

MUNI



CEITEC

# Enhancing Data Accessibility: Strategies and Tools for Effective (Meta)data Management *in life-sciences*

Adrián Rošinec, Tomáš Svoboda

IBERGRID 2023, 28.9.2023, Benasque

# CEITEC

- **Central European Institute of Technology**
- Established 2011
- Fields of science centre
  - Life sciences
  - Advanced materials
  - Nanotechnology
- Establishing institutions (Brno universities/institutes)



MUNI

VYSOKÉ UČENÍ  
TECHNICKÉ  
V BRNĚ

Mendelova  
univerzita  
v Brně



VUVeL

MUNI

CEITEC

# CEITEC in numbers (2022)

- 28 research groups
- **13 core facilities (CF)**
  - 720 users (from 22 countries)
- 270 research FTEs
- 241 publications (108 Q1, 24 T5)
- Ph.D. program (212 students)

# Core facilities



MUNI MU

**Cryo-Electron Microscopy and Tomography Core Facility**



Jiří Nováček, Ph.D.



MUNI MU

**Josef Dadok National NMR Centre**



Radovan Fiala



MUNI MU

**Proteomics Core Facility**



Prof. Zbyněk Zdráhal



MUNI MU

**Biomolecular Interactions and Crystallography Core Facility**



Josef Houser, Ph.D.



MUNI MU

**Nanobiotechnology Core Facility**



Jan Příbyl, Ph.D.



MUNI MU

**Multimodal and Functional Imaging Laboratory**



Michal Mikl, Ph.D.



MUNI MU

**Genomics Core Facility**



Boris Tichý, Ph.D.



MUNI MU

**Cellular Imaging Core Facility**



Milan Ešner, Ph.D.



MUNI MU

**Plant Sciences Core Facility**



Natallia Madzia Valasevich, Ph.D.



MUNI MU

**Bioinformatics Core Facility**



Vojtěch Bystrý, Ph.D.



MUNI MU

**Biological Data Management and Analysis Core Facility**



Radka Svobodová, Ph.D.

# What is role of CF

- Almost 700 special instruments
  - Cutting-edge technologies and instruments
  - MRIs, CryoEM, Xray Crystallography, ...
- Provide services to research groups, scientists
  - **Conduct experiments and provides data**
  - accessible to internal/external users from academy and industry at national and international levels

# Data management in CEITEC

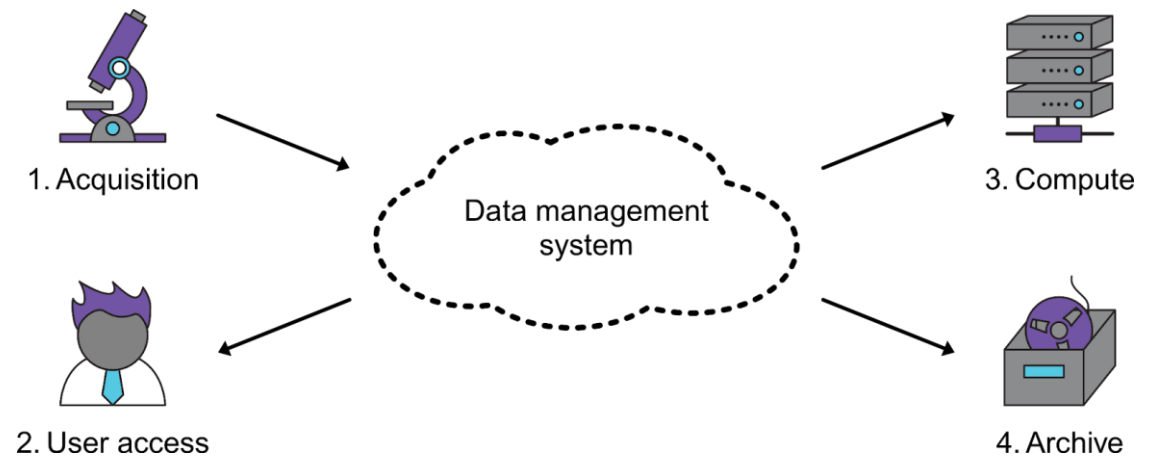
- Wide spectrum of produced data types
  - Big – CryoEM images, very small – CSVs, sensitive data, ...
- Currently no unified solution
  - No central data repository or storage facility
  - CFs/Labs may have large storage solutions or NAS or only laptop disk as storage solution
- Some CFs don't manage data
  - Passing data to the users using USB disks, CDs, random sharing services
- Institutional agreement to provide *unified solution across all CFs*

# Data management system

## – Requirements for DMS

- Handle data acquisition – from instruments
- Run data workflows (metadata ext., preprocessing)
- Ability to provide user access and sharing
- Ability to prepare dataset for publication
- Guarantee provenance
- Archiving and long-term preservation
- Mounting to processing / collaboration tools

– All this should be as  
“automated” as possible



# Metadata - integral part of each dataset

**Tuto stranu vyplňuje personál laboratoře!**

Měření prováděl: KR Projekt, popř. verze: test 5053B

#	Rozpis použitých měřících sekvencí:	Stim. protokol	Stim. log file	Fyziolog RAW file (pro Siemens <input checked="" type="checkbox"/> )	EEG <input checked="" type="checkbox"/>	ET <input checked="" type="checkbox"/>
1.	loc					
2.	t1-upr-seq-gal					
3.	SpinEcho FieldMap-4p					
4.	HCP_bold					
5.						
6.						
7.						
8.						
9.						
10.						
11.						
12.						
13.						
14.						
15.						
16.						
17.						
18.						

BP ExG: EKG Resp. GSR ACC \_\_\_\_\_ SIEMENS: EKG Respir. PT

Další záznamy o měření:  
Konzultace k Be. param

5053B strana 4 z 4

metadata - Notepad

File Edit Format View Help

User:  
Name: Prokop Buben  
ResearchGroup:  
Institution: CEITEC MU  
Affiliation: CF Plants

Project  
Name: Test project  
Number: 42

Operator:  
Name: Markéta Pernisová

Experiment:  
Name: phenotype  
Number:  
QR-codes: 107-108  
Mask: 5x4 full  
Sowing: 11.08.2022  
StartDate: 11.08.2022  
EndDate:  
FirstPhoto: 22. 08. 2022  
LastPhoto: 05. 09. 2022



# Why to bother with metadata?

- Metadata help to understand why datasets was created
  - context/equipment/method/source – patient
- Right metadata ontology helps with interoperability
  - Datasets might be reused within the community
- Finding datasets
  - Functions like searching/filtering/categorizing/identification
- Provenance – how the dataset was created
- Licensing

# Obtaining metadata

## – Manual annotations

- Human labor, error-prone, but very popular
- Goal to „ make annotating work enjoyable“ – nice GUIs, autocomplete, validators, context recognition

## – Automated annotations

- The base are metadata from instrument settings, administrative ones (owner of dataset), information about experiment, ...
- e.g. image recognition, capture/filter outputs from analytics tools and simulations

## – Lab entries and lab notebooks

# Example: Extraction of metadata from molecular dynamics

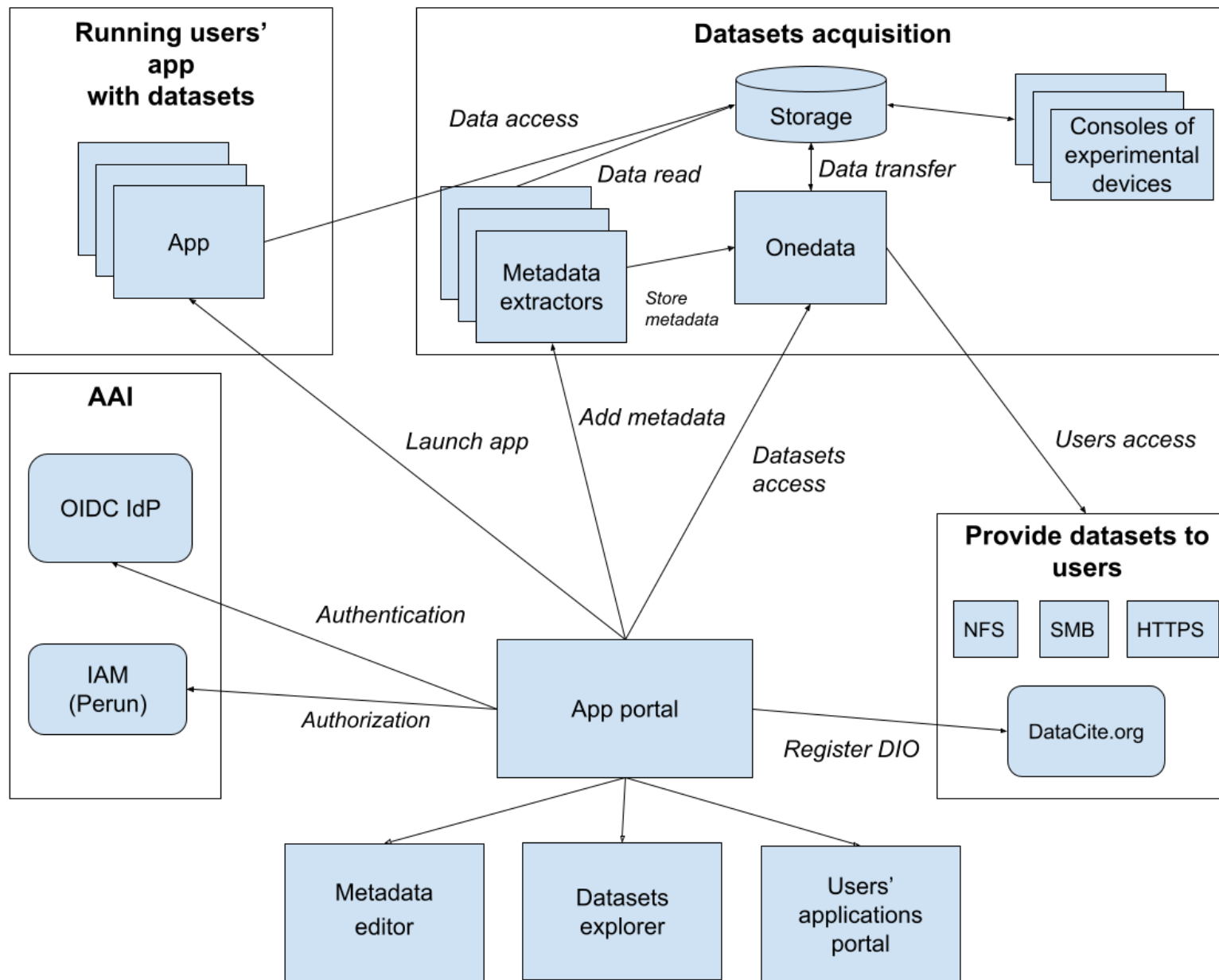
- Created parser of output from `gmx dump`
- Could be placed as a step of data workflow
  - Simulations ends -> output is stored on storage -> took by the DMS -> run preprocessing -> creating dataset
- Parsing on specific keywords and values important for dataset findability
  - based on the specification from the subset of the community
  - this community need accepted standard
- Results are outputted in json/yaml format stored in DB

# Metadata standards and ontologies

- Making CEITEC compliant with the world
- Community schemas
  - Some communities have rich ontologies – EDAM, Open Microscopy, EMPIAR
  - Some communities don't have anything
    - e.g. Molecular dynamics, plants/crops experiments
    - Motivate and provide support to establish partnership across CZ and for now Max Planck Institute
- Administrative and technical schemas  
(for CEITEC/National e-infrastructure)
- Institutional and national requirements for datasets
  - Administrative metadata for funders (project IDs, ...)
  - Part of EOSC CZ initiative

# Data management platform

- Aim to provide user interface to support data management life-cycle
- Ability to manage datasets self-service in a transparent manner
  - hide complexity of backend services
- Interconnecting many backend services
  - Onedata
  - Raw storage system
  - Metadata catalogue
  - Project control and IAM
  - DataCite
  - Public (meta)data repositories
  - Metadata schemas repositories and editors
  - FAIR checkers



# Challenges

- Make the solution accepted at least nation-wide
  - We have support of national e-infrastructure
- Development costs and sustainability of a project
  - Making use of what's available, minimum coding