# Sharing & reusing molecular dynamics data: what did we miss?

**Pierre Poulain**
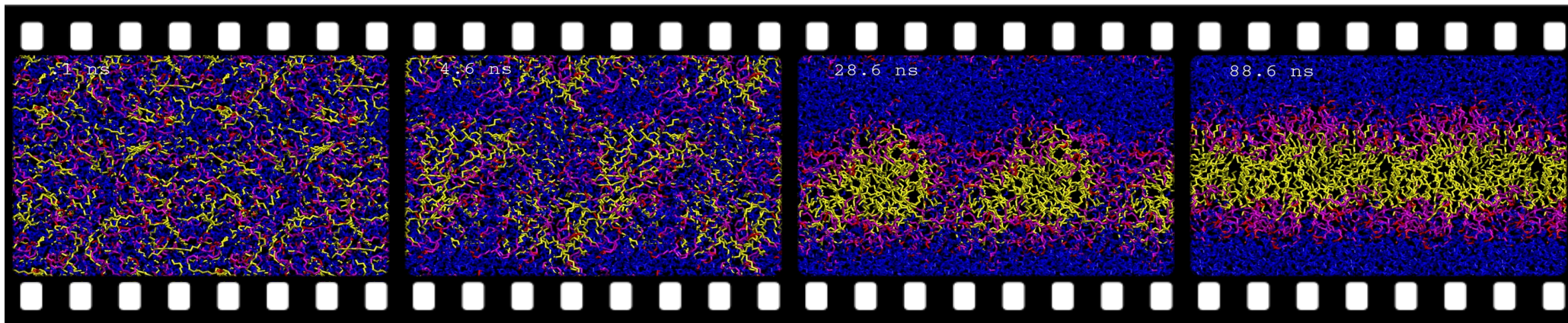
pierre.poulain@u-paris.fr

Laboratoire de Biochimie Théorique, Université Paris Cité

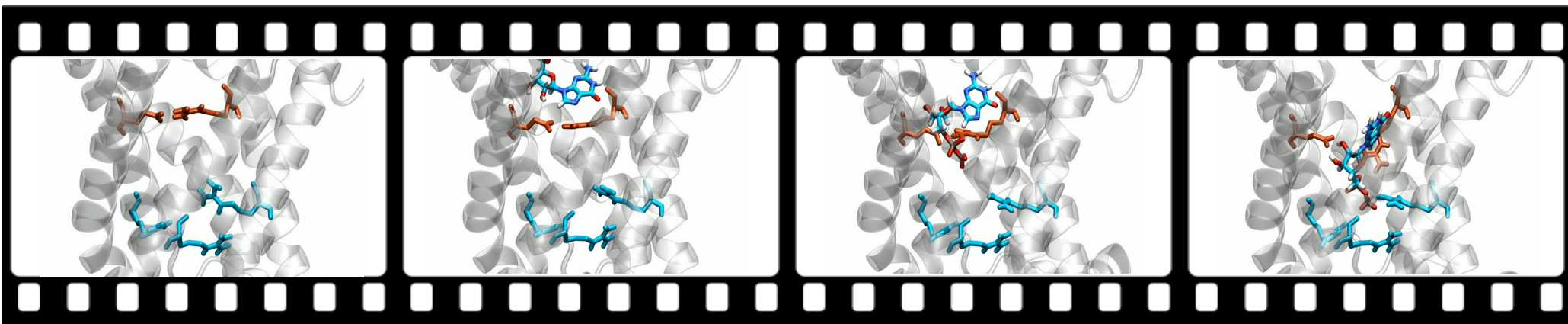**Open Science Days @UGA, Grenoble, 2024**

# What is molecular dynamics (MD)?



water + detergent

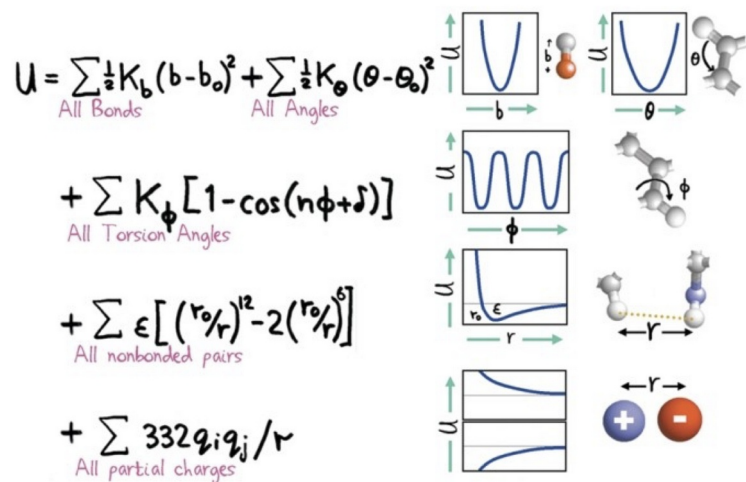Senac et al, Langmuir, 2017. Movie by Patrick Fuchs



Gagelin et al, Nature communications, 2023.

# Molecular dynamics simulations require resources

## expertise



M. Levitt

**GROMACS**
*fast, flexible & free*

## computer power



Source ; Photothèque CNRS/Cyril Frésillon (droits réservés)

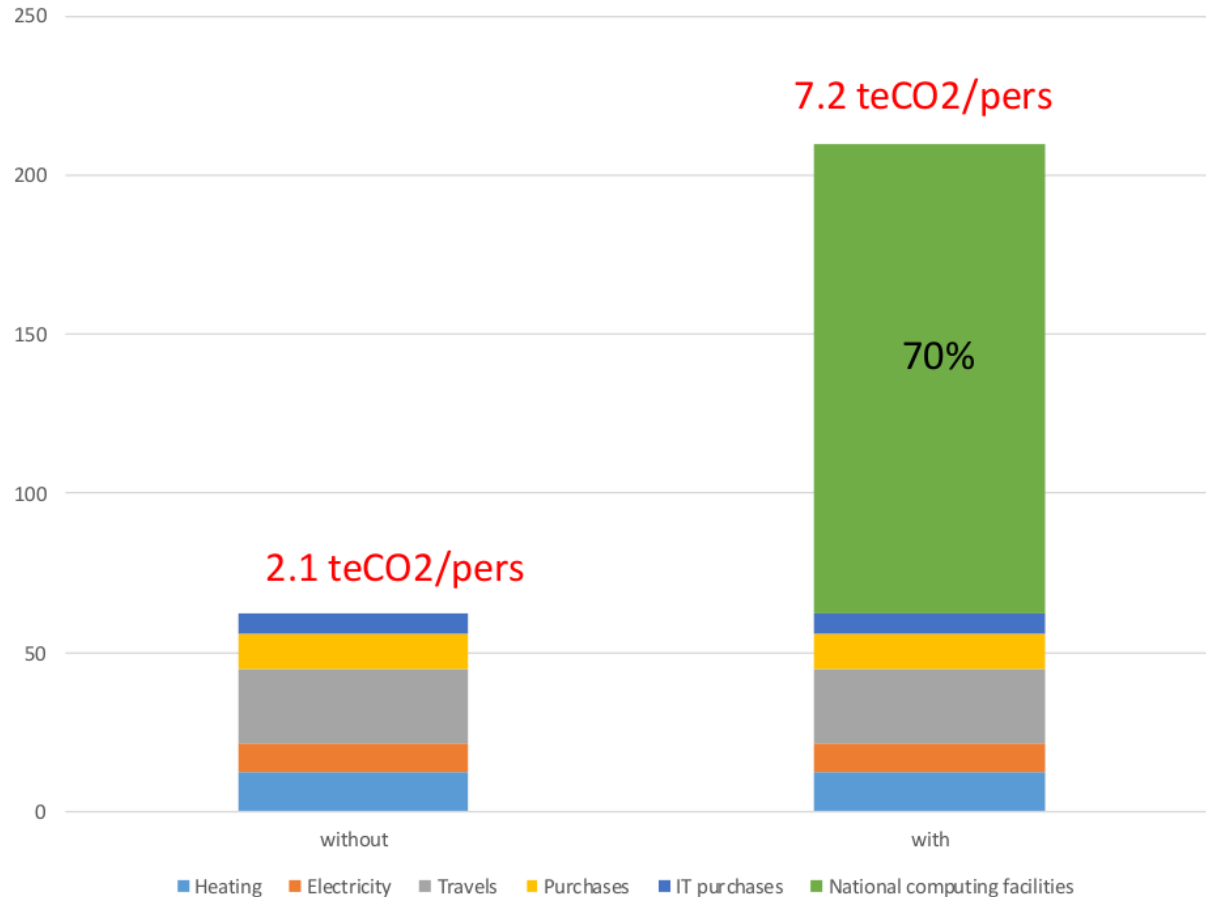# MD simulations require large computational resources → high cost

In 2023:

- GENCI, CT 7

  (simulation in biology)

- 160 Mh CPU

- 10 Mh GPU

- **Total cost: 8.5 M€**

Patrick Fuchs, 2024.

# MD simulations require large computational resources → high environmental cost

Élise Duboué-Dijon & Antoine Taly, LBT, 2024.

# Sharing research data is important

- Requirements from funders, institutions, or journals

- Open science

- Reproducibility

Editorial | Open access | Published: 21 February 2020

## No raw data, no science: another possible source of the reproducibility crisis

Tsuyoshi Miyakawa ✉

Molecular Brain **13**, Article number: 24 (2020) | Cite this article

JCIM JOURNAL OF CHEMICAL INFORMATION AND MODELING

**Using Open Data to Rapidly Benchmark Biomolecular Simulations: Phospholipid Conformational Dynamics**

Hanne S. Antila*, Tiago M. Ferreira, O. H. Samuli Ollila, and Markus S. Miettinen*

Cite this: *J. Chem. Inf. Model.* 2021, 61, 2, 938−949
Publication Date: January 26, 2021
https://doi.org/10.1021/acs.jcim.0c01299
Copyright © 2021 The Authors. Published by American Chemical Society. This publication is licensed under CC-BY.
Open Access

| Article Views | Altmetric | Citations |
|---|---|---|
| 2405 | 6 | 17 |

LEARN ABOUT THESE METRICS

# **Sharing research data in biology**

Community-approved data repositories for experimental (omics) data:

- SRA, GEO, ENA: genomics & transcriptomics

- PRIDE, MassIVE: proteomics & metabolomics

# Sharing MD simulation files

# Sharing MD simulation files

- **Many initiatives**: MoDEL, GPCRmd, NMRLipids…

- **No consensus data repository** for MD simulation files (yet)

# Sharing MD simulation files 😢

- **Many initiatives**: MoDEL, GPCRmd, NMRLipids…

- **No consensus data repository** for MD simulation files (yet)

- Use of generic **non-moderated** data repositories: Zenodo, Figshare, OSF, Dryad…

# Welcome to the Dark Matter of MD

Data that is **technically accessible**,

but neither **indexed**, **curated**,

or easily **searchable**.

# A universal search engine for MD open data

# MDverse

1. Find and index **scattered** MD data

2. Extract, enhance and explore **metadata**

3. Assess **accessibility** and **reusability**
(from FAIR principles)

# Step 1: Find and index scattered MD data

# Step 1: Find and index scattered MD data

# Step 1: Find and index scattered MD data



**250k files**
**2k datasets**
**14 TB of data**

# Step 1: Find and index scattered MD data



250k files
2k datasets
14 TB of data

1% of data in Zenodo = MD data

# Step 1: Find and index scattered MD data

# Step 2: Metadata = context

# Step 2: Metadata = context



No metadata (bad)

Natural language metadata (good) 🧑‍🔬

Controlled vocabulary metadata (better) 🤖

# Step 2: Metadata = context



No metadata
(bad)

Natural language 🧑‍🔬
metadata (good)

Controlled vocabulary 🤖
metadata (better)

Metadata for MD simulation: identity of simulated molecules,
temperature, length, force field, software…

Soup can idea from "Metadata for Effective
Research Data Management", Chiu et al., 2019

# Step 2: Sources of MD metadata



February 27, 2019

## SLIPID POPG-POPE 1:3 Bilayer Simulation (Last 100 ns, 150 mM NaCl, 310 K )

PEÓN, Antonio

Simulation of a POPG-POPE 1:3 bilayer of 500 lipids (126 L-POPG lipids and 374 POPE lipids, 250 per leaflet) is simulated for 500 ns using Gromacs v5.1.2 in water solution with Na+ counterions and 150 mM of NaCl. The SLIPID model is employed for lipids and TIP3P Water Model .

Trajectory (.xtc) is for the last 100 ns of a simulation of 500 ns with data saved every 10 ps. Additionally, the topology (.top) , simulation parameter file (.mdp), index file (.ndx),  portable binary run input file (.tpr) and the energy output file (.edr) are provided.

Files (5.2 GB)

| Name | Size | |
|---|---|---|
| index.ndx | 5.5 MB | Download |
| md5:5a8c0d796996b9432ec7a15919544a17 | | |
| m_400_500_PG_PE_1_3_slipid.xtc | 5.2 GB | Download |
| md5:177e0b0c513b1910a9bbe6329ba7a814 | | |
| md_02.tpr | 3.9 MB | Download |
| md5:3b00b35301d3079a63f7055e09f4089a | | |
| md_400_500.edr | 33.0 MB | Download |
| md5:9496927277c7b890c0a41974a3e6c746 | | |
| topol.top | 554 Bytes | Download |
| md5:1be0f066d982838ef60156ec2449d117 | | |

30 views    23 downloads

See more details...

Indexed in

OpenAIRE

**Publication date:**
February 27, 2019
**DOI:**
DOI  10.5281/zenodo.2579224
**Keyword(s):**
POPG-POPE 1:3 , SLIPID
**License (for files):**
Creative Commons Attribution 4.0 International

**Versions**

| Version 1 | Feb 27, 2019 |
|---|---|
| 10.5281/zenodo.2579224 | |

Cite all versions? You can cite all versions by using the DOI 10.5281/zenodo.2579223. This DOI represents all versions, and will always resolve to the latest one. Read more.

# Step 2: Sources of MD metadata

# Step 2: Sources of MD metadata

# Step 2: Uneven amount of metadata

# Step 2: Uneven amount of metadata



February 27, 2019    Journal article   Open Access

## SLIPID POPG-POPE 1:3 Bilayer Simulation (Last 100 ns, 150 mM NaCl, 310 K )

PEÓN, Antonio

Simulation of a POPG-POPE 1:3 bilayer of 500 lipids (126 L-POPG lipids and 374 POPE lipids, 250 per leaflet) is simulated for 500 ns using Gromacs v5.1.2 in water solution with Na+ counterions and 150 mM of NaCl. The SLIPID model is employed for lipids and TIP3P Water Model .

Trajectory (.xtc) is for the last 100 ns of a simulation of 500 ns with data saved every 10 ps. Additionally, the topology (.top) , simulation parameter file (.mdp), index file (.ndx), portable binary run input file (.tpr) and the energy output file (.edr) are provided.

Files (5.2 GB) ⌄

| Name | Size | |
|------|------|---|
| index.ndx | 5.5 MB | ⬇ Download |
| md5:5a8c0d796996b9432ec7a15919544a17 ❓ | | |
| m_400_500_PG_PE_1_3_slipid.xtc | 5.2 GB | ⬇ Download |
| md5:177e0b0c513b1910a9bbe6329ba7a814 ❓ | | |

March 22, 2021    Dataset   Open Access

## Simulations of a pulmonary surfactant monolayer with additional compounds

To be added.

Files (47.1 GB) ⌄

| Name | Size | |
|------|------|---|
| benzaldehyde.ndx | 4.8 MB | ⬇ Download |
| md5:987f07a7e7fab150db9713558d7dd8ad ❓ | | |
| benzaldehyde.top | 492 Bytes | ⬇ Download |
| md5:e90dbc7f37ff3ea109ca4aba127cd3ca ❓ | | |

# Step 2: Uneven amount of metadata



February 27, 2019 — Journal article, Open Access

## SLIPID POPG-POPE 1:3 Bilayer Simulation (Last 100 ns, 150 mM NaCl, 310 K )

PEÓN, Antonio

Simulation of a POPG-POPE 1:3 bilayer of 500 lipids (126 L-POPG lipids and 374 POPE lipids, 250 per leaflet) is simulated for 500 ns using Gromacs v5.1.2 in water solution with Na+ counterions and 150 mM of NaCl. The SLIPID model is employed for lipids and TIP3P Water Model .

Trajectory (.xtc) is for the last 100 ns of a simulation of 500 ns with data saved every 10 ps. Additionally, the topology (.top) , simulation parameter file (.mdp), index file (.ndx),  portable binary run input file (.tpr) and the energy output file (.edr) are provided.

Files (5.2 GB)

| Name | Size | |
| --- | --- | --- |
| index.ndx | 5.5 MB | Download |
| md5:5a8c0d796996b9432ec7a15919544a17 | | |
| m_400_500_PG_PE_1_3_slipid.xtc | 5.2 GB | Download |
| md5:177e0b0c513b1910a9bbe6329ba7a814 | | |

March 22, 2021 — Dataset, Open Access

## Simulations of a pulmonary surfactant monolayer with additional compounds

To be added.

Files (47.1 GB)

| Name | |
| --- | --- |
| benzaldehyde.ndx | |
| md5:987f07a7e7fab150db9713558d7dd8ad | |
| benzaldehyde.top | |
| md5:e90dbc7f37ff3ea109ca4aba127cd3ca | |

# Step 2: Uneven amount of metadata



February 27, 2019 — Journal article — Open Access

## SLIPID POPG-POPE 1:3 Bilayer Simulation (Last 100 ns, 150 mM NaCl, 310 K )

PEÓN, Antonio

Simulation of a POPG-POPE 1:3 bilayer of 500 lipids (126 L-POPG lipids and 374 POPE lipids, 250 per leaflet) is simulated for 500 ns using Gromacs v5.1.2 in water solution with Na+ counterions and 150 mM of NaCl. The SLIPID model is employed for lipids and TIP3P Water Model .

Trajectory (.xtc) is for the last 100 ns of a simulation of 500 ns with data saved every 10 ps. Additionally, the topology (.top) , simulation parameter file (.mdp), index file (.ndx), portable binary run input file (.tpr) and the energy output file (.edr) are provided.

Files (5.2 GB)

| Name | Size | |
|------|------|---|
| index.ndx | 5.5 MB | Download |
| md5:5a8c0d796996b9432ec7a15919544a17 | | |
| m_400_500_PG_PE_1_3_slipid.xtc | 5.2 GB | Download |
| md5:177e0b0c513b1910a9bbe6329ba7a814 | | |

March 22, 2021 — Dataset — Open Access

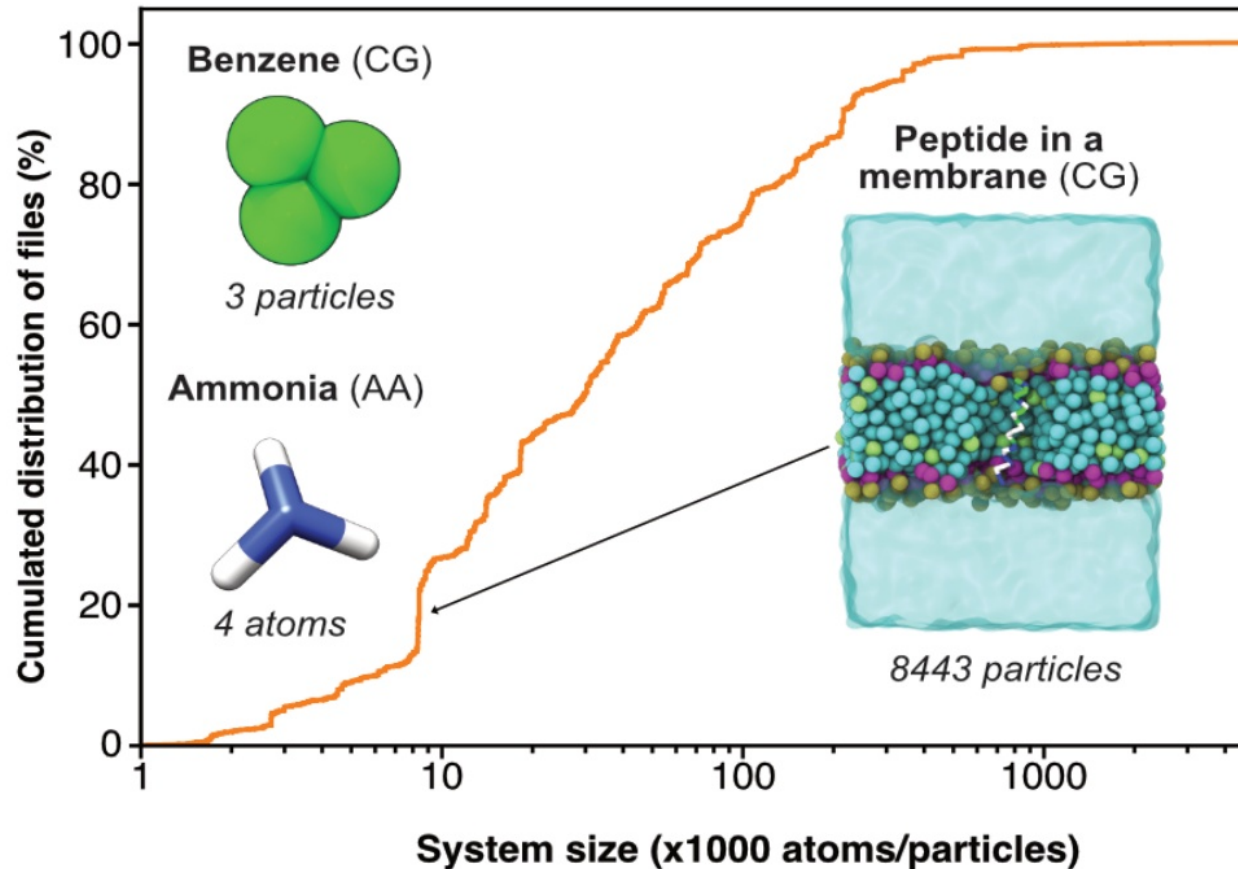## Simulations of a pulmonary surfactant monolayer with additional compounds

To be added.

Files (47.1 GB)

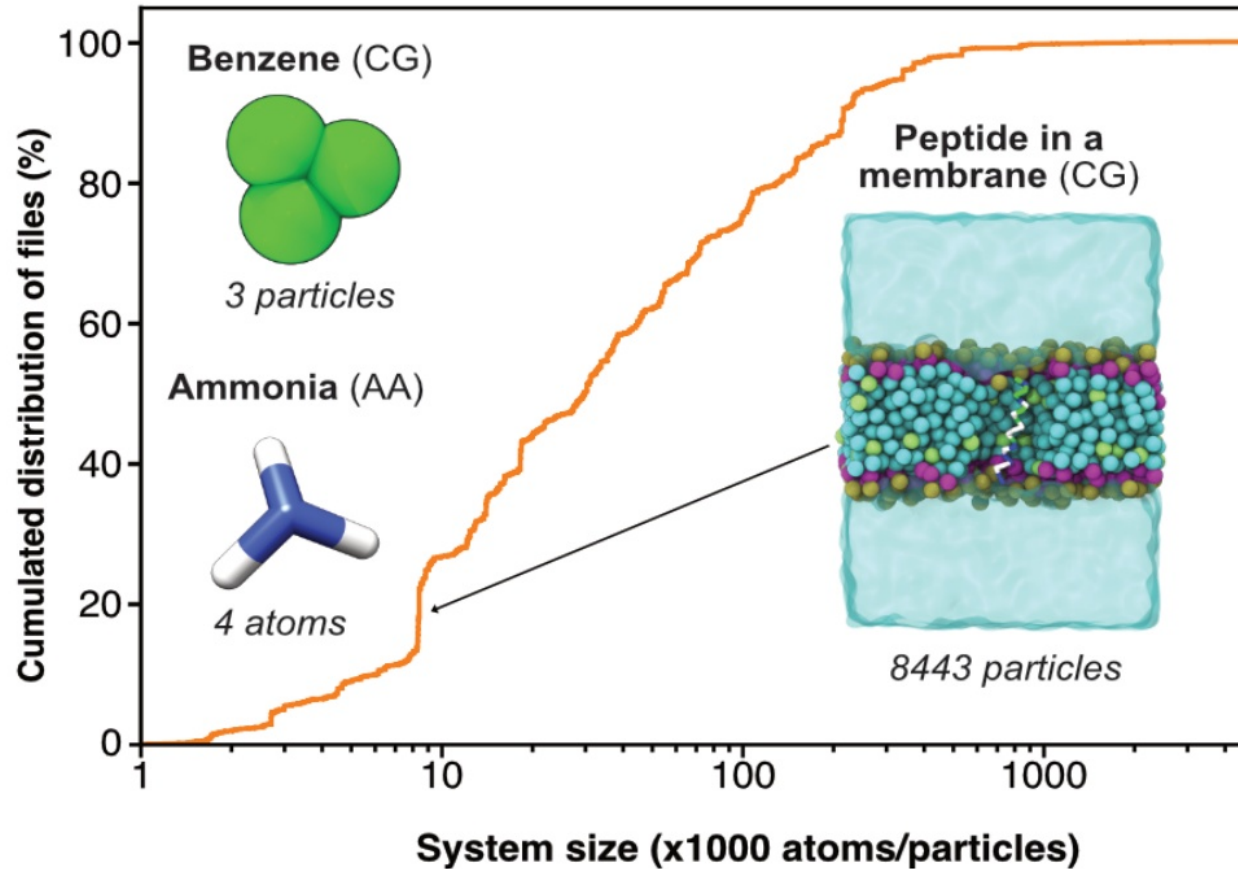| Name | |
|------|---|
| benzaldehyde.ndx | |
| md5:987f07a7e7fab150db9713558d7dd8ad | |
| benzaldehyde.top | |
| md5:e90dbc7f37ff3ea109ca4aba127cd3ca | |

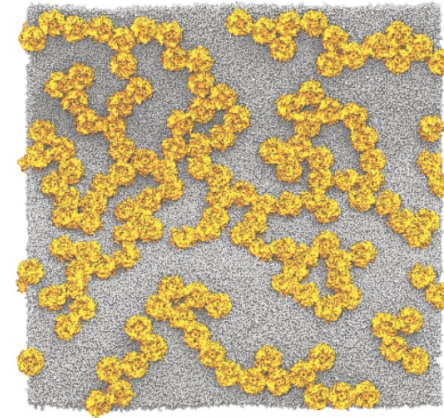Open the can and taste the soup
## Open the file and extract metadata

# Step 2: Metadata from Gromacs (.gro) files
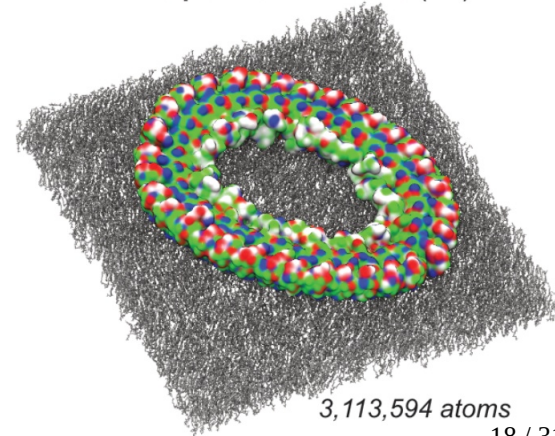
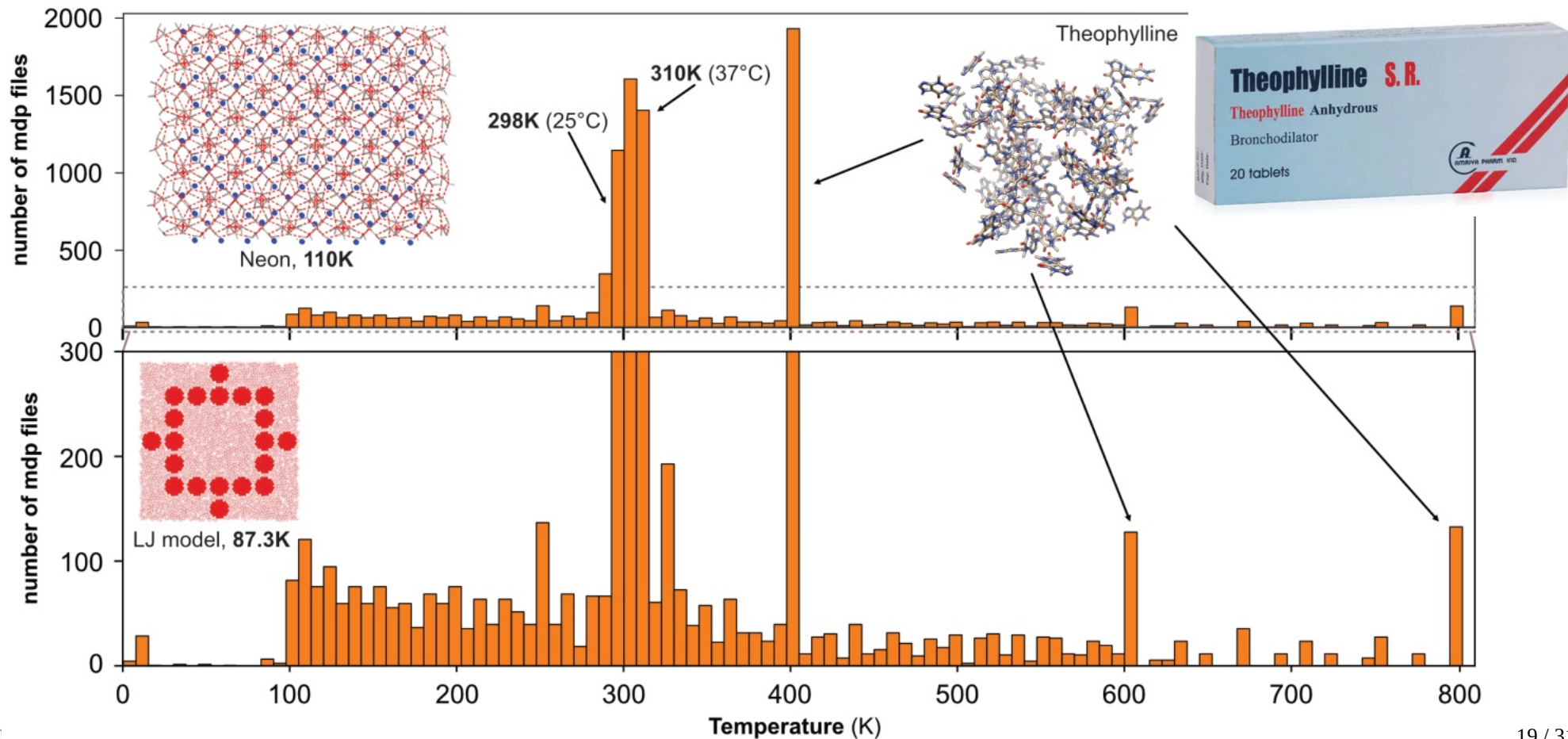# Step 2: Metadata from Gromacs (.gro) files

# Step 2: Metadata from Gromacs (.mdp) files

# Step 2: Explore metadata 🔍



https://mdverse.streamlit.app/

# Step 3: Back to the FAIR principles

**F**indable
**A**ccessible
**I**nteroperable
**R**eusable

# Step 3: Back to the FAIR principles

**F**indable
**A**ccessible
**I**nteroperable
**R**eusable

*"[…] the FAIR Principles put specific emphasis on enhancing the ability of **machines to automatically** find and **use the data**, in addition to supporting its reuse by individuals."*

Wilkinson et al., Scientific Data, 2016.

# Step 3: Assess accessibility (zip files)

## All-atom molecular dynamics simulations of SARS-CoV-2 envelope protein E

Kuzmin Alexander; Orekhov Philipp; Astashkin Roman; Gordeliy Valentin; Gushchin Ivan

The trajectories of all-atom (AA) MD simulations (NoPTM-1;2;3;4_POPC;Mix_CHARMM36m: 0.1x3 µs) were obtained based on 4 starting representative conformations from the coarse-grained simulation (10.5281/zenodo.4740706). For each starting structure, there are six trajectories of the E protein: 3 with the protein embedded in the membrane containing POPC, and 3 with the membrane mimicking the natural ERGIC membrane (Mix: 50% POPC, 25% POPE, 10% POPI, 5% POPS, 10% cholesterol).

Simulations have been performed using the CHARMM36m (AA) force field, running with the GROMACS 2019.5 package on the supercomputer JURECA at Forschungszentrum Jülich under the conditions reported in bioRxiv 2021.03.10.434722.

https://doi.org/10.1002/prot.26317

Preview ›

Files (45.6 GB) ⌄

| Name | Size | | |
|------|------|------|------|
| NoPTM-1_Mix_CHARMM36m_0.1x3mks.zip | 6.0 GB | 👁 Preview | ⬇ Download |
| md5:3f3854a3de4a10489e9895b8ba5d368b ❓ | | | |
| NoPTM-1_POPC_CHARMM36m_0.1x3mks.zip | 5.1 GB | 👁 Preview | ⬇ Download |
| md5:e9ef288e489689c70956952091bba5bf ❓ | | | |

# Step 3: Assess accessibility (zip files)

## All-atom molecular dynamics simulations of SARS-CoV-2 envelope protein E

Kuzmin Alexander; Orekhov Philipp; Astashkin Roman; Gordeliy Valentin; Gushchin Ivan

The trajectories of all-atom (AA) MD simulations (NoPTM-1;2;3;4_POPC;Mix_CHARMM36m: 0.1x3 μs) were obtained based on 4 starting representative conformations from the coarse-grained simulation (10.5281/zenodo.4740706). For each starting structure, there are six trajectories of the E protein: 3 with the protein embedded in the membrane containing POPC, and 3 with the membrane mimicking the natural ERGIC membrane (Mix: 50% POPC, 25% POPE, 10% POPI, 5% POPS, 10% cholesterol).

Simulations have been performed using the CHARMM36m (AA) force field, running with the GROMACS 2019.5 package on the supercomputer JURECA at Forschungszentrum Jülich under the conditions reported in bioRxiv 2021.03.10.434722.

https://doi.org/10.1002/prot.26317

Preview                                                    >

Files (45.6 GB)                                            ⌄

| Name | Size | | |
|------|------|--|--|
| NoPTM-1_Mix_CHARMM36m_0.1x3mks.zip | 6.0 GB | 👁 Preview | ⬇ Download |
| md5:3f3854a3de4a10489e9895b8ba5d368b ❓ | | | |
| NoPTM-1_POPC_CHARMM36m_0.1x3mks.zip | 5.1 GB | 👁 Preview | ⬇ Download |
| md5:e9ef288e489689c70956952091bba5bf ❓ | | | |

**88 %** of indexed files were in zip files

| | |
|--|--|
| 📄 NoPTM-2_Mix_CHARMM36m_0.1x3mks.zip | |
| 📄 NoPTM-2-1_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2-2_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2-3_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2_Mix_CHARMM36m.pdb | 11.8 MB |
| 📄 NoPTM-2_Mix_CHARMM36m.tpr | 4.1 MB |

# Step 3: Assess accessibility (zip files)

## All-atom molecular dynamics simulations of SARS-CoV-2 envelope protein E

Kuzmin Alexander; Orekhov Philipp; Astashkin Roman; Gordeliy Valentin; Gushchin Ivan

The trajectories of all-atom (AA) MD simulations (NoPTM-1;2;3;4_POPC;Mix_CHARMM36m: 0.1x3 μs) were obtained based on 4 starting representative conformations from the coarse-grained simulation (10.5281/zenodo.4740706). For each starting structure, there are six trajectories of the E protein: 3 with the protein embedded in the membrane containing POPC, and 3 with the membrane mimicking the natural ERGIC membrane (Mix: 50% POPC, 25% POPE, 10% POPI, 5% POPS, 10% cholesterol).

Simulations have been performed using the CHARMM36m (AA) force field, running with the GROMACS 2019.5 package on the supercomputer JURECA at Forschungszentrum Jülich under the conditions reported in bioRxiv 2021.03.10.434722.

https://doi.org/10.1002/prot.26317

Preview >

**Files (45.6 GB)** ∨

| Name | Size | | |
|------|------|---|---|
| NoPTM-1_Mix_CHARMM36m_0.1x3mks.zip | 6.0 GB | 👁 Preview | ⬇ Download |
| md5:3f3854a3de4a10489e9895b8ba5d368b ❓ | | | |
| NoPTM-1_POPC_CHARMM36m_0.1x3mks.zip | 5.1 GB | 👁 Preview | ⬇ Download |
| md5:e9ef288e489689c70956952091bba5bf ❓ | | | |

**88 %** of indexed files were in zip files

📄 NoPTM-2_Mix_CHARMM36m_0.1x3mks.zip

| | |
|---|---|
| 📄 NoPTM-2-1_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2-2_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2-3_Mix_CHARMM36m_0.1mks.xtc | 1.9 GB |
| 📄 NoPTM-2_Mix_CHARMM36m.pdb | 11.8 MB |
| 📄 NoPTM-2_Mix_CHARMM36m.tpr | 4.1 MB |

# Step 3: Assess reusability

A Gromacs trajectory file could be reused to:

- Analyze a simulation (.xtc + .pdb/.gro/.tpr)

- Continue a simulation (.gro/.trr/.cpt + .mdp/.tpr)

# Step 3: Assess reusability

A Gromacs trajectory file could be reused to:

- Analyze a simulation (.xtc + .pdb/.gro/.tpr)

- Continue a simulation (.gro/.trr/.cpt + .mdp/.tpr)

Do we have any proof we can actually reuse the data?

# MD data sharing highlights

# MD data sharing highlights

Depositing MD simulations data in a FAIR-enabled repository does **not guarantee** your data is actually FAIR

# MD data sharing highlights

Depositing MD simulations data in a FAIR-enabled repository does **not guarantee** your data is actually FAIR

- Provide metadata (context)

# MD data sharing highlights

Depositing MD simulations data in a FAIR-enabled repository does **not guarantee** your data is actually FAIR

- Provide metadata (context)

- Avoid zip files (no .tgz)

# MD data sharing highlights

Depositing MD simulations data in a FAIR-enabled repository does **not guarantee** your data is actually FAIR

- Provide metadata (context)

- Avoid zip files (no .tgz)

Why is it important?

- Sharing and storing data cost **money** and **energy**

- Dynamic generative deep-learning models?
  (Need for high quality, curated data)

# What's next? Keep digging!

# What's next?                    Keep digging!

Explore other "data" repositories:

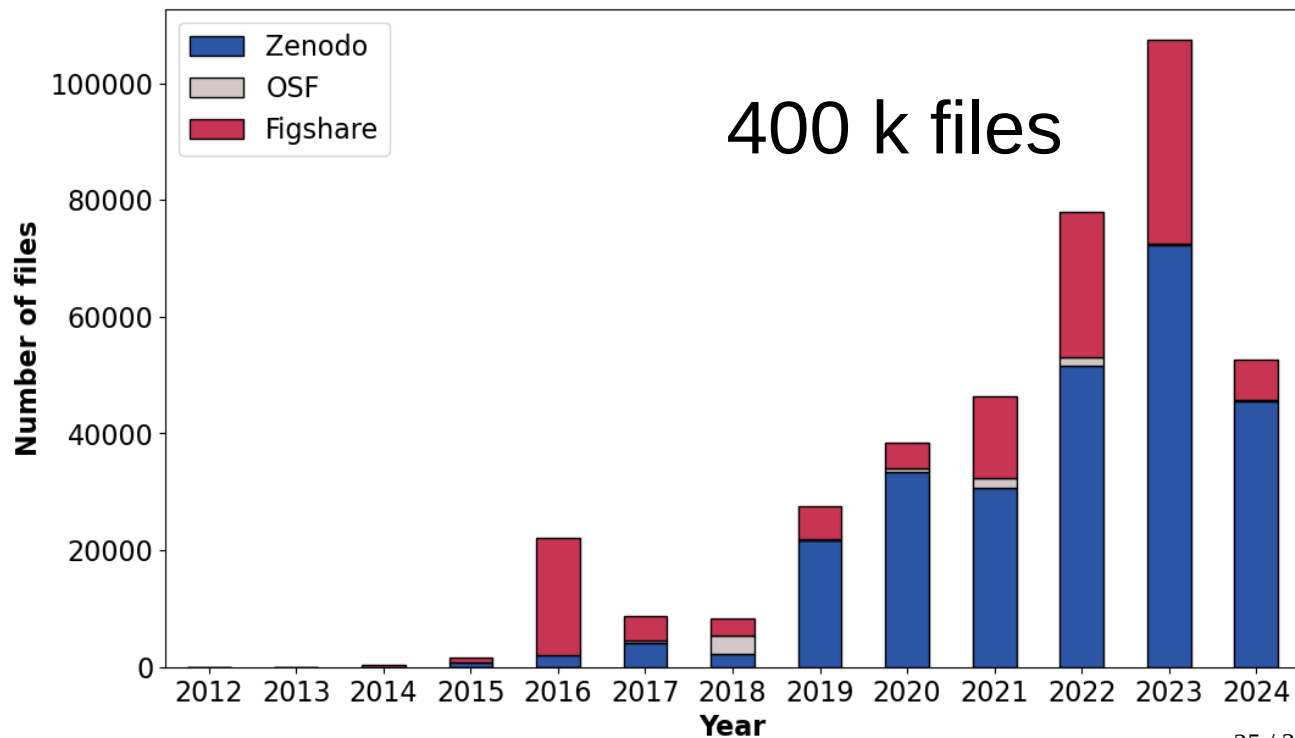- Dryad, Science Data Bank

- ATLAS, MDDB

# What's next?                    Keep digging!

Explore other "data" repositories:

- Dryad, Science Data Bank

- ATLAS, MDDB



400 k files

# What's next?

Extract structured metadata from **raw text**

# What's next?

Extract structured metadata from **raw text**

Dataset   🔓 Open

## Simulations of DLPC, DOPC, and DEPC membranes

Matti Javanainen[1] 🆔

Show affiliations

Simulations of three different lipid bilayers using the CHARMM36 lipid model. The bilayers were built using CHARMM-GUI and simulated for 500 ns using Gromacs 2019. The systems contain 200 lipids (DLPC, DOPC, or DEPC), 50 lipids per water, and 130 mM NaCl. For each simulation, the run input file (tpr), topology (top), index (ndx), trajectory (xtc), log file (log), energy file (edr), continue point (cpt), and the final structure (gro) are given. The topology refers to the Gromacs-compatible CHARMM36 force field, available at http://mackerell.umaryland.edu/charmm_ff.shtml#gromacs and to other itp files provided in this upload. The simulation parameters are given in the md.mdp file.

# What's next?

Extract structured metadata from **raw text**

Text mining / Named Entity Recognition



Simulations of DLPC SELECTED , DOPC MOL , and DEPC MOL membranes

Simulations of three different lipid MOL bilayers using the CHARMM36 FFM lipid MOL model. The bilayers were built using

CHARMM-GUI SOFT and simulated for 500 ns STIME using Gromacs 2019 SOFT . The systems contain 200 lipids MOL ( DLPC

MOL , DOPC MOL , or DEPC MOL ), 50 lipids MOL per water MOL , and 130 mM NaCl MOL . For each simulation, the run

input file (tpr), topology (top), index (ndx), trajectory (xtc), log file (log), energy file (edr), continue point (cpt), and the final structure (gro) are

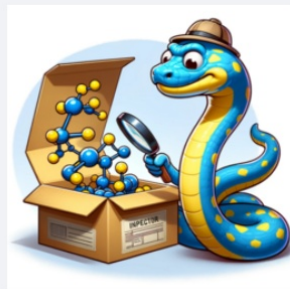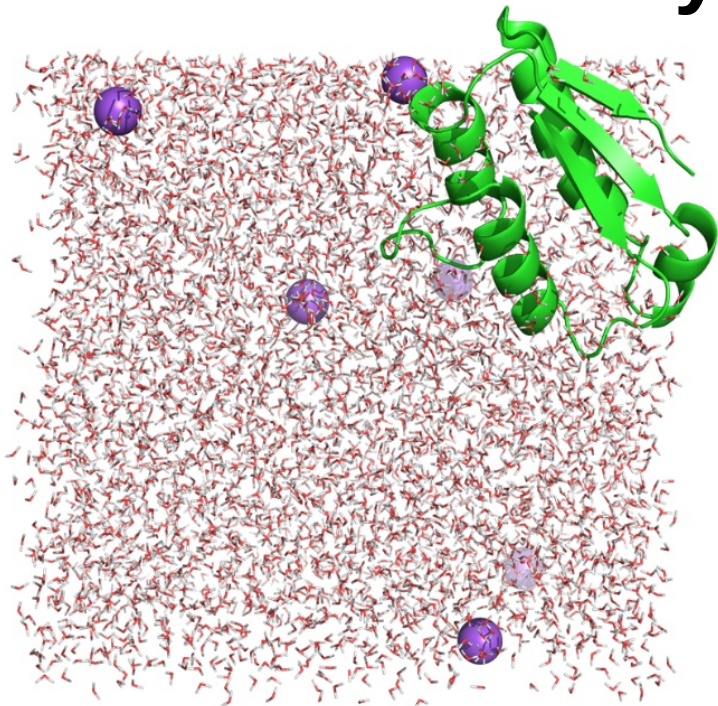given. The topology refers to the Gromacs SOFT -compatible CHARMM36 FFM force field, available at https://removed

ff.shtml#gromacs and to other itp files provided in this upload. The simulation parameters are given in the md.mdp file.

Mohamed Oussaren

# What's next?

Extract structured metadata
from **molecular systems**



https://grodecoder.streamlit.app/

Karine Duong

# What's next?

Improve discoverability and data(sets) exploration



*Linked User-driven Multidisciplinary Exploration Network*
2025 → 2027

Involved communities: SSH, Maths, Earth System, Molecular Dynamics

# Thanks 🙏

**IJM, Paris, France**
Lisa Bouarroudj
Mohamed Oussaren

**CBI, Toulouse, France**
Magdalena Szczuka
*Matthieu Chavent*

**LBT, Paris, France**
Marc Baaden
Karine Duong

**Amsterdam, Netherlands**
Steven Garcia

**Univ. Copenhagen, Denmark**
*Johanna K. S. Tiemann*
Kresten Lindorff-Larsen

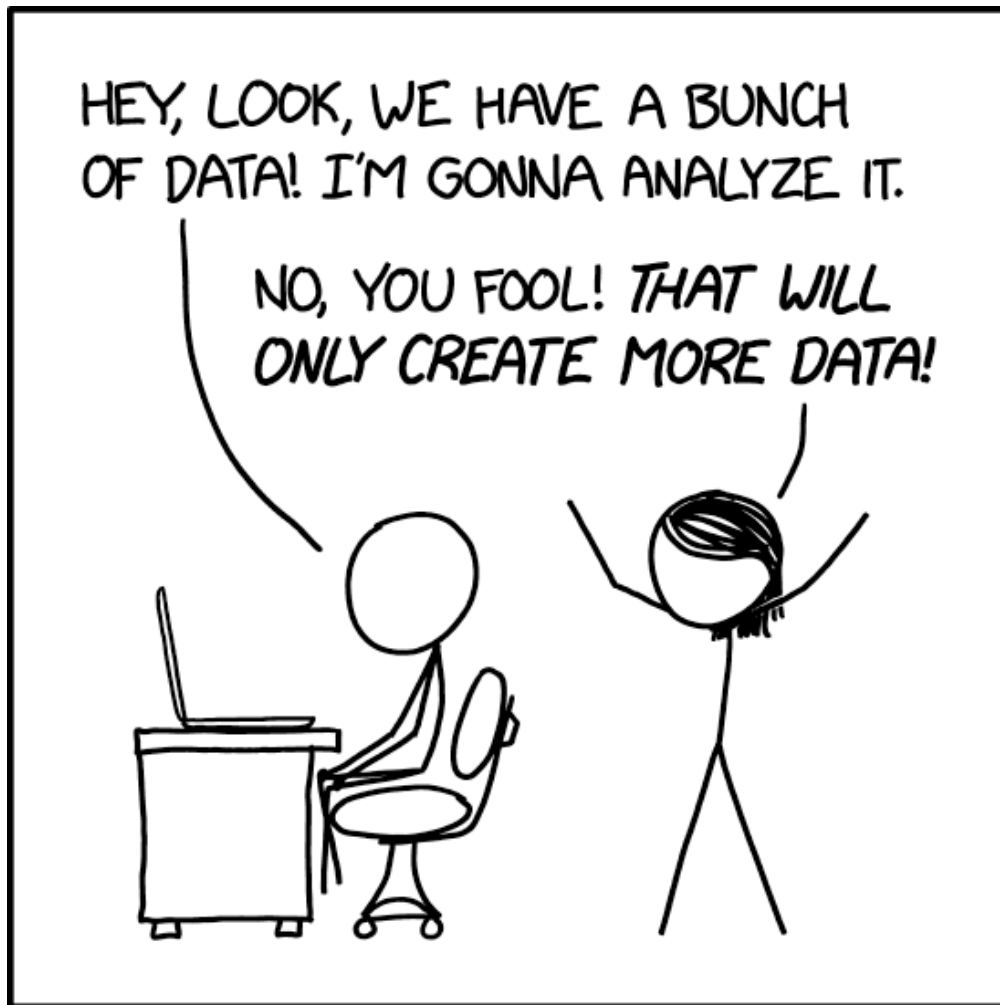**KTH Royal Inst. Tech. Stockholm, Sweden**
Lucie Delemotte
Erik Lindahl

**Stockholm Univ. , Sweden**
Rebecca J. Howard

# Questions? ✨



▲ This presentation ▲

DOI 10.5281/zenodo.14264571



HEY, LOOK, WE HAVE A BUNCH OF DATA! I'M GONNA ANALYZE IT.

NO, YOU FOOL! *THAT WILL ONLY CREATE MORE DATA!*