# M U N I

# FAIR Molecular Dynamics

Adrián Rošinec

CEITEC CF Biodata

adrian@muni.cz

# Why?

— MD has evolved in last decades

— System size by a factor of 10^9

— Trajectory length by a factor of 10^9

— Ensemble size increased by a factor of 10^11

— From dozens of groups to thousands

    — Approx. 15% of all HPC use dedicated to MD

MUNI

# MD matured, but…

in terms of data management behind other fields

- Simulation efforts are lost
- No build-up on existing research
- No metanalysis possible
- Lack of reproducibility and quality checks
- Integration of AI / ML methods
- Poor interaction with other fields

MUNI

# ELIL5: FAIR

— Principles on using and sharing (scientific) data

— FAIR

- **Findable**: easy and transparent to search data
- **Accessible**: clear rules how to access the data
- **Interoperable**: common file formats, software inputs/outputs
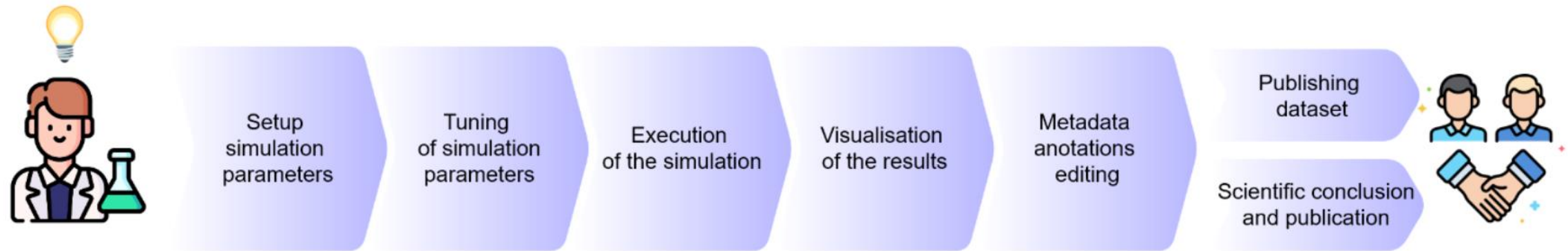- **Reusable**: data usable for more purposes

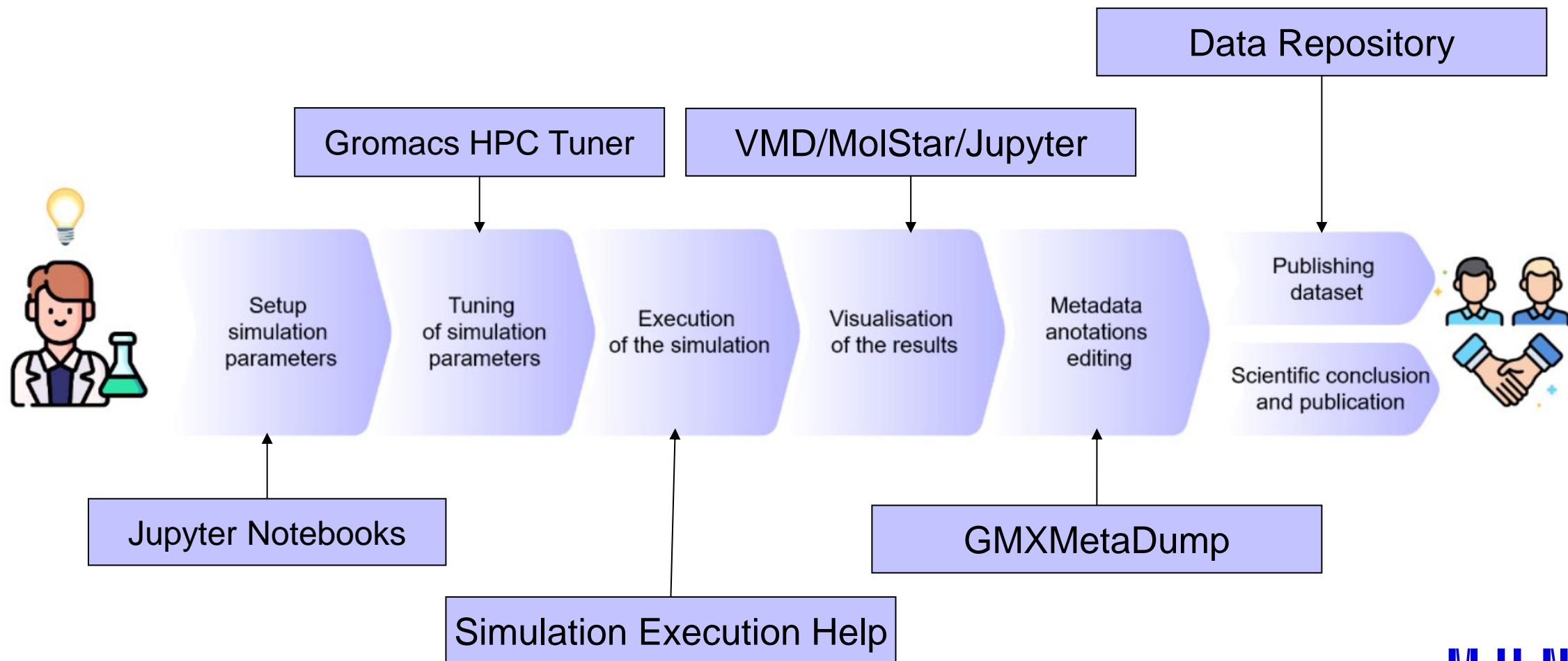More at: https://www.go-fair.org/fair-principles/

MUNI

# Challenges

– Standard MD data exchange formats

- Trajectory and traj. compression formats
- Trajectory identification – atom/residue names
- Full simulation settings-parameters of sim.

– Establish metadata ontologies and semantics

- Search based on contents (biomolecules, …) / parameters (thermal, …) / purpose of simulation

– Provenance – how was trajectory generated, hashes of files, steps to create trajectory, attached additional files

- Custom named / missing residues, non-standard force-fields or molecules,

– Quality control mechanisms and metrics

– Sharing – PIDs, community repositories

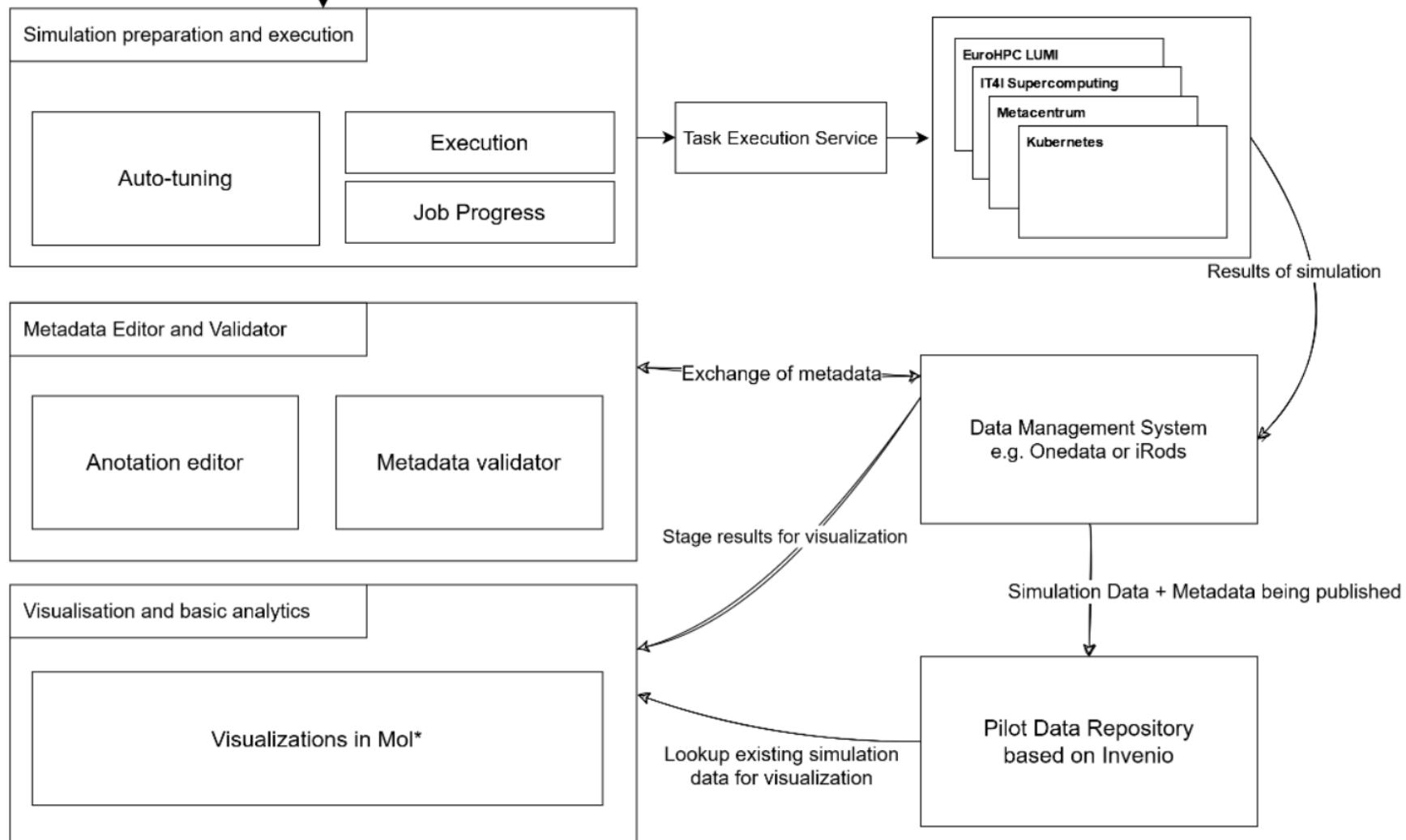– Cost of storing large amounts of data
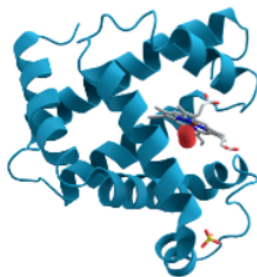
MUNI

# Broader context

— The typical MD scientific process

# Broader context

# "FAIRification" tools

– Automatization

– Remove burden of manual annotations

– Ensuring completeness

  – Automatic harvesting from Gromacs/… file formats
  – Linking to the biomolecule databases (e.g. PDB)
  – Administrative metadata such as publishing institution, authors, funding information, …

– Validations – quality control

  – Meeting repositories and community requirements for accepting datasets
  – e.g. has subject of simulation – biomolecule, has environmental conditions (temp, press), has used force field + attached custom FF
  – Based on ontologies

MUNI

# Gromacs MetaDump

A tool to describe molecular dynamics simulations with powerfull metadata

## About

This tool is designed to help you analyze and edit the metadata of a GMX file. You can upload a TPR file to analyze its metadata and download the metadata in JSON or YAML format.
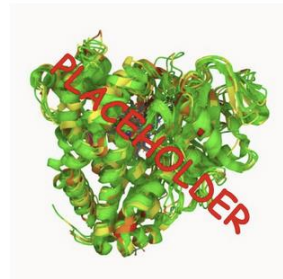
## Upload File

⬆

Upload files here
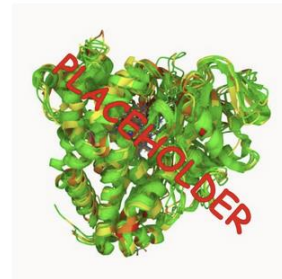
*(Only [*.tpr *.json *.yaml] are accepted)*
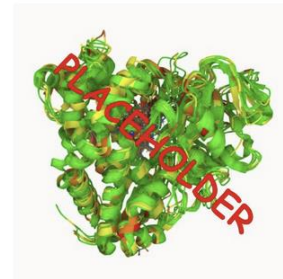
## Examples

**Protein 1**



VIEW METADATA

**Protein 2**



VIEW METADATA

**Protein 3**



VIEW METADATA

# Gromacs MetaDump

A tool to describe molecular dynamics simulations with powerfull metadata

## Selected File

📄 md_0.tpr   ⊖

## Analyze Metadata

FILE IDENTIFICATION    MAIN INFORMATION    **DETAILED INFORMATION**    OTHER

### Detailed information

nstcomm [step]
```
100
```

comm-mode
```
linear
```

lincs-iter
```
1
```

lincs-order
```
4
```

fourierspacing
```
0.16
```

constraint-algorithm
```
lincs
```

### van der Waals interactions

rvdw [nm]
```
1
```

dispcorr
```
enerpres
```
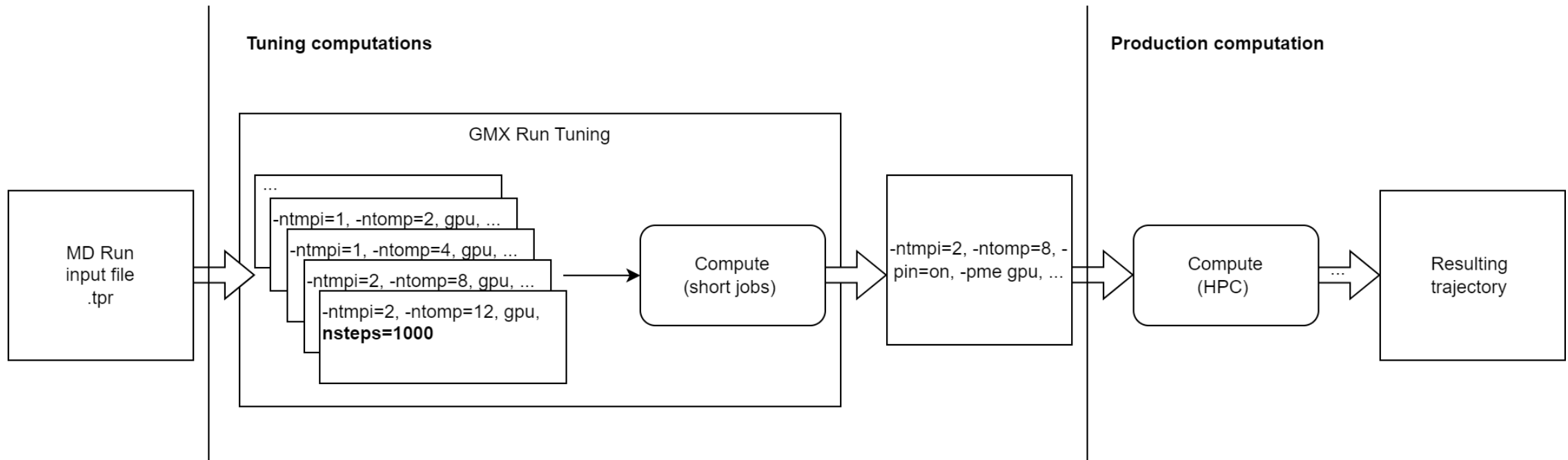
rvdw-switch [nm]
```
0.9
```

vdw-modifier
```
force-switch
```

## Simulation Preview

# Tuning and execution

# Experiment Tuning and Exec

– Interactive experiment setup

  – Choosing biomolecule, Equilibration, sampling process, setting simulation parameters

– (Auto)Tuning production run parameters

  – Avoiding common mistakes:
  – e.g. when running in paralell slows the computation
  – too small time-steps, not exploiting GPU or overexploiting them, right force fields, not minimized/equilibrated system, and many more
  – Running several small jobs (several seconds per job) to find better parameters before production run (several weeks)

– Remove burden of interacting with various computing interfaces

  – k8s, batch systems, …

– Provenance/Protocols – capture operations used to get trajectories

– Controlled => Simple => Correct

M U N I

# EU Efforts on building MDREPO

– Federated architecture

  – Method of sustainable funding to hold increments of tens of PBs per year

– Central component is metadata catalogue

  – One (features rich) interface to browse all MD data
  – Enables findability of datasets and points to "downstream" storage a.k.a repository (several nodes within nation, one national, …)

– Implemented quality control mechanisms

– Built as a custom solution

  – Is the CZ variant of repo platform (within NDI) more sustainable?
  – Is it worth extend functionality of CZ repo platform?
  – …

MUNI

Project MD-A00001 > Overview page ℹ

**DATA IN THIS PAGE**

### SARS spike receptor binding domain bound with FERR ACE2

`Trajectory` Classical MD

Theoretical model generated using Modeller from the PDB 6vw1

Authors: Vito Genna

Groups: IRB Barcelona, Orozco lab

Node: IRB Barcelona, MMB

Program: GROMACS

Version: 2019.1

### Spike glycoprotein

Gene: S

Organism: Severe acute respiratory syndrome coronavirus 2

UniProt ID: P0DTC2

### Angiotensin-converting enzyme 2

Gene: ACE2

Organism: Homo sapiens

UniProt ID: Q9BYF1

PDB Accession: 6VW1

**Structure of SARS-CoV-2 chimeric receptor-binding domain complexed with its receptor human ACE2**

Experimental method: x-ray

Organisms:   Homo sapiens; severe acute respiratory syndrome coronavirus 2

homo sapiens; human sars coronavirus; severe acute respiratory syndrome coronavirus 2
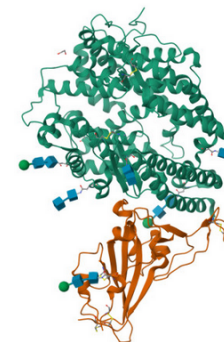
Keyword:   Cell invasion

Publication date: Tue Feb 18 2020

🔗 PDBE WEBSITE    🔗 RCSB WEBSITE    🔗 3DBIONOTES    🔗 PDBBIND

https://mdposit.mddbr.eu

OVERVIEW  TRAJECTORY  ANALYSES  DOWNLOADS

Project MD-A00001 > Trajectory page ⓘ

DATA IN THIS PAGE

Domains  Overall  Spike glycoprotein - Spike protein S1  Spike glycoprotein - BetaCoV S1-CTD  Spike glycoprotein - Disordered
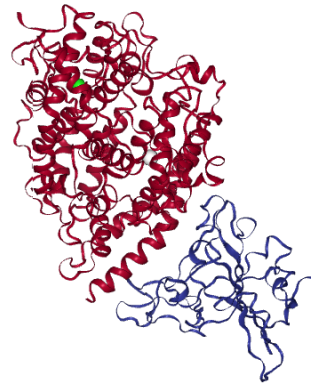
Spike glycoprotein - Receptor-binding domain (RBD)  Spike glycoprotein - Integrin-binding motif;  Spike glycoprotein - Receptor-binding motif; binding to human ACE2

Spike glycoprotein - Immunodominant HLA epitope recognized by the CD8+; called NF9 peptide  Angiotensin-converting enzyme 2 - Processed angiotensin-converting enzyme 2

https://mdposit.mddbr.eu

▼ Trajectory metadata

Counts

| System atoms | Proteins atoms | Proteins residues | Solvent molecules | Positive ions |
|---|---|---|---|---|
| 12584 | 12582 | 790 | 0 | 0 |

| Negative ions |
|---|
| 1 |

System box

| Type | Size X | Size Y | Size Z | Volume |
|---|---|---|---|---|
| Triclinic | 12.15 nm | 11.45 nm | 9.92 nm | 1379.09 nm³ |

Simulation

| Length | Timestep | Snapshots | Frequency | Force fields |
|---|---|---|---|---|
| 200.01 ns | 2 fs | 20001 | 10 ps | Not available |

# National Node of MDREPO

— Need to join the EU activity and build national node

— WIP mdrepo.eu

— Several groups publishing MD results

   — Robert Vácha at CEITEC, Vojtěch Spiwok at UOCHB, Michal Kolář at VŠCHT, Karel Berka at UPOL, …

— Need for (shared) data curator, quality control, support (steward) and tools

— Storage for CZ MD data (PBs/year)

— Publishing to the metadata catalogues – EU MDDB node

M U N I

# Sources

- [https://mddbr.eu/first-mddb-webinar-recording-now-online/](https://mddbr.eu/first-mddb-webinar-recording-now-online/)
- [http://arXiv:2407.16584](http://arXiv:2407.16584)
- Own project proposal

MUNI