



# PSDI

PHYSICAL SCIENCES  
DATA INFRASTRUCTURE

**PSDI – Pathfinder 4**

FAIR Data for the Biomolecular Simulation Community

James Gebbie-Rayet & Jas Kalayan

# Webinar Structure

- ▶ Overview of domain (James)
- ▶ Aims of the project (James)
- ▶ Technical overview of platform so far (Jas)
- ▶ Technical demonstration (Jas)
- ▶ Questions (You)

# Biologists are brilliant with data!

Other areas of biology are very organised about data!



<https://www.rcsb.org/>



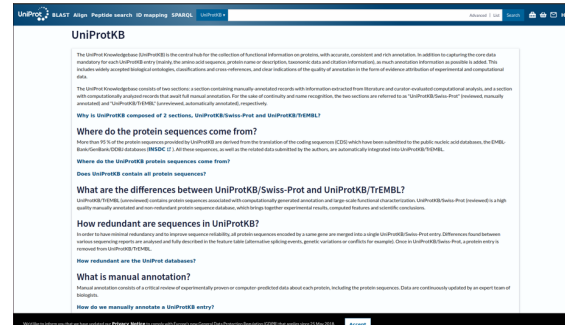
<https://www.ebi.ac.uk/emdb/>



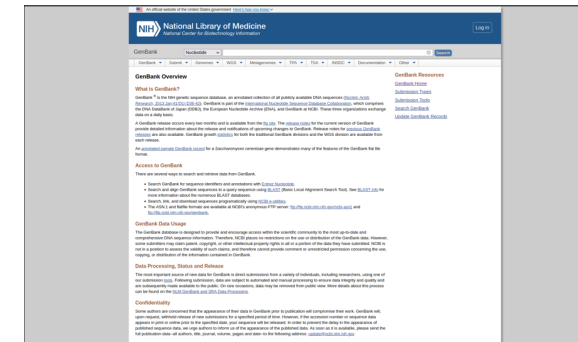
<https://gpcrdb.org/>



<https://www.ebi.ac.uk/empair/>



<https://www.uniprot.org/help/uniprotkb>



<https://www.ncbi.nlm.nih.gov/genbank/>

# What about Biomolecular Simulation?

Not much in the way of established production databases or services in UK!!





# Existing Generic Sharing Platforms

Some existing ways we as researchers share our data:

The Zenodo logo is the word "zenodo" in a bold, lowercase, sans-serif font.

<https://zenodo.org/>

The Figshare logo consists of a circular arrangement of small, multi-colored dots (red, green, blue, yellow, grey) forming a ring, followed by the word "figshare" in a lowercase, sans-serif font.

<https://figshare.com/>

- ▶ Great for sharing files
- ▶ Great for getting persistent identifiers
- ▶ No requirement to share data that is FAIR
- ▶ Researchers left to own devices on what is shared!



Open Science Framework

<https://osf.io/>

The Mendeley logo features a red stylized icon of three interconnected shapes above the word "Mendeley" in a red, lowercase, sans-serif font.

<https://www.mendeley.com/>

# What's the Main Issues in Domain?

The main issues that our pathfinder aims to tackle are:

- ▶ No consistent approach to storing or sharing simulation data across the community.
- ▶ No real infrastructure or route to make simulation data accessible.
- ▶ Usage of existing or emerging platforms not FAIR
- ▶ The funding, research, publish cycle in the field currently discourages researchers to think or do anything about sharing data.
- ▶ Research papers in biomolecular simulation do not often contain information to allow fully reproducible studies.
- ▶ Even publications where simulations are well described, are difficult to reproduce due to full provenance of model creation not being present.
- ▶ Difficult to know exactly which experimental data sets/sources involved in research studies.

# What is PSDI Pathfinder Doing?

Objective is to establish data infrastructure prototype and tools to improve data practices in field without major cultural shifts:

- ▶ Main aim is to improve data practices in domain – align with FAIR principles
- ▶ Prototype tools to capture full data provenance for model creation, simulation and analytics (FAIR)
- ▶ Prototype infrastructure tools to store, access, find and share data (FAIR)
- ▶ Establish long term collaboration with other data initiatives (EBI, EU and US funded)
- ▶ Establish hard data links to experimental data sources
- ▶ “I” (FAIR) not yet in scope of this pathfinder (excellent projects in wider community)



# We want data to be easy!

Our philosophy is that data practice should be baked into the way we work, if it adds extra burden then we won't do it!

- ▶ We want to eliminate having to fill in web forms with huge numbers of fields
- ▶ Our tooling is designed to harvest data and metadata in an automated way
- ▶ We do this in a way that avoids you having to change much about how you work
- ▶ You should be in control of your data at all times and choose if and when to share
- ▶ It should be easy to generate overviews of your data to make publishing simpler
- ▶ Much more is under development!

# What could it enable?

We think that better data practices unlock enormous potential!

- ▶ As a research field we can avoid duplication of effort
- ▶ Having our full data out in the open will enable exploitation by other domains (AI/ML)
- ▶ Collecting and sharing full model provenance will enable full reproducibility
- ▶ Sharing simulation data can improve statistical/physical understanding of systems when studied with different codes/methods (we saw this during SARS-COV2)
- ▶ Sharing simulation data is climate friendly – simulations on HPC are energy intensive!
- ▶ We can use high quality datasets to inform novel method development and evaluate the quality of existing and emerging forcefield, physics etc
- ▶ Knowledge will bleed between research groups since model creation can be shared
- ▶ Councils and other funders will be able to see the true cost of research, by linking with experimental databases and collecting contributors in the chain

Over to Jas!



# PSDI Webinar: FAIR Data for the Biomolecular Simulation Community

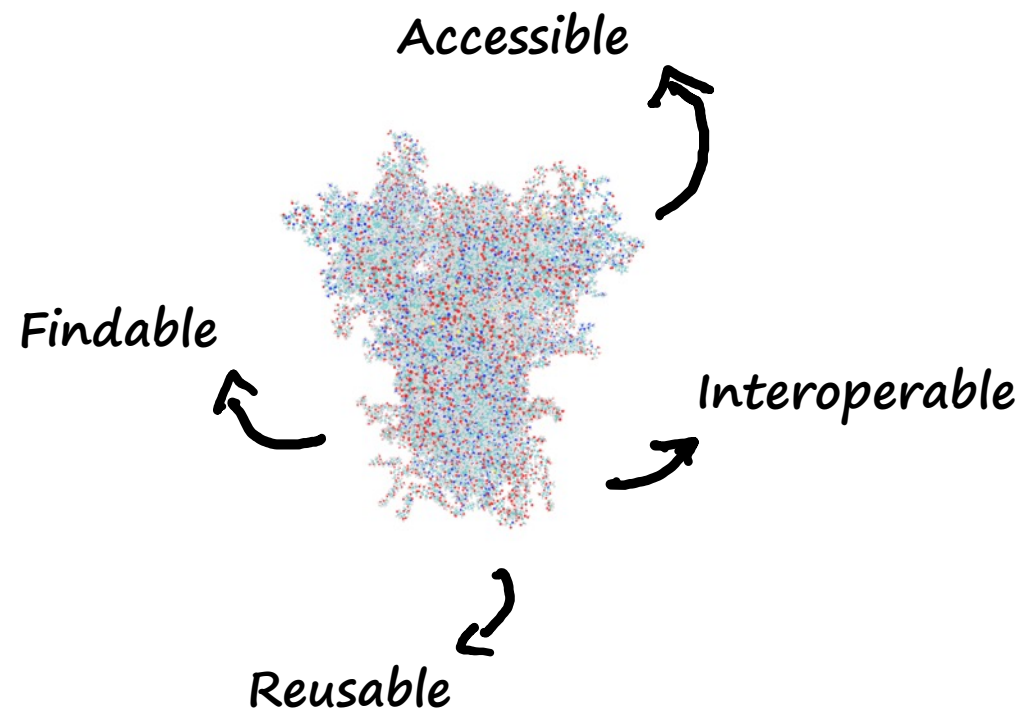
Acknowledgments to:

Kin Chao,  
Joel Greer, Tom Burnley,  
Martyn Winn

Imperial College  
London



18.10.2023



\*DESRES COVID19 spike protein



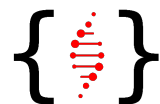
**PSDI**  
PHYSICAL SCIENCES  
DATA INFRASTRUCTURE



Science and  
Technology  
Facilities Council

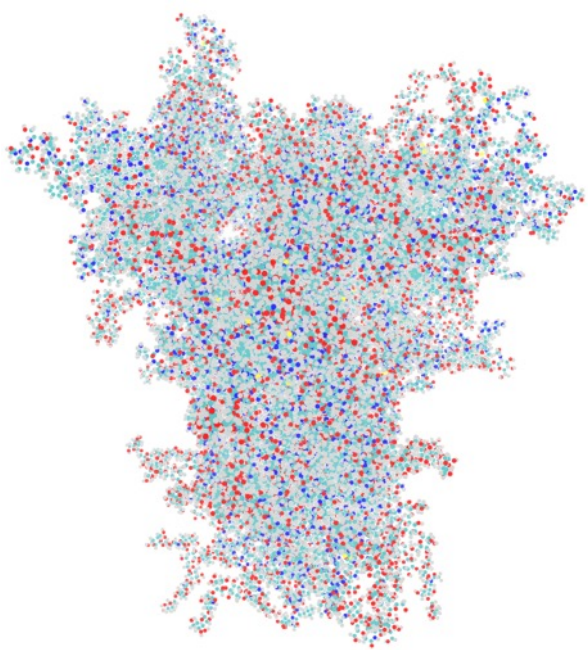


Engineering and  
Physical Sciences  
Research Council



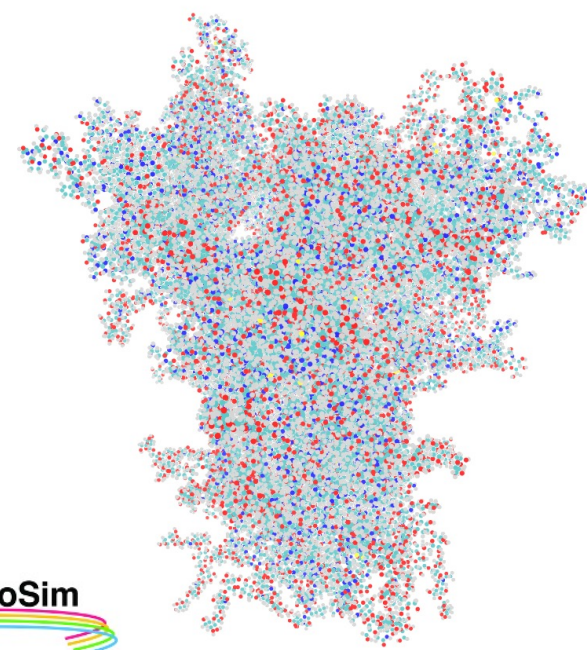
**CCP**BioSim

# The Biomolecular Simulation Community



## Data we produce?

xyz coordinates for each atom  
 $\approx 10^5$  atoms,  
 $\approx 10^3$ - $10^4$  frames,  
 $\approx$  weeks to perform  
 $\approx$  TBs of data per simulation



Start

*...various  
protocols...*

Simulation

# Various Protocols in MD Simulation

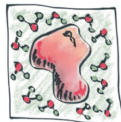
1. Get crystal Structure



2. System Preparation



3. Parameterisation



PACKMOL

Initial configurations for Molecular Dynamics Simulations by packing optimization



CHARMM-GUI

Effective Simulation Input Generator and More

GROMACS  
FAST. FLEXIBLE. FREE.



APBS & PDB2PQR

Software for biomolecular electrostatics and solvation



6. Simulation



5. Equilibration



4. Minimisation

Using the same MD engine

# Biomolecular Simulation Data Provenance using AiiDA

*Mainly used in computational materials science*



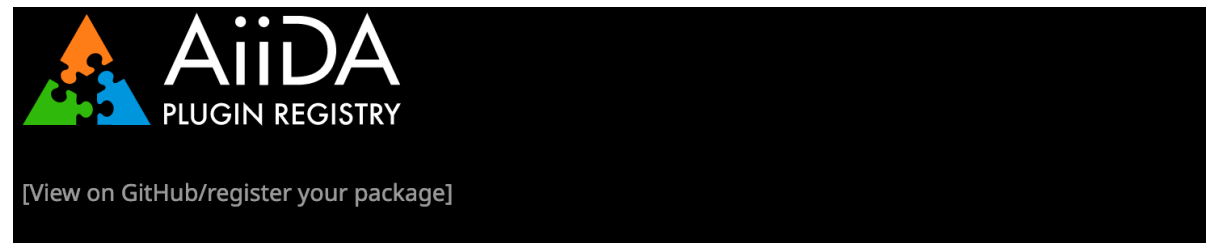
**MATERIALSCLOUD**



AiiDA is a Python infrastructure that helps track complex workflows used in computational science.  
<https://github.com/aiidateam/aiida-core>

# AiiDA-GROMACS: A plugin for Data Provenance with GROMACS Simulations

*GROMACS is used by 70% of HECBioSim users, But more plugins to come!*



## AiiDA plugin package "aiida-gromacs"

[< back to the registry index](#)

### General information

Current state: status alpha

Short description: A plugin for using GROMACS with AiiDA for molecular dynamics simulations.

How to install: `pip install git+https://github.com/jimboid/aiida-gromacs`

Source code: [Go to the source code repository](#)

Documentation: [Go to plugin documentation](#)

### Detailed information

Author(s): James Gebbie-Rayet

Contact: [james.gebbie@stfc.ac.uk](mailto:james.gebbie@stfc.ac.uk)

How to use from python: `import aiida_gromacs`

Most recent version:

Compatibility: AiiDA >=2.0,<3

### Plugins provided by the package

Calculations 7 Parsers 7 Data 6 Workflows 1



# Basics of using AiiDA for Data Provenance

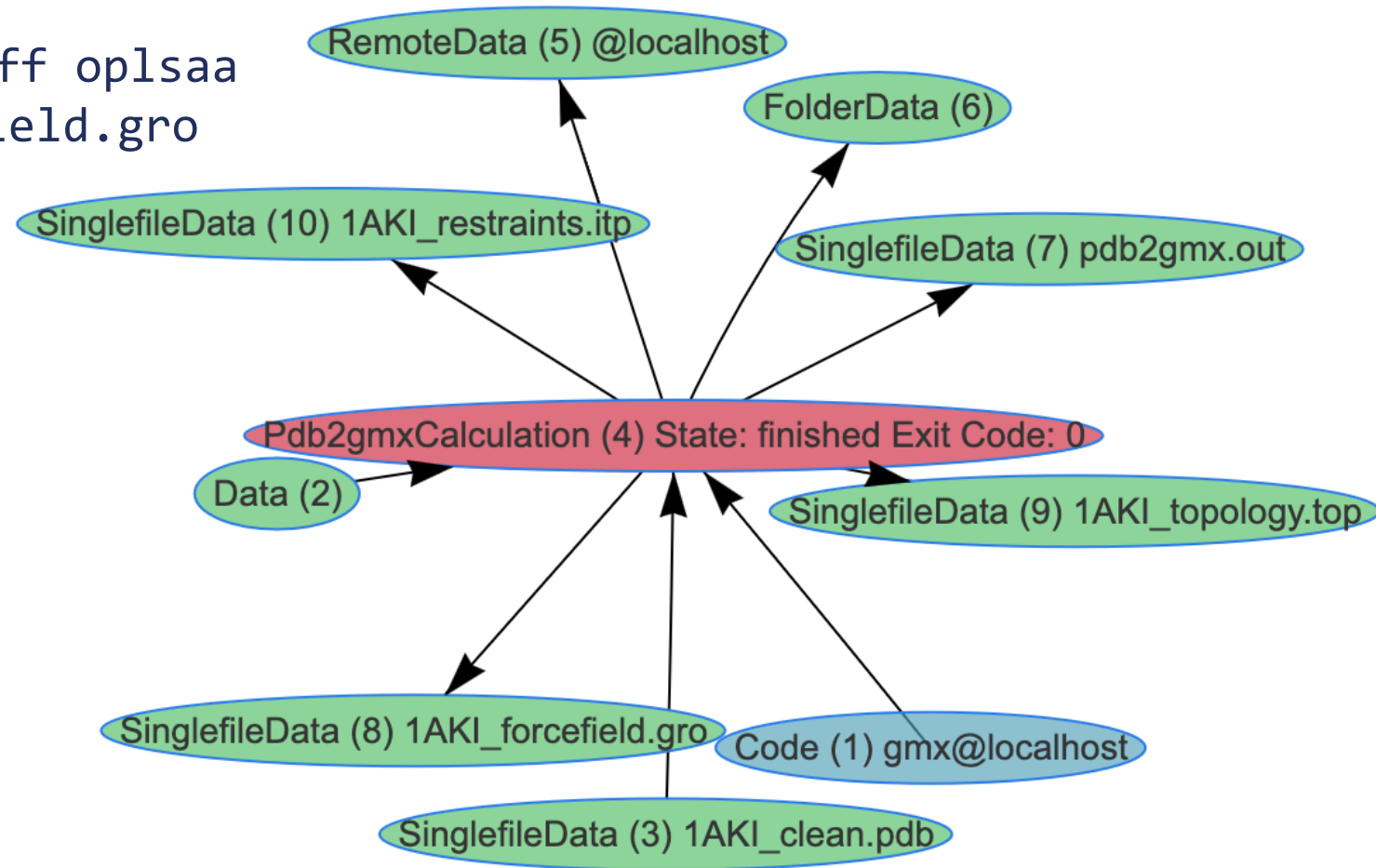


1. *What is the command you want to run?*
2. *Where do you want to run it?*
3. *What are your inputs?*
4. *What are your outputs?*

# Example of Data Provenance with aiida-gromacs

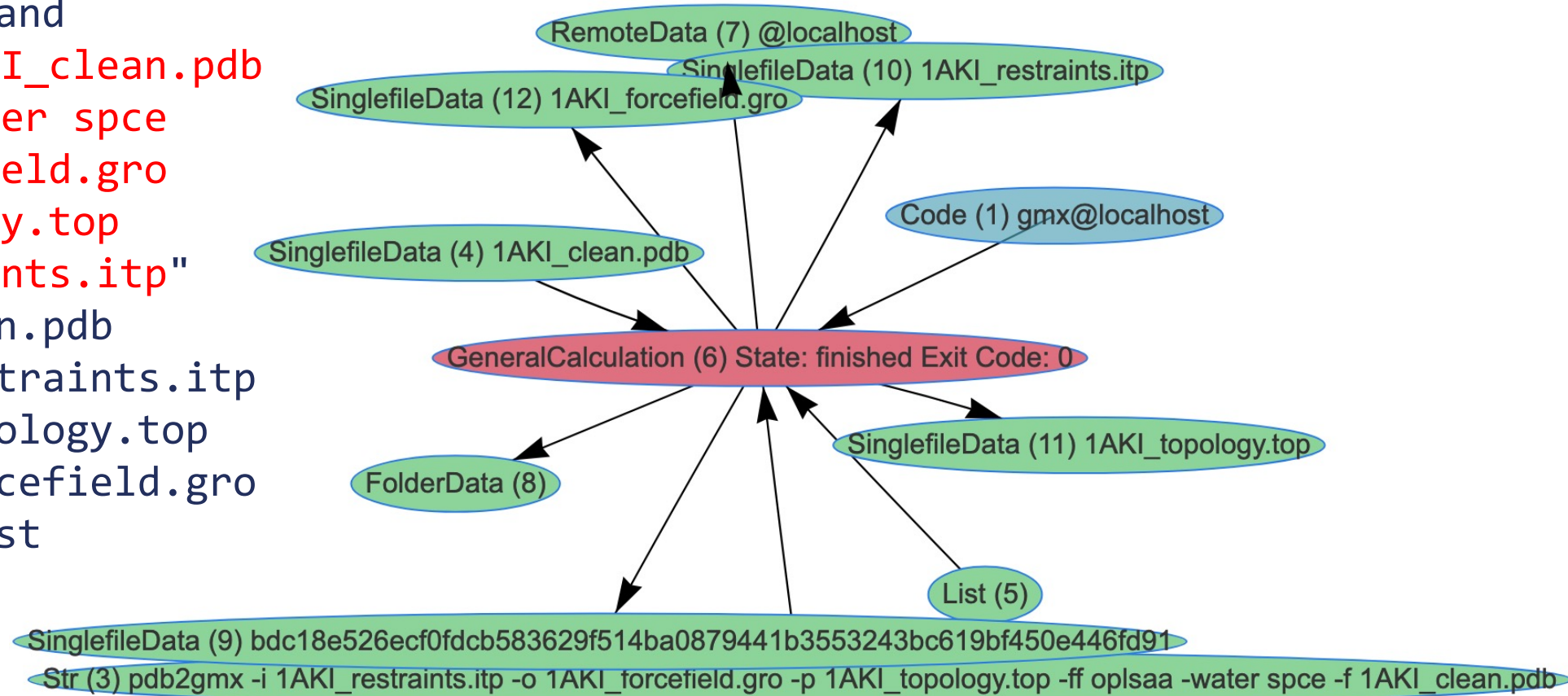
```
$ gm_x_pdb2gmx -f 1AKI_clean.pdb -ff oplsa  
-water spce -o 1AKI_forcefield.gro  
-p 1AKI_topology.top  
-i 1AKI_restraints.itp
```

*Note the underscore!*



# Example of Data Provenance with aiida-gromacs

```
$ genericMD --command  
  "pdb2gmx -f 1AKI_clean.pdb  
  -ff oplsaa -water spce  
  -o 1AKI_forcefield.gro  
  -p 1AKI_topology.top  
  -i 1AKI_restraints.itp"  
--inputs 1AKI_clean.pdb  
--outputs 1AKI_restraints.itp  
--outputs 1AKI_topology.top  
--outputs 1AKI_forcefield.gro  
--code gmx@localhost
```



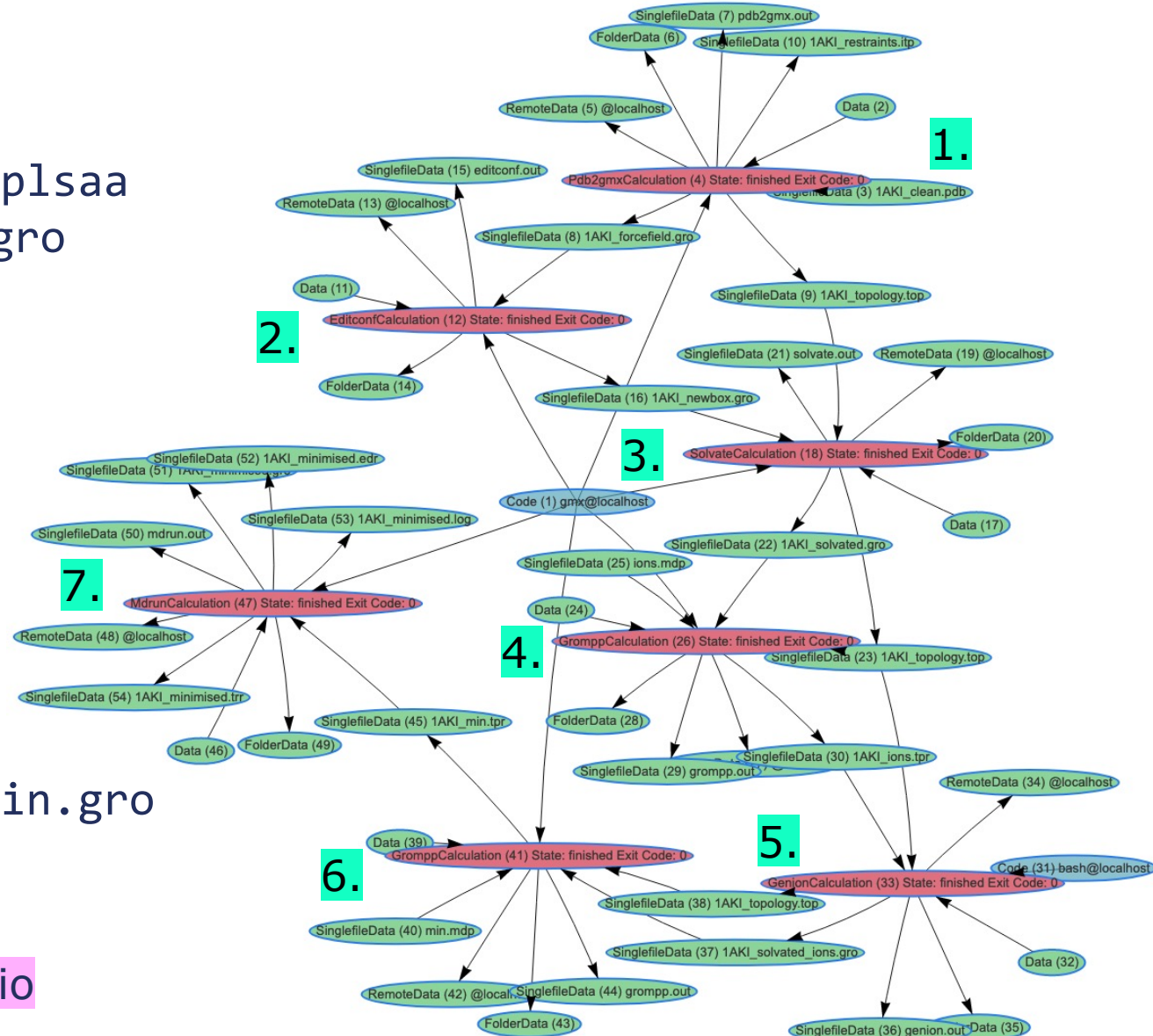
# Example of Data Provenance with aiida-gromacs

1. `gmx_pdb2gmx -f 1AKI_clean.pdb -ff oplsa -water spce -o 1AKI_forcefield.gro -p 1AKI_topology.top -i 1AKI_restraints.itp`

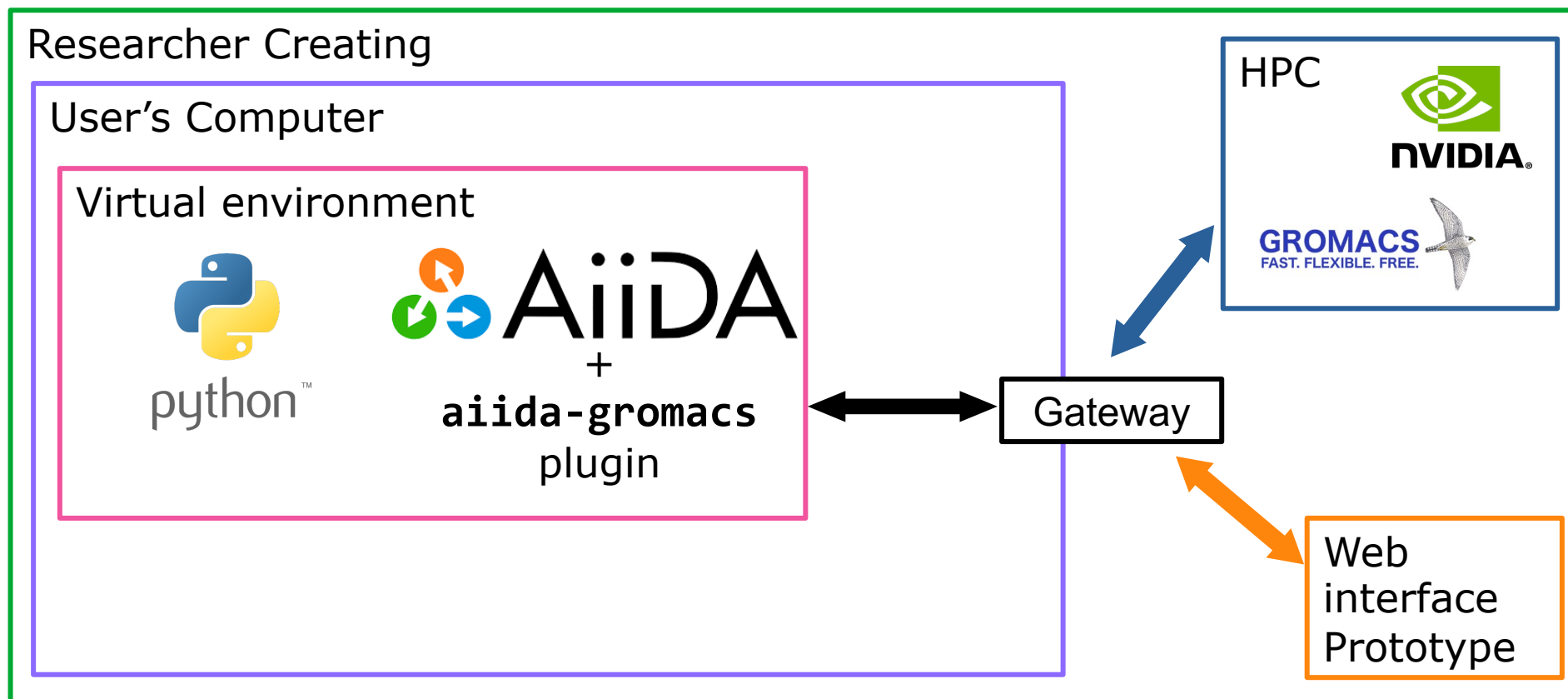
2. `gmx_editconf,`  
3. `gmx_solvate,`  
4. `gmx_grompp,`  
5. `gmx_genion,`  
6. `gmx_grompp,`

7. `gmx_mdrun -s 1AKI_min.tpr -c 1AKI_min.gro -e 1AKI_min.edr -g 1AKI_min.log -o 1AKI_min.trr`

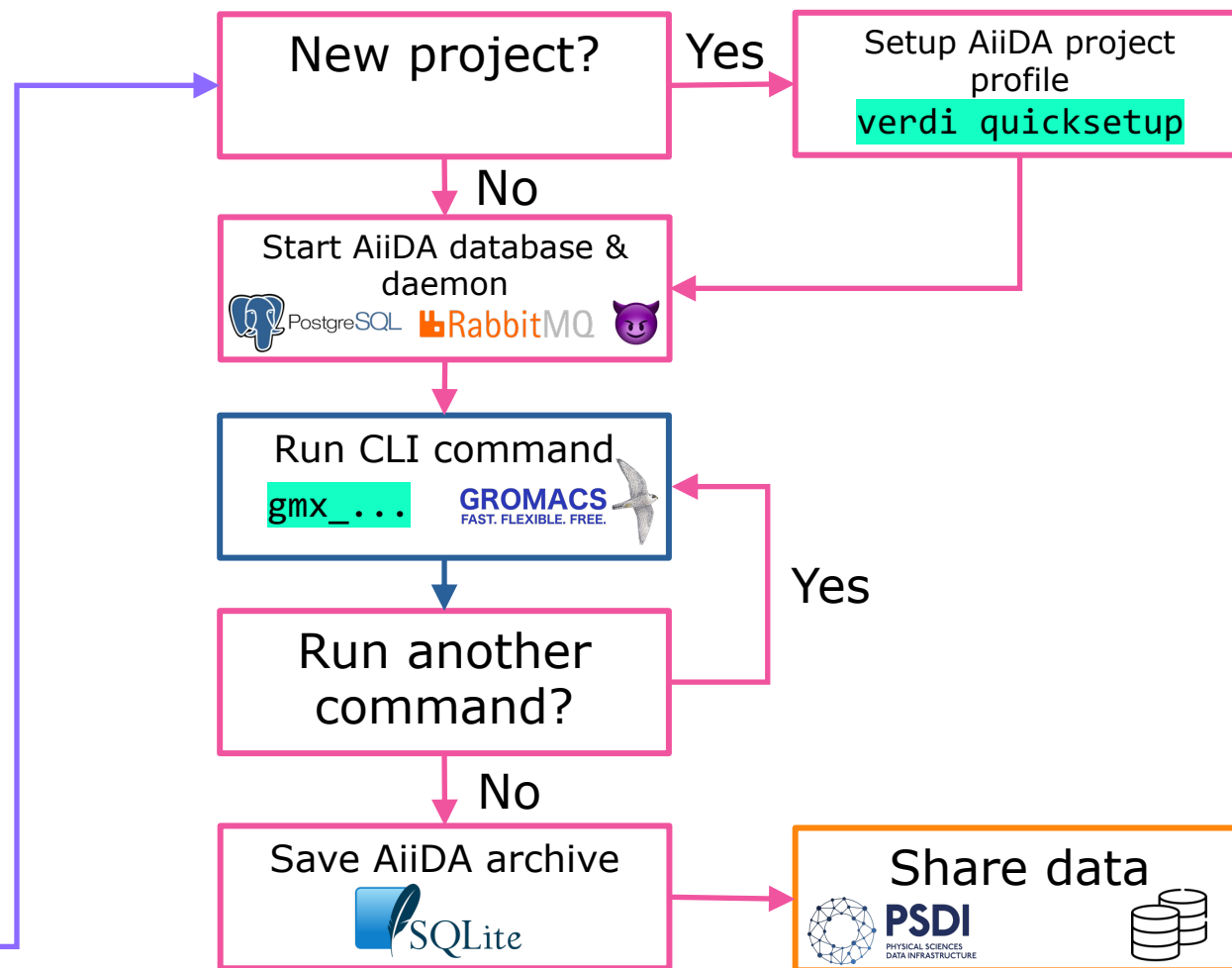
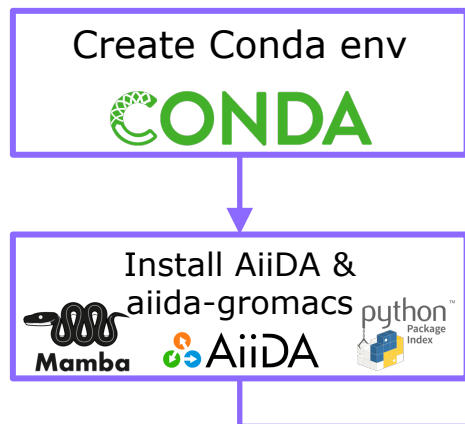
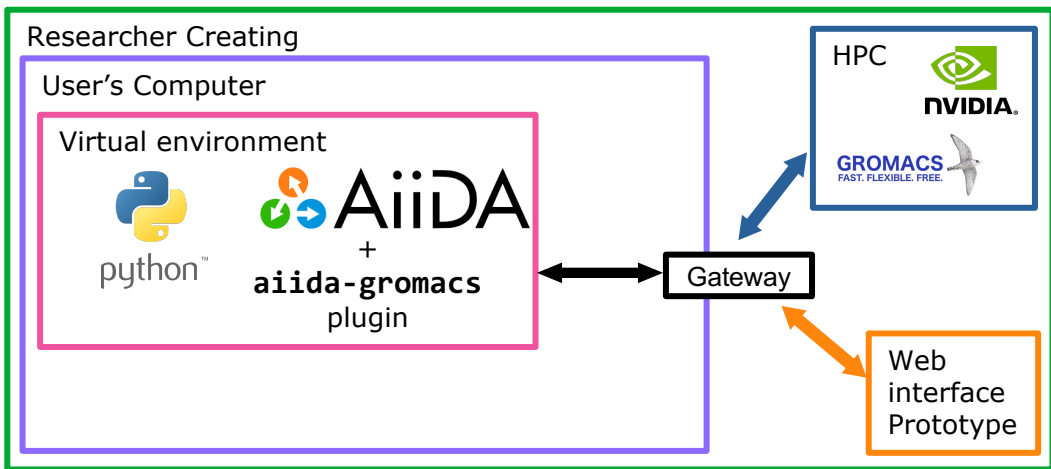
<https://aiida-gromacs.readthedocs.io>



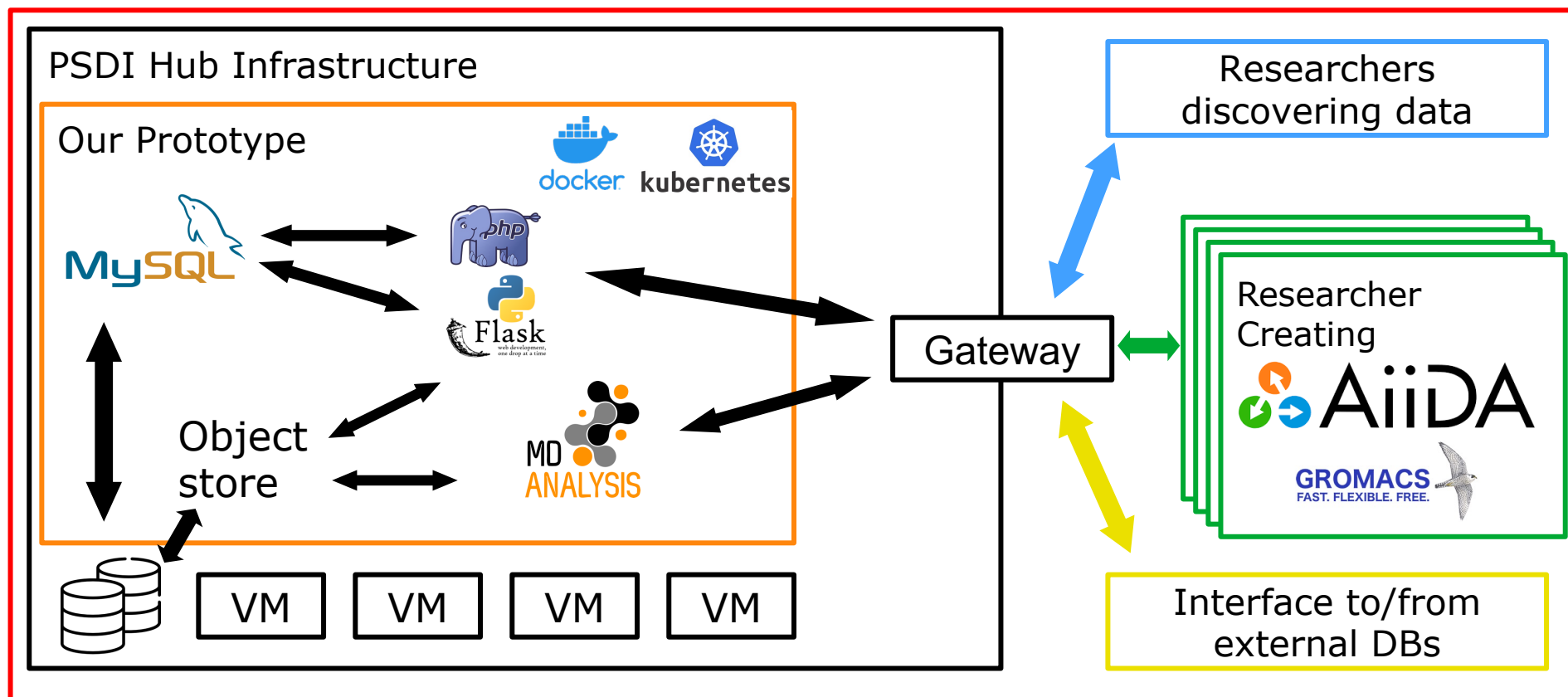
# Our User Environment Prototype



# Our User Environment Prototype



# Our Infrastructure Prototype



# Demo: Lysozyme Minimisation with aiida-gromacs

