, PubMed Central ID, ADS Bibliographic Code, arXiv, Life Science Identifiers (LSID), EAN-13, ISTC, URNs and URLs.

Related identifiers			
e.g. 10.1234/foo.bar.56789		N/A	Optional. Resource type of the related identifier.
e.g. 10.1234/foo.bar.56789		N/A	Optional. Resource type of the related identifier.
e.g. 10.1234/foo.bar.56789		N/A	Optional. Resource type of the related identifier.

Metadata close to DataCite standard

- Metadata and data can be uploaded / entered in „staging mode“ → can be freely changed and saved
- Uploads are tied to the account of the uploader (= only account that can edit / correct metadata)
- After publishing:
 - metadata can still be changed and corrected (e.g. add ORCID, fix a typo in the description, add a new publication where the dataset / software/ ... was used, etc)
 - data is permanently fixed → changing data creates a **new version** (changes are saved incrementally) with a new DOI
 - Zenodo provides a so-called concept DOI



License required ▾

Access right *

- Open Access
- Embargoed Access
- Restricted Access
- Closed Access

Required. Open access uploads have considerably higher visibility on Zenodo.

License *

- GNU General Public License v1.0 only
- GNU General Public License v1.0 or later
- GNU General Public License v1.0 only
- GNU General Public License v1.0 or later
- GNU General Public License v2.0 only
- GNU General Public License v2.0 or later
- GNU General Public License v2.0 only

Funding

Zenodo is integrated into reporting lines for res funding agency know!

- broad set of licence options
- four access options

Friedrich-Alexander-Universität Erlangen-Nürnberg - FAU

The Friedrich-Alexander University Erlangen-Nürnberg (FAU) is a strong research university with an international perspective and one of the largest universities in Germany, with around 40,000 students, 260 degree programmes, 4,000 academic staff (including over 580 professors), 200 million euros third-party funding, and 500 partnerships with universities all over the world. Teaching at the University is closely linked to research and focuses on training students in both theory and practice to enable them to think critically and work independently. The research itself also strikes the perfect balance between a theoretical approach and practical application. Only fast, direct and ideally free access to academic publications and primary research data unlocks the full potential of this research environment. Visit <https://www.fau.de/> to learn more about the FAU.

Zenodo FAU community collection

Feel free to use this Zenodo collection FAU as a central electronic archiving and publication platform for research data which can't be archived or published on OPUS FAU (<http://opus4.kobv.de/opus4-fau/home>), RRZE basis storage (https://www.rrze.fau.de/files/2017/06/Betreuungsvereinbarung_Basis-Storage.pdf) or institutional websites. Qualified scientific research data of all members of FAU may be published here free of charge. Corresponding full texts to the research data shall be published on OPUS FAU or disciplinary repositories. Thus all publications and research data sets are permanently available to the global public and are searchable and citable via catalogues and search engines.

For the creation of data management plans the tools RDMO (<https://rdmorganiser.github.io/>) or DMPonline (<https://dmponline.dcc.ac.uk/>) can be useful. For further information visit e.g. the Digital Curation Centre (DDC) (<http://www.dcc.ac.uk/resources/data-management-plans>). FAU also offers advice on legal and organizational matters to all researchers who would like to publish their research in open access media, via its Open Access Policy (<urn:nbn:de:bvb:29-opus4-68651>).

Only free and direct access to academic publications and primary research data unlocks the full potential of a digital research environment, academic networks and integrated research databases. The University thus wants to promote open access in the long term and has established an open access publishing fund to cover these costs (<https://ub.fau.de/en/writing-publishing/open-access/>).

Since this community is operated by the 'Referat Open Access', opening and making the data available (i.e., maximum restricted access, please no closed access) is desired.

Subject-based research data repositories

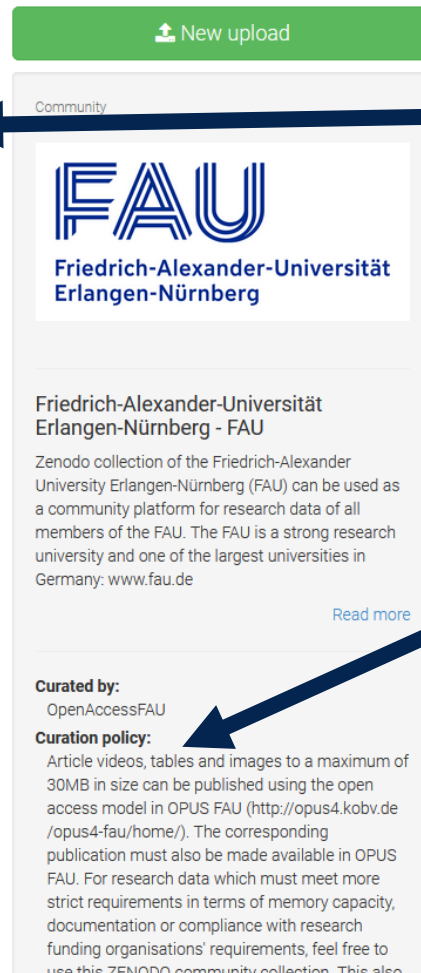
Other research data repositories can be found sorted by subjects on the web-based directory of the project re3data.org (Registry of Research Data Repositories, <http://www.re3data.org/browse/by-subject/>).

Data citation standards

Using the citation standard offers proper recognition to authors as well as permanent identification through the use of global, persistent identifiers in place of URLs. Use of universal numerical fingerprints (UNFs) guarantees to the scholarly community that future researchers will be able to verify that data retrieved is identical to that used in a publication decades earlier, even if it has changed storage media, operating systems, hardware, and statistical program format.

Following are two authentic examples of replication data citations:

From International Studies Quarterly, King and Zeng, 2006, p. 209:



Long description what the community is about.

Curation policy:

What documents are accepted, who decides etc.

jhrohrwild/Public_Repository OFF

GitHub / Releases [Create release ...](#)

Get started!

- 1 Flip the switch**

Toggle the switch below to turn on/off automatic preservation of your repository.

OFF
- 2 Create a release**

Go to GitHub and create a release. Zenodo will automatically download a .zip-ball of each new release and register a DOI.

[jhrohrwild/Public_Repository](#)

jhrohrwild Update citation.cff

1 contributor

27 lines (27 sloc) | 673 Bytes

```
1  cff-version: 1.2.0
2  title: My test software
3  message: >-
4    If you use this software, please cite it using the
5    metadata from this file.
6  type: software
7  authors:
8    - given-names: Jürgen
9      family-names: Rohrwild
10     email: juergen.rohrwild@fau.de
11     affiliation: >-
12       Friedrich-Alexander-Universität
13       Erlangen-Nürnberg
14     orcid: 'https://orcid.org/0000-0002-1167-0339'
15 abstract: >-
16   This is a short, but very informative abstract. It
17   describes what the software does and where it can
18   be used to great effect.
19 keywords:
20   - data management
21   - git
22   - data analysis
23   - python
24 license: Apache-2.0
25 commit: Number
26 version: '1.0'
27 date-released: '2022-11-01'
```

[Give feedback](#)

- publications are tied to the account that uploaded the file
 - get messy if someone changes affiliations / jobs and the metadata has to be fixed
 - consider a team account with a shared mailbox (ub-zenodo@fau.de)
- Communities are nice to present the collected publications (but accounts need to be 14 days old to be able to found a community)
- Zenodo is very user-friendly when it comes to metadata, but you cannot change the data files (or delete them) once published
- Linking Zenodo and GitHub is useful:
 - profit from minted DOI and persistent citability
 - all advantages of GitHub
 - recommended by DFG for published research software in mathematics