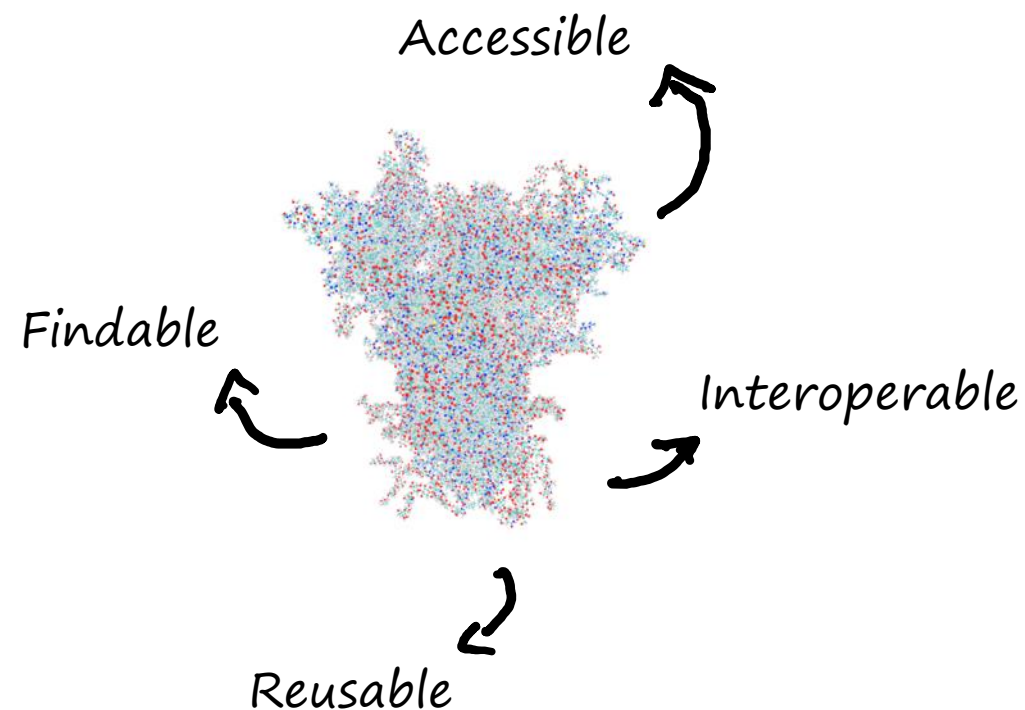


# FAIR Data for the Biomolecular Simulation Community

Jas Kalayan

Senior Computational Scientist, STFC UKRI

[jas.kalayan@stfc.ac.uk](mailto:jas.kalayan@stfc.ac.uk)



\*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,  
 $\approx 10^3$ - $10^4$  snapshots,  
 $\approx$  weeks to perform,  
 $\approx$  TBs of data per simulation

 **JADE**  
Tier 2 HPC

 archer2

 **N8 Bede**  
COMPUTATIONALLY INTENSIVE RESEARCH

 **HEC BioSim**

Start

*...various  
protocols...*

Simulation

# Various Protocols in MD Simulation

1. Get crystal Structure



2. System Preparation



3. Parameterisation



AlphaFold



CHARMM-GUI  
Effective Simulation Input Generator and More



PACKMOL  
Initial configurations for Molecular Dynamics Simulations by packing optimization

GROMACS  
FAST. FLEXIBLE. FREE.



APBS & PDB2PQR

Software for biomolecular electrostatics and solvation

AMBER MD

6. Simulation



5. Equilibration



4. Minimisation

Using the same MD engine

# How to get FAIR Biomolecular MD Simulation Data?

1. Get crystal Structure



2. System Preparation



3. Parameterisation

 Track all steps of complex simulation protocols (data provenance) (**R**)

6. Simulation



5. Equilibration



4. Minimisation

Using the same MD engine



# How to get FAIR Biomolecular MD Simulation Data?

1. Get crystal Structure



2. System Preparation



3. Parameterisation

 Consistent and automated approach to gather simulation metadata (**F**)

6. Simulation



5. Equilibration



4. Minimisation

Using the same MD engine



# PSDI: A Programme for National Infrastructure in the Physical Sciences

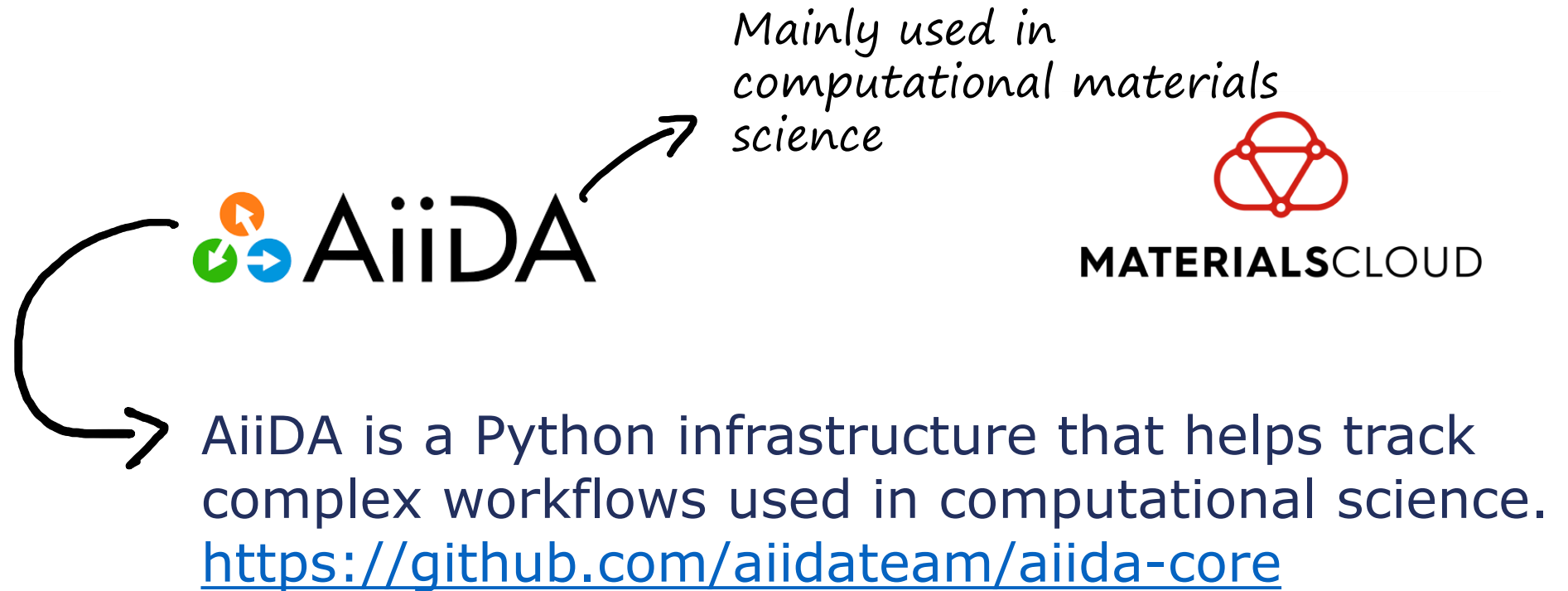
 Providing infrastructure that **connects** experimental and computational data services within Physical Sciences.

 Various initial pathfinders to scope areas of development (catalysis/environment, **biomolecular science**, artificial intelligence).



<https://www.psdi.ac.uk/>

# Biomolecular Simulation Data Provenance using AiiDA



# AiiDA Plugins for Data Provenance with Amber & GROMACS Simulations

Used by 90% of HECBioSim applicants



**PSDI:Biomolecular-team** Unfollow

This github organisation is for the code repositories that are part of the Physical Sciences Data Infrastructure (PSDI) pathfinder for Biomolecular Simulation.

👤 2 followers 📍 United Kingdom 🔗 <https://www.psd.ac.uk/>

🐦 @PSDI\_UK 📄 company/psdiuk

📺 channel/UCd16A90vVYkFNUnWLTl...

🔗 <https://www.jiscmail.ac.uk/cgi-bin/...>

---

👁 View as: Public ▾

You are viewing the README and pinned repositories as a public user.

You can [create a README file](#) or [pin repositories](#) visible to anyone.

### Popular repositories

- aiida-gromacs**  
A GROMACS plugin for AiiDA  
Python ⭐ 5 🍴 2
- aiida-amber**  
An Amber plugin for AiiDA  
Python 🍴 1

**AiiDA**  
PLUGIN REGISTRY

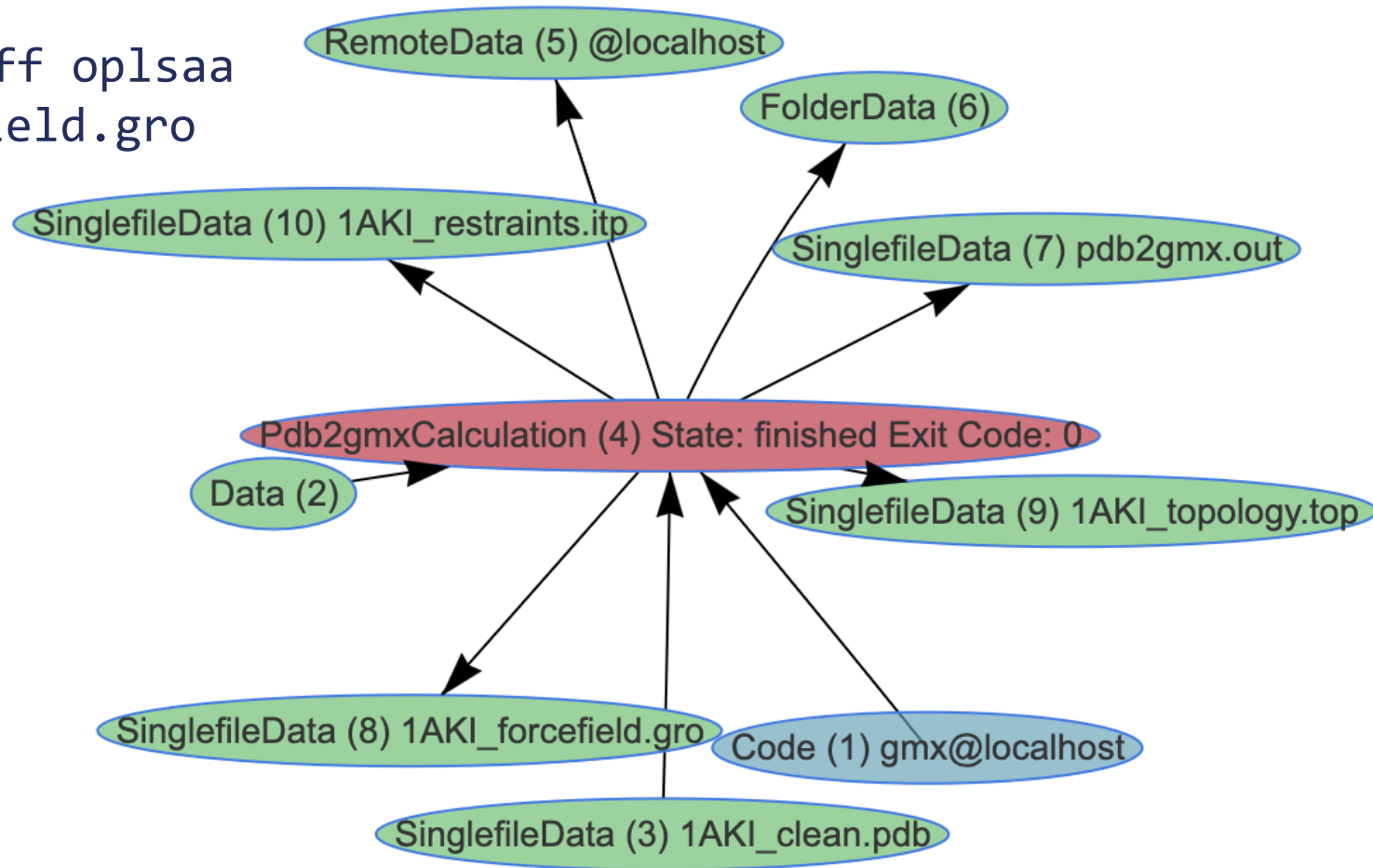
[View on GitHub/register your package]



# Example of Data Provenance with aiida-gromacs

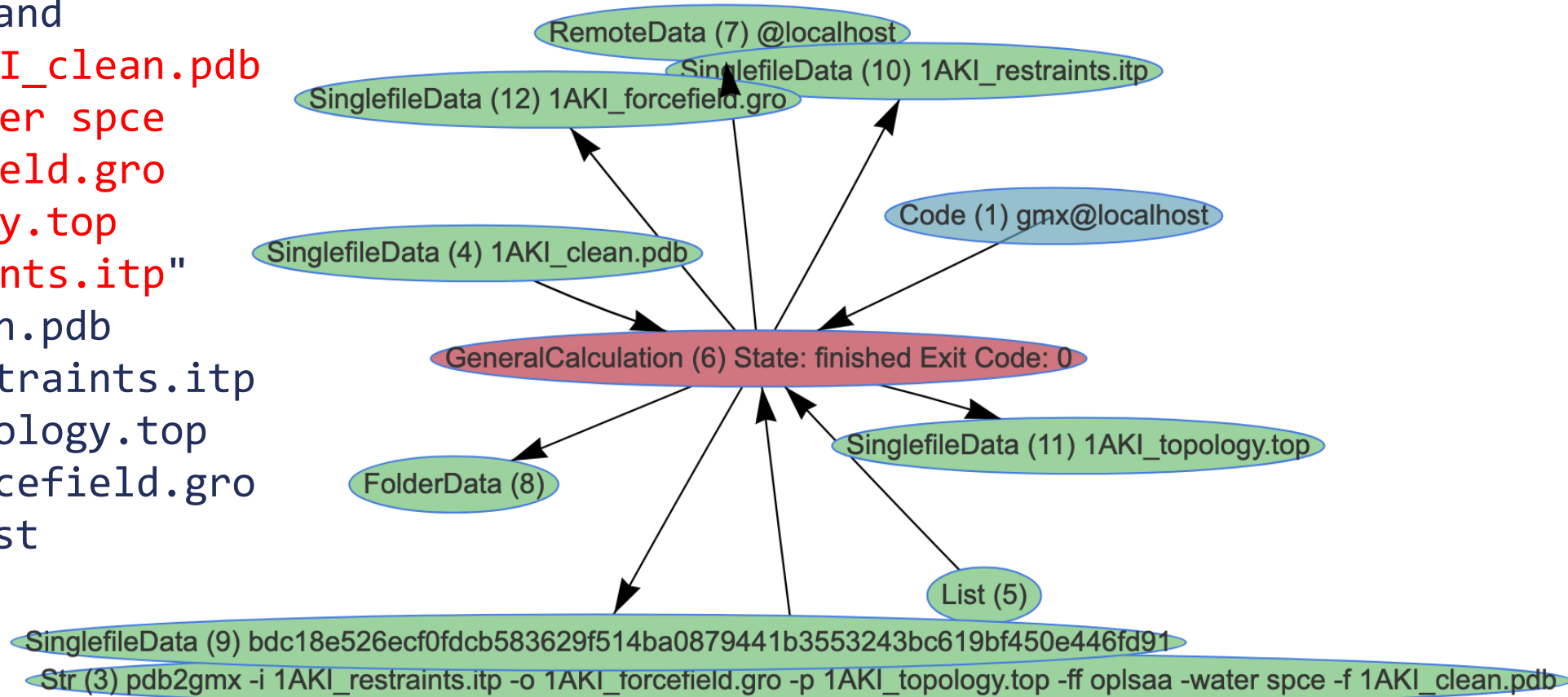
```
$ gm_x_pdb2gmx -f 1AKI_clean.pdb -ff oplsa  
-water spce -o 1AKI_forcefield.gro  
-p 1AKI_topology.top  
-i 1AKI_restraints.itp
```

*Note the underscore!*



# Example of Data Provenance with any CLI command

```
$ 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
```



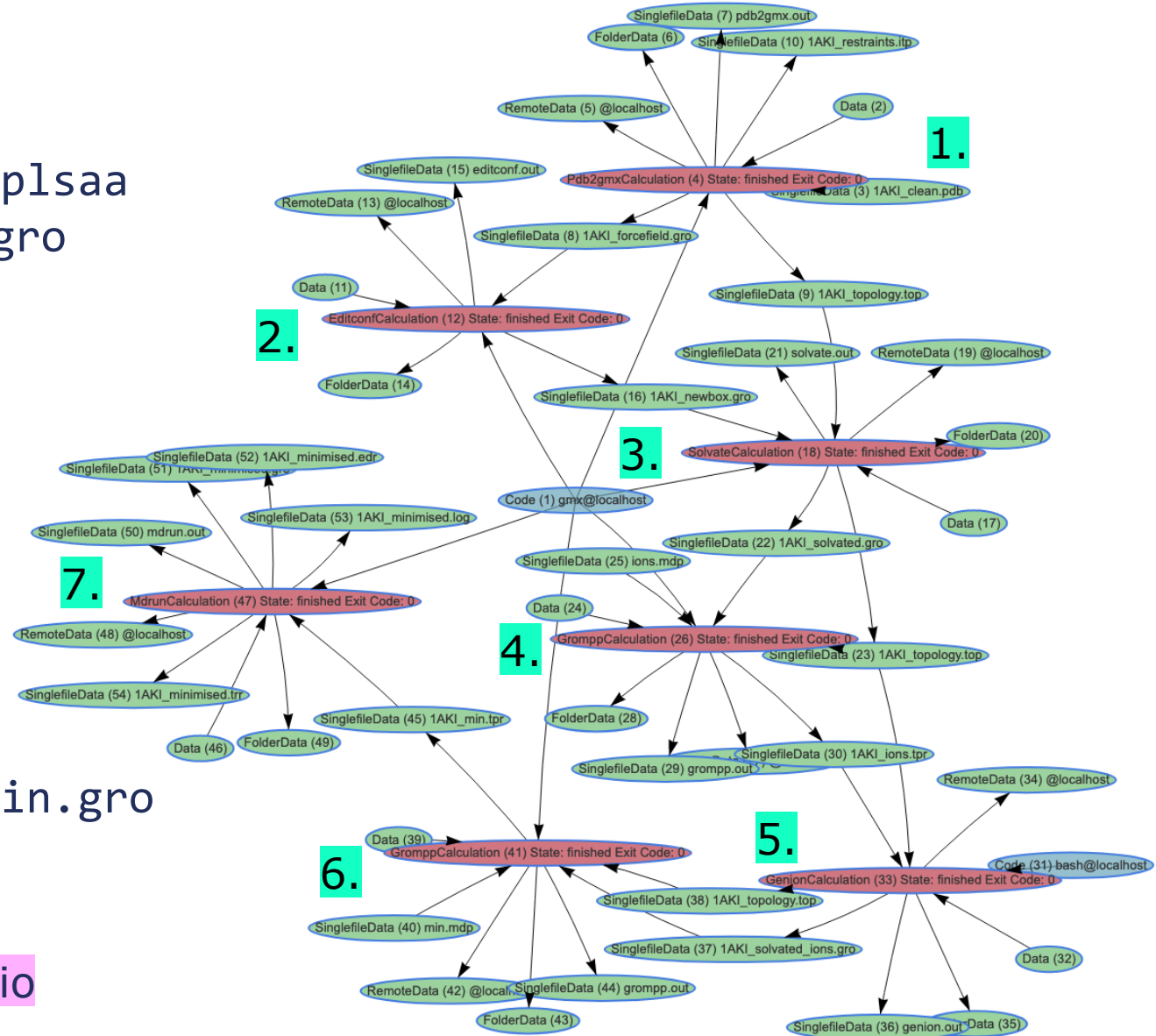
# Example of Data Provenance with aiida-gromacs

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,`

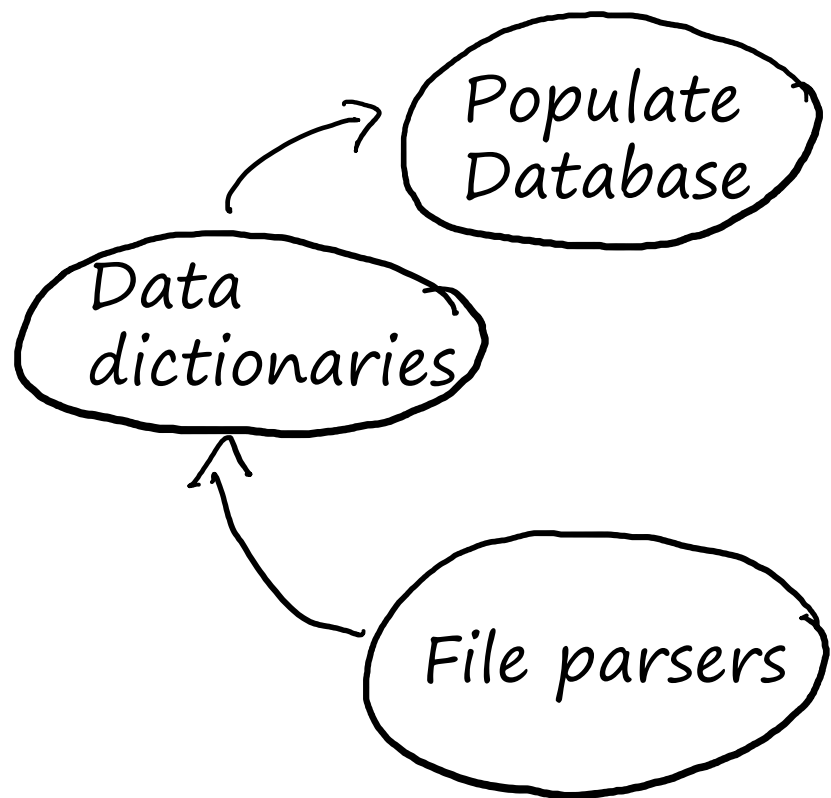
7. `gmx_mdrun -s 1AKI_min.tpr -c 1AKI_min.gro -e 1AKI_min.edr -g 1AKI_min.log -o 1AKI_min.trr`

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





# Automation of Metadata Extraction



 Input and log files are usually excluded / discarded,

 But offer insight to quality of simulations, input parameters, compute,

 Parsers needed for various file types with no defined formats,

 Consensus needed on metadata standards.

# MD Simulation Input File Examples

## Lysozyme NPT equilibration

### Amber

```
&cntrl
  imin=0,      ! No minimization
  ntx=5,      ! Read coords and velocities from coordinate file
  irest=1,    ! Restart simulation
  nstlim=10000, ! Number of MD steps
  dt=0.002,   ! Time step (ps)
  ntf=2,      ! Bond interactions involving H omitted
  ntc=2,      ! Use SHAKE to constrain bond distances to H
  temp0=300.0, ! Reference system temperature (K)
  ntp=5000,   ! Print progress of min steps every ntp steps
  ntwx=5000,  ! Print coordinates every ntp steps
  cut=10.0,   ! Non-bonded cutoff value (Angstrom)
  ntb=2,      ! Constant pressure, PBCs for nb interactions
  ntp=1,      ! MD with isotropic position scaling
  ntt=3,      ! Langevin thermostat
  barostat=1, ! Berendsen barostat
  gamma_ln=2.0, ! Collision frequency (ps^-1)
  ig=-1,     ! random seed
/
```

#### Group together shared terms:

total steps: nstlim, nsteps = 10000

temperature: temp0, ref\_t = 300 K

timestep: dt = 0.002 ps

coordinate steps: ntwx, nstxout = 5000

non-bonded cutoffs: cut, rcoulomb, rvdw = 1 nm

...

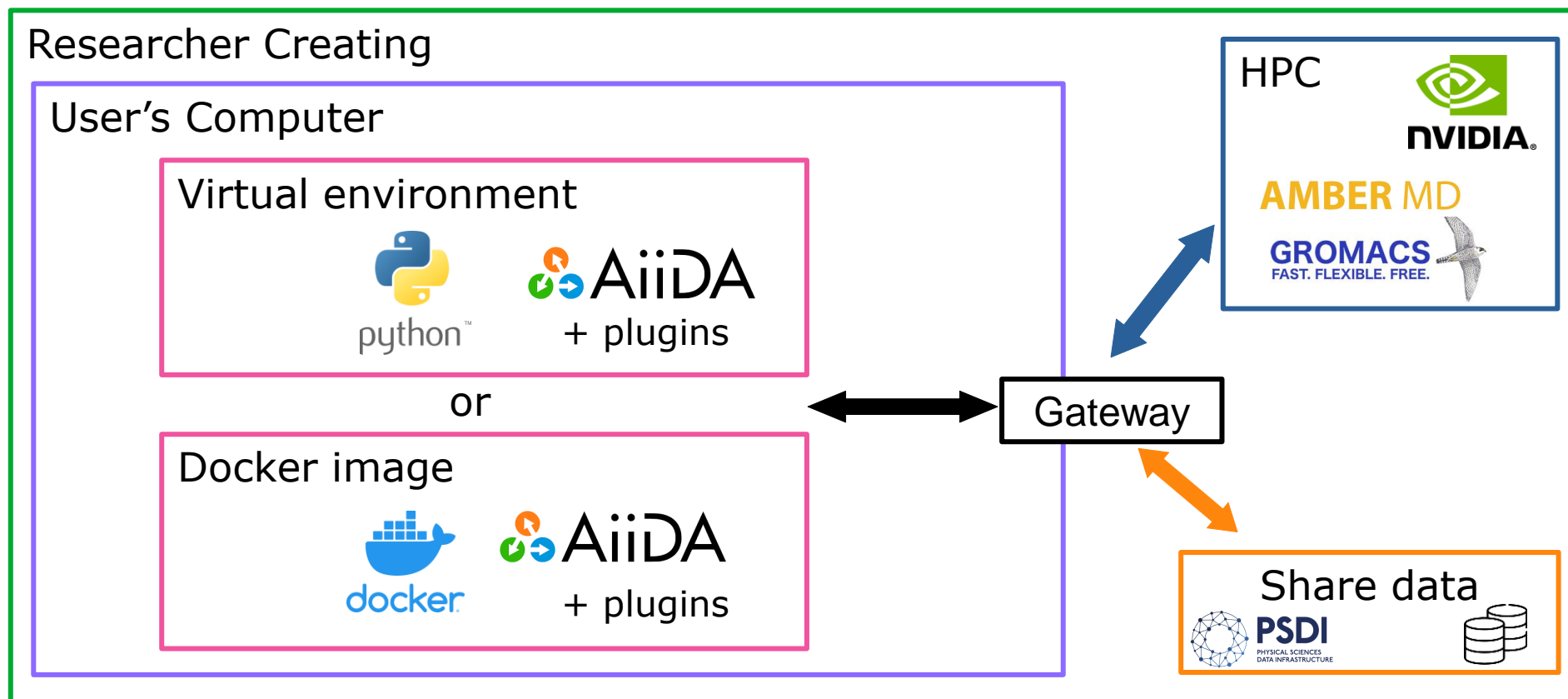
.

.

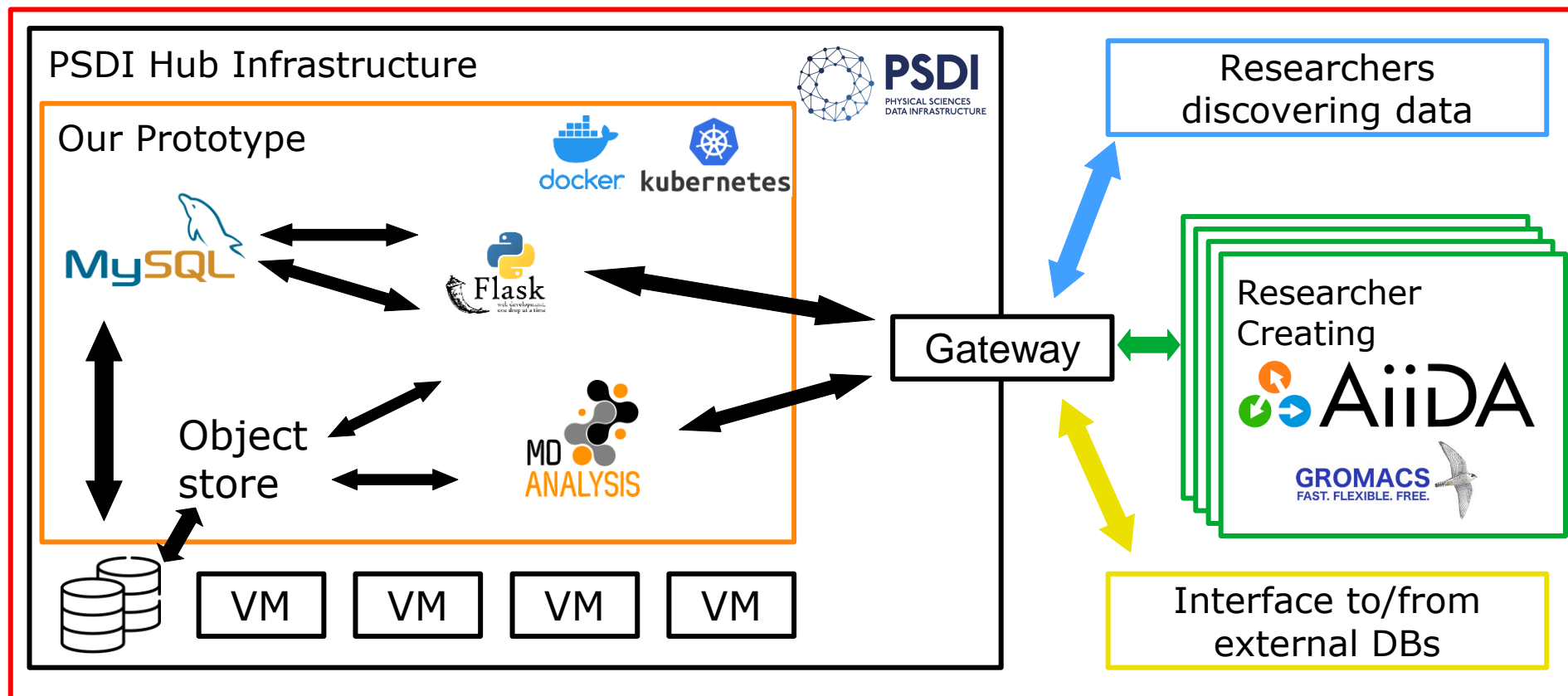
### GROMACS

```
title = OPLS Lysozyme NPT equilibration
; Run parameters
integrator = md ; leap-frog integrator
nsteps = 10000 ; 2 * 500000 = 1000 ps (1 ns)
dt = 0.002 ; 2 fs
; Output control
nstxout = 5000 ; suppress bulky .trr file by specifying
nstvout = 5000 ; 0 for output frequency of nstxout,
nstfout = 5000 ; nstvout, and nstfout
nstenergy = 5000 ; save energies every 10.0 ps
nstlog = 5000 ; update log file every 10.0 ps
; Bond parameters
continuation = yes ; Restarting after NPT
constraint_algorithm = lincs ; holonomic constraints
constraints = h-bonds ; bonds involving H are constrained
lincs_iter = 1 ; accuracy of LINCS
lincs_order = 4 ; also related to accuracy
; Neighborsearching
cutoff-scheme = Verlet ; Buffered neighbor searching
ns_type = grid ; search neighboring grid cells
nstlist = 10 ; 20 fs, largely irrelevant with Verlet scheme
rcoulomb = 1.0 ; short-range electrostatic cutoff (in nm)
rvdw = 1.0 ; short-range van der Waals cutoff (in nm)
; Electrostatics
coulombtype = PME ; Particle Mesh Ewald for long-range electrostatics
pme_order = 4 ; cubic interpolation
fourierspacing = 0.16 ; grid spacing for FFT
; Temperature coupling is on
tcoupl = V-rescale ; modified Berendsen thermostat
tc-grps = Protein Non-Protein ; two coupling groups - more accurate
tau_t = 0.1 0.1 ; time constant, in ps
ref_t = 300 300 ; reference temperature, in K
; Pressure coupling is on
pcoupl = Parrinello-Rahman ; Pressure coupling on in NPT
pcoupltype = isotropic ; uniform scaling of box vectors
tau_p = 2.0 ; time constant, in ps
ref_p = 1.0 ; reference pressure, in bar
```

# Our User Environment Prototype



# Our Infrastructure Prototype





# Thank you for your attention!

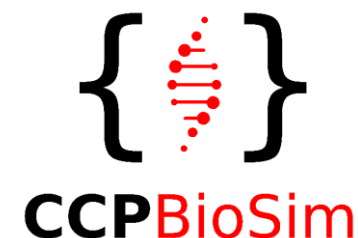
Acknowledgments to:

Kin Chao, Sarah Rouse 

Joel Greer, Tom Burnley, Martyn Winn



James Gebbie-Rayet (PI)



PSDI acknowledges the funding support by the EPSRC grants EP/X032701/1, EP/X032663/1 and EP/W032252/1