# Data Management Makes Machine Learning Easier

Till Korten, Helene Hoffmann, Peter Steinbach, Özlem Özkan

Data management days, HZDR
October 29, 2024

# Helmholtz AI Consulting for Matter



- Based at Helmholtz-Zentrum Dresden-Rossendorf

- Working in research field matter

- **Exploration voucher: 80 working hours of our time**

- Data management makes our collaboration much more effective

- A few practical examples for each letter in FAIR Data

# **F**indable data gets used

**F**AIR

- AI benchmarks

- Training material for foundational models

- As teaching material in machine learning courses

- For AI research projects

- If your data is findable, you don't need to find AI experts to collaborate with - they will find you!

**Collaborators find you**

**F**AIR

# **F**indable data gets used

Example (>4000 citations) :

> *Deng, L.* ***The MNIST Database of Handwritten Digit Images for Machine Learning Research*** *[Best of the Web]. IEEE Signal Processing Magazine 2012, 29 (6), 141–142. https://doi.org/10.1109/MSP.2012.2211477.*

**Used data gets cited**

# **A**ccessible data saves time and effort

FAIR

| Where is your data? | | How can you share it? |
|---|---|---|
| On paper | ➜ | Digitise it ~ hours to weeks |
| On a computer | ➜ | Upload to cloud ~ 1h/GB |
| On a portable harddisk | ➜ | Borrow hard disk ~ min to days |
| On a Fileserver | ➜ | Upload to cloud ~ 20min/GB |
| On cloud storage | ➜ | Send link ~ 10 s |
| In a public repository | ➜ | Send DOI ~ 10 s |

# Common example

**FAIR**

- 6 GB on a fileserver:

  - 2h: find a suitable cloud storage

  - 2h: file upload

  - 4h: write a data loader

  - 2h: file download

- 10h total

# Why machine learning experts love **MNIST**

**FAIR**

- Data is ready for machine learning in 2 min -> see notebook

- You can achieve this by uploading a simple python package with a few lines of code to https://pypi.org

```python
def load_data_from_cloud(target_path: Path):
    from nc_py_api import Nextcloud
    username, password = read_credentials()
    nc = Nextcloud(nextcloud_url='https://syncandshare.desy.de',
                   nc_auth_user=username, nc_auth_pass=password)
    # download file from the cloud
    files = nc.files.find(['like', 'name', target_path.name])
    nc.files.download2stream(files[0], target_path)
```
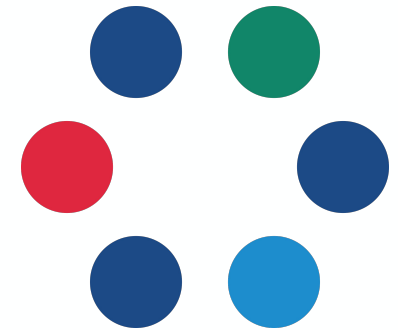
- Example: https://github.com/psteinb/b3get

# Interoperable data is easy to open

FAIR

- **Open** file format

  - **Open** source (ideally python) libraries for opening the data exist

  - Metadata is **well structured** and **machine readable**

  - **Good Practice**: have script on how to load small example of your dataset
    *(this can also be automatically tested)*

- For supervised machine learning:
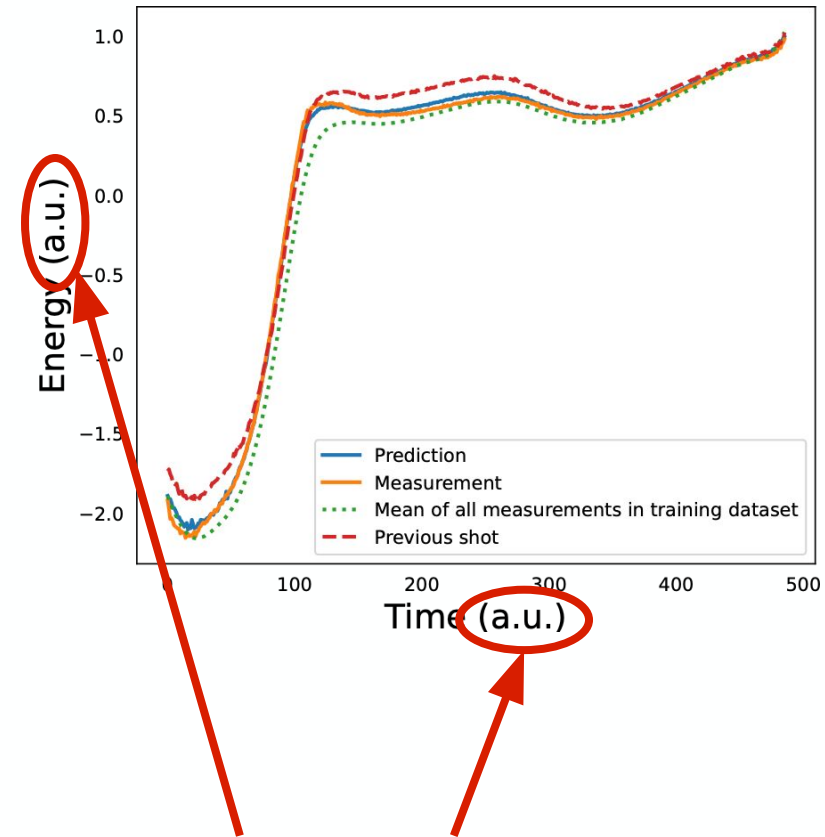
  - Data is **labeled**

Writing a workflow to read and convert an obscure format can take days

# Interoperable data is easy to open

- **Open** file format

  - **Open** source (ideally python) libraries for opening the data exist

  - Metadata is **well structured** and **machine readable**

  - **Good Practice**: have script on how to load small example of your dataset
    *(this can also be automatically tested)*

- For supervised machine learning:

  - Data is **labeled**

**Ex: OME Zarr**
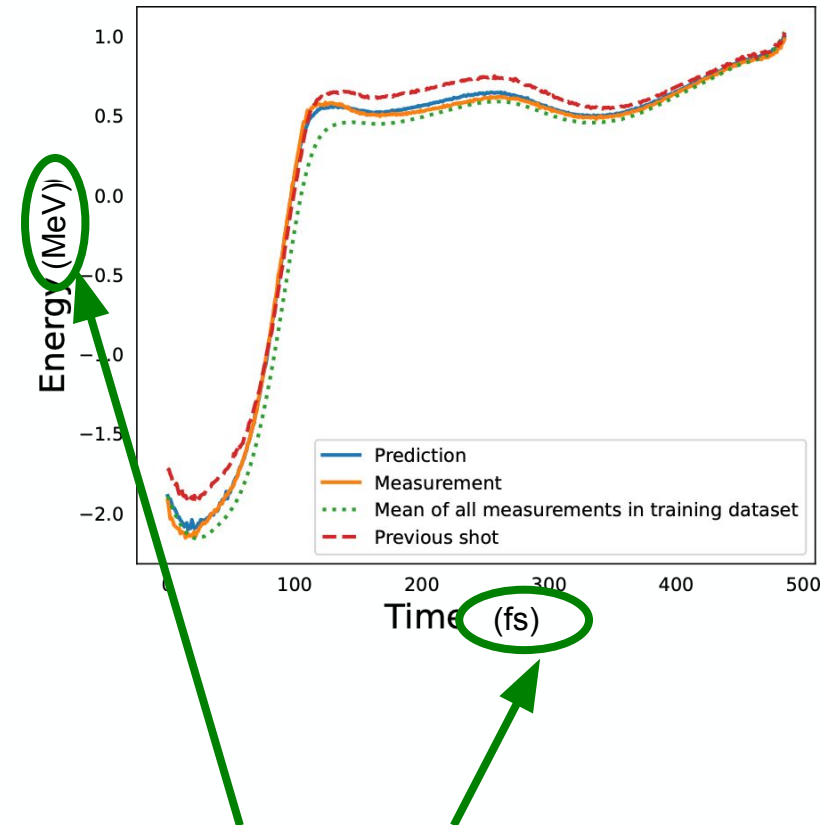*for Microscopy data*

# **R**eusable data is easy to work with FAI**R**

## Documentation

- The data is **well described** and **easy to understand**

- Descriptions are **detailed enough** that data can be used for more than the originally intended purpose

# **R**eusable data is easy to work with

**FAIR**

Documentation

- The data is **well described** and **easy to understand**

- Descriptions are **detailed enough** that data can be used for more than the originally intended purpose

# **R**eusable data is easy to work with     FAIR

## Licensing

- Data is **licensed** *(ideally with a permissive creative commons license)*
  - I once paid <u>35 dollars for my own paper</u> because of a **restrictive license**

# **R**eusable data is easy to work with FAIR

## Licensing

- Data is **licensed** *(ideally with a permissive creative commons license)*
  - My most cited paper is open access
    - ideally under a permissive license like CC BY 4.0

      → the lower the hurdles, the more citations

# **R**eusable data is easy to work with

FAI**R**

## Citable

- Data has a **DOI** (DFG counts data citations)

## Standardised

- Data and metadata **follow community standards**

Control.gif
sample-1.jpg
sample_two.jpg
anothersample.png

📂Control
 wildtype.tif
📂Samples
 <gene_code>-.tif
 <gene_code>+.tif

# Summary and Conclusion

- Good Data Management enables Good Machine Learning
  (reduce 80/20 split, FAIR data enables model building, can massively reduce time to solution, availability of data ensures transparency and progress)

- Curating a Data Set entails software and data science skills
  (collaborate where you can, the higher the load on the ML engineer - the less ML is done, well curated data ensures transparency and progress)

## Slides on figshare (CC-BY 4.0)



**Thank you for attention!**
We are happy to take questions, feedback or concerns.

**Shout out to our collaborators!**
Helmholtz AI Consulting Team HZDR
Helmholtz Metadata Collaboration (HMC)