

MD Simulations – day I



Toni Giorgino

National Research Council of Italy

toni.giorgino@cnr.it

www.giorginolab.it



for Prof. Fuxreiter's course @ University of Padova

May 2024

<https://github.com/giorginolab/MD-Tutorial-Data>

This class

- Molecular dynamics is a powerful tool for studying molecular systems
- OpenMM is a software library that allows for efficient and customizable MD simulations
- It's exemplary of a modern well-maintained open-source library:
 - CI infrastructure, developed on GitHub
 - C++ w/ Python bindings
- We'll use the latter, testing *live* on Google Colab.

Molecular Dynamics

What is MD?

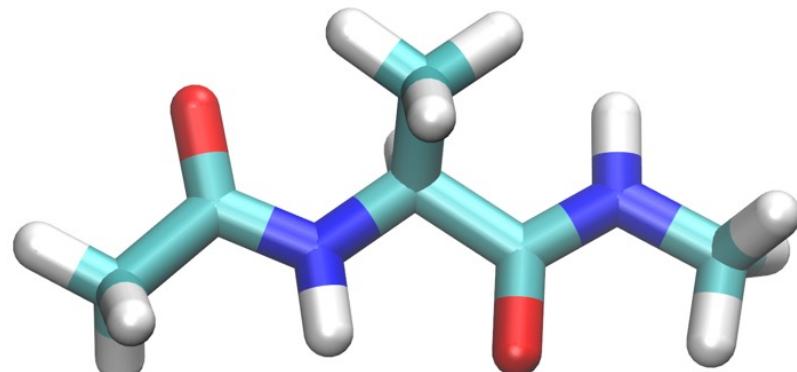
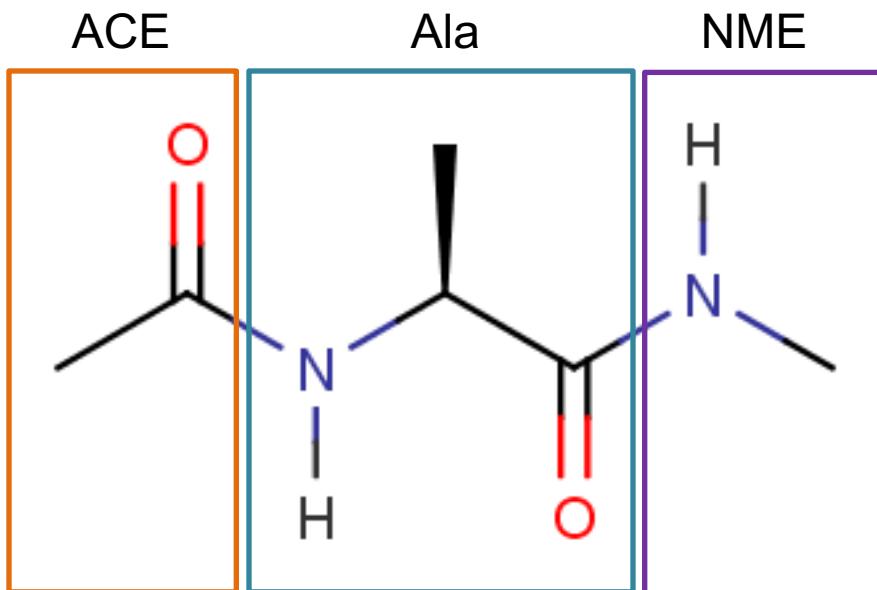
- Attempt the most detailed description of a system which is
 1. atomistic
 2. classical
- Model the internal *forces*...
- ...in order to *integrate* the motion
- Hope in convergent *sampling*

$$\vec{F}_i(\mathbf{x}) = m_i \ddot{\mathbf{x}}_i$$

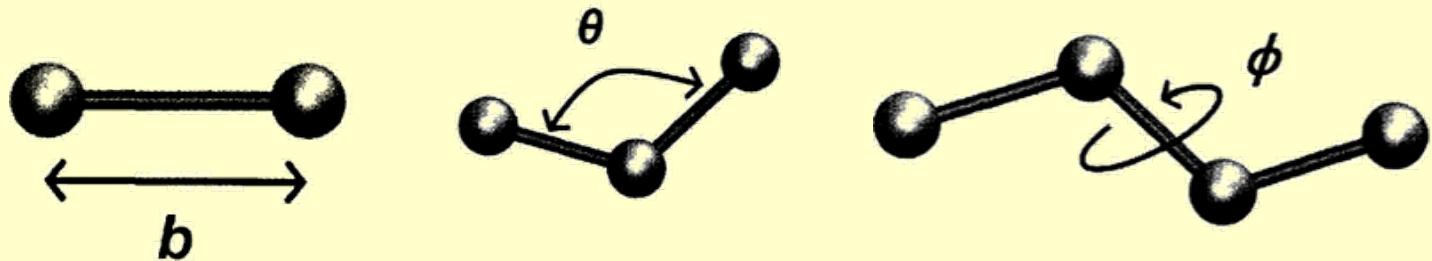
Assumptions

- In this tutorial we shall deal with **unbiased** sampling approaches with **explicit** solvent, i.e.
 - no added forces except the "physical" ones in your system;
 - all of the system (including water molecules) have atomic resolution.
- Also, current classical MD does not address, by design, the following:
 - Chemical reactions, e.g. catalysis, phosphorylation, ubiquitination etc.
 - Protonation changes
- Finally, small molecules pose distinct challenges and need a separate, expensive **parameterization** step.

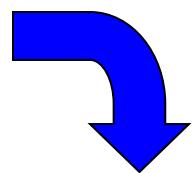
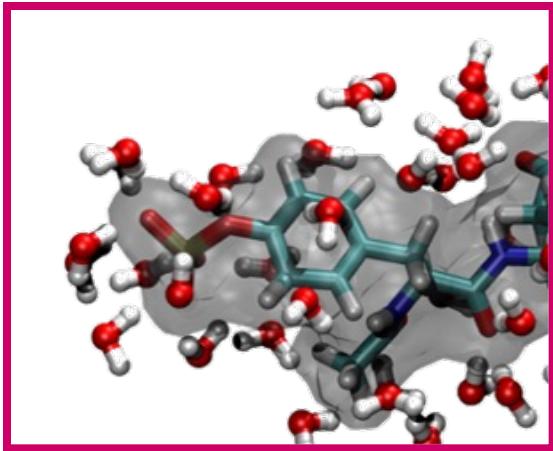
Alanine “dipeptide”



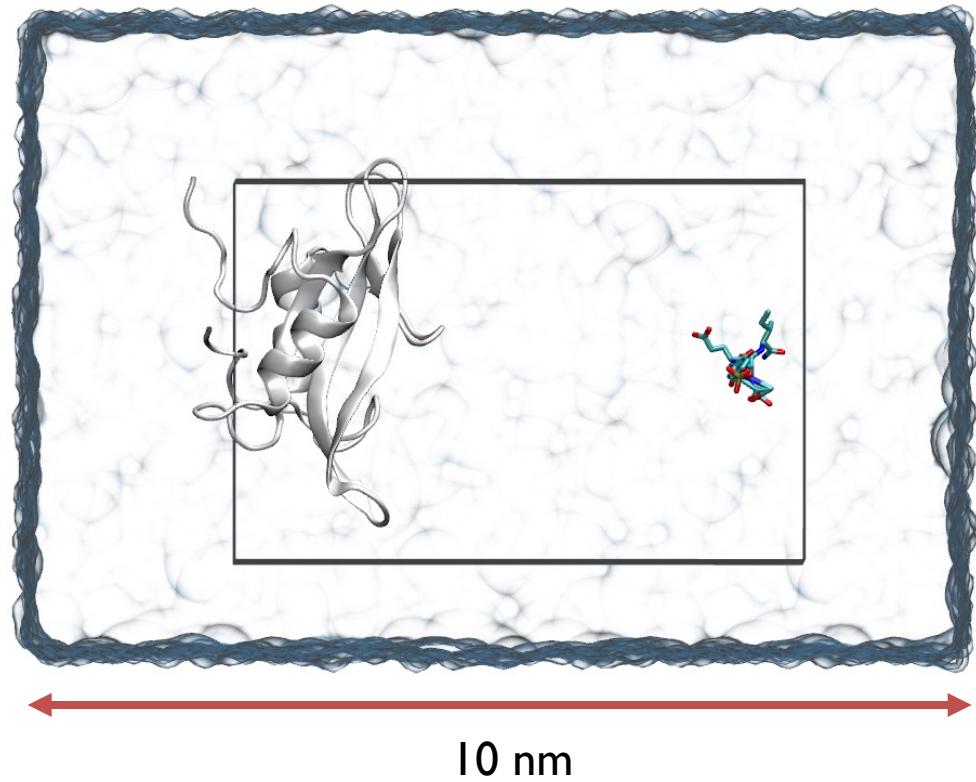
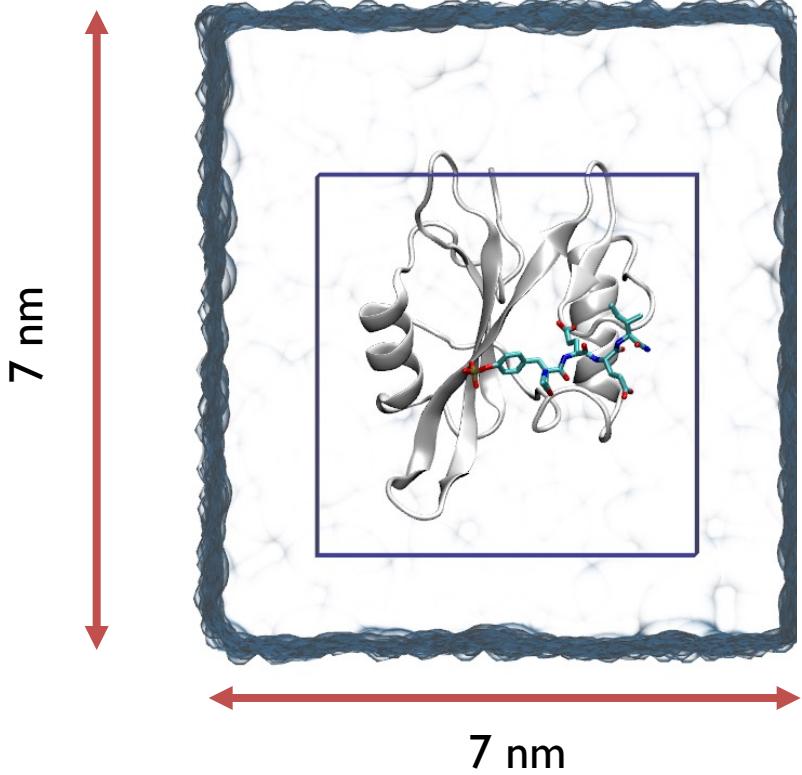
Bonded
energy terms
+ Electrostatics
+ VdW



The forcefield is a database of interatomic parameters



- **Explicit solvation**
- $\rightarrow \mathcal{O}(10^5)$ atoms
- **Unbiased dynamics**
- Update every 10^{-15} s (1 fs)



Event ≡ Binding / Unbinding / Folding / Unfolding / ...

$$* \frac{1}{t_{\text{on}}} = \text{association rate of SH2-pYEEI} \times [\text{pYEEI}]$$

Large gain

Ability to “play” biomolecular processes at
all-atom resolution in silico

Molecular bases of folding, binding, selectivity, gating...

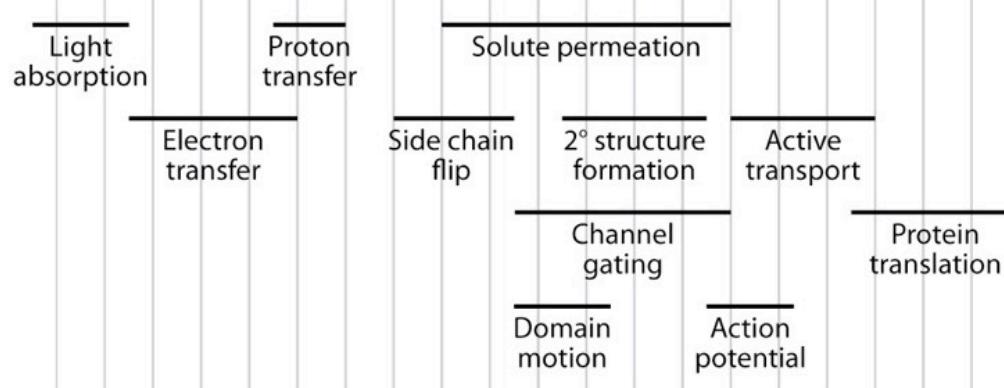
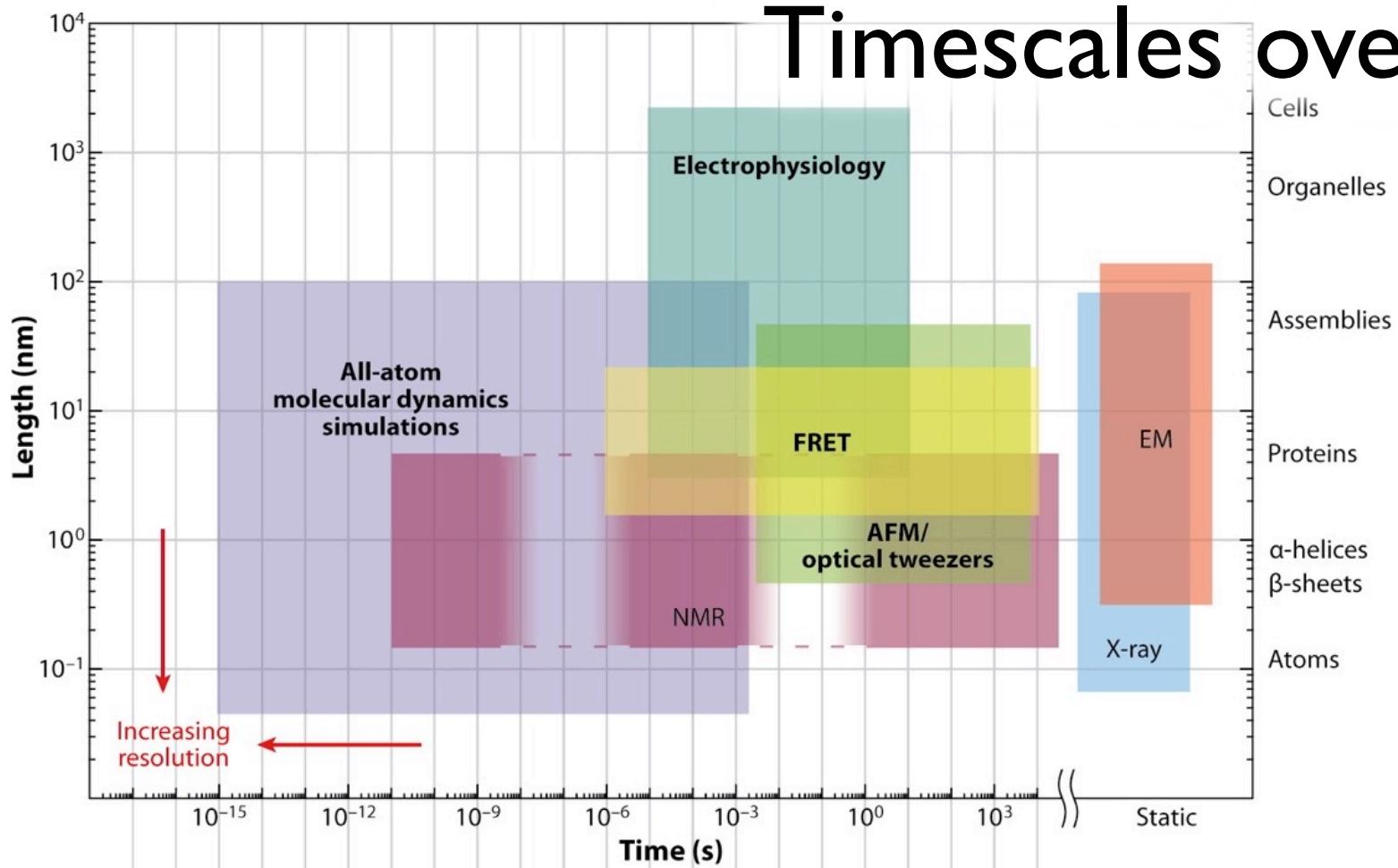
Large cost

E.g.*: $t_{\text{on}} \sim 30 \mu\text{s} \rightarrow$
→ **10¹⁰** integration timesteps →
→ 15 years single-CPU compute time

MD is entirely about timescales

- Your ability to obtain quantitative results is severely limited by the sampling ability you have. You will only be able to reach phenomena occurring on the sampled timescales, or shorter.
 - Sidechain rearrangements, diffusion-limited processes: usually possible *
 - Local flexibility: usually possible *
 - Membrane environments: ok-ish
 - Binding: hard but not impossible
 - Folding: very hard but not impossible
 - [*] Unless there are significant barriers.

Timescales overview



Dror RO, et al. 2012.

Annu. Rev. Biophys. 41:429–52

Patience and other limits

- The following factors affect the running speed (usually expressed in ns per simulation day, ns/day)
 - System size. Reasonable is 100 AA \sim 30,000 atoms.
 - Computer speed. Forget laptops.
 - Definitely use GPUs.
 - Software.

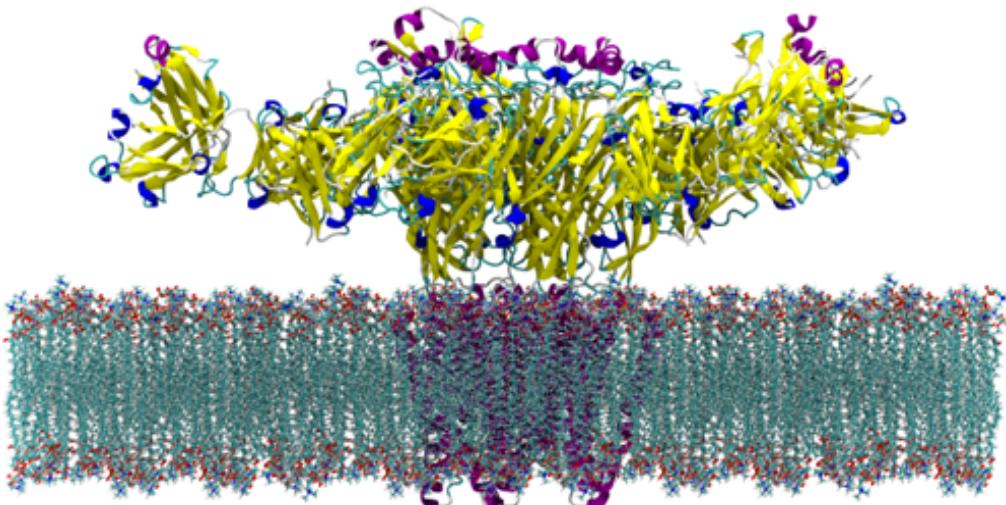
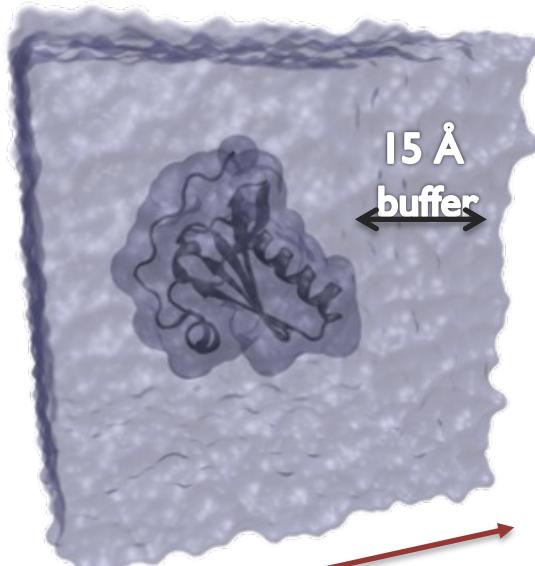
"Reasonable" system sizes:

Compact

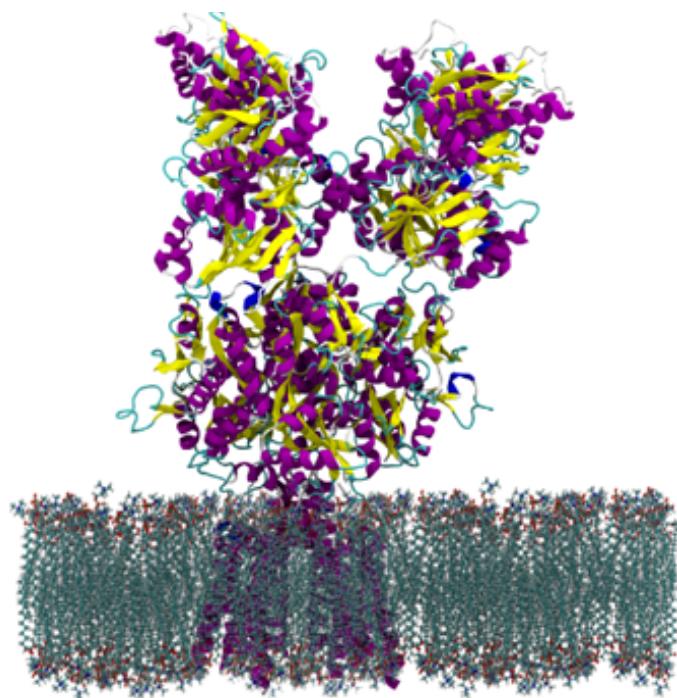
< 500 residues

< $(150 \text{ \AA})^3$

< 1 M atoms



(a) 3RHW



(b) 4U1W

Roughly in order of trouble

- Biological assembly / asymmetric unit
- Disordered regions
- Transmembrane?
- Chains
- Non standard residues
- Ligands
 - of interest
 - of no interest
- Metal ions
- Missing loops
- Termini
- Missing atoms
- Water molecules



Forcefields and software

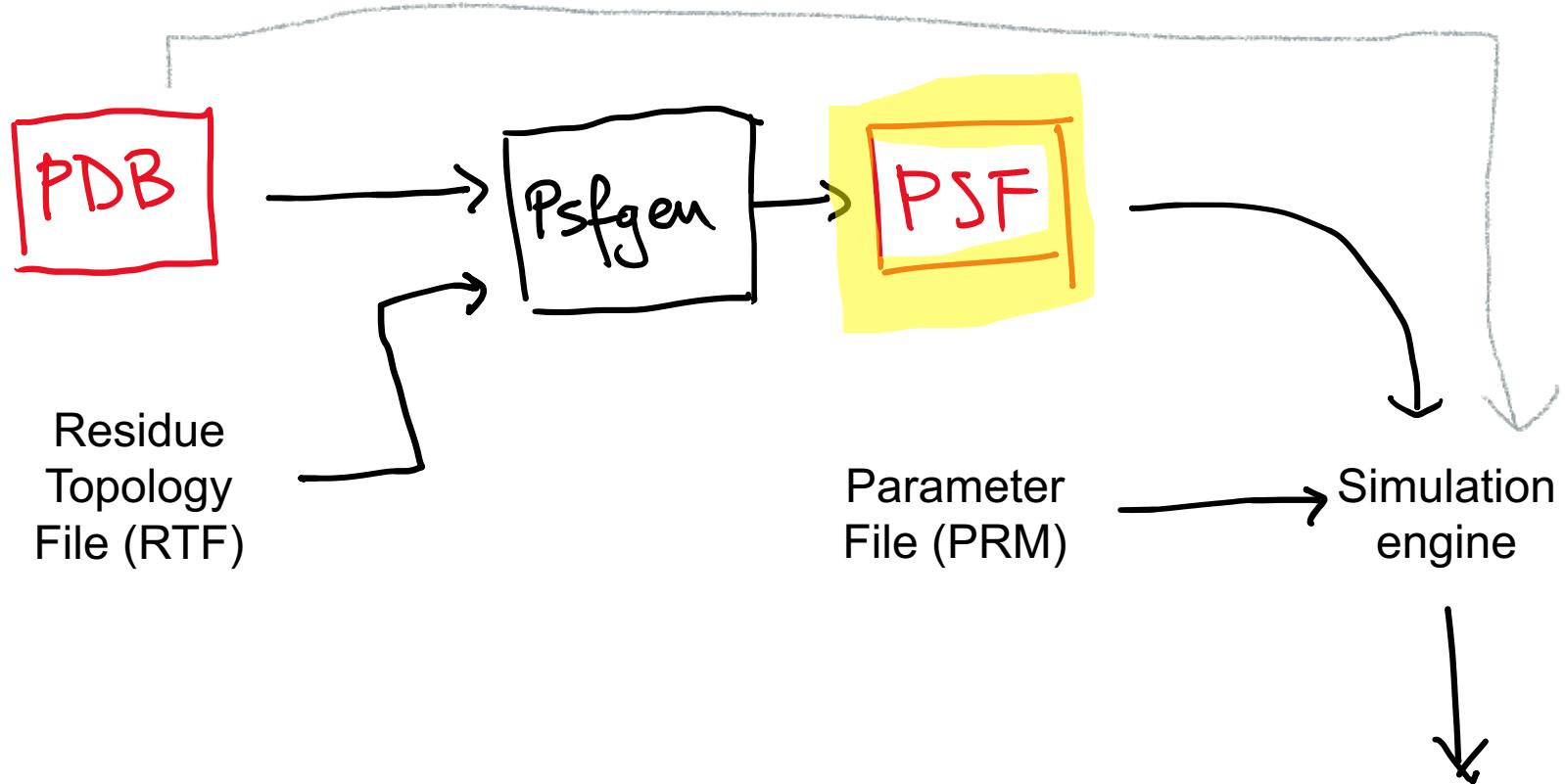
They run in “families”

- They are “databases” of atomic parameters
- We deal with non-polarizable all-atom ones:
AMBER and **CHARMM**
 - There are many others, these are s.o. art
- Most notable difference:
molecular types supported (lipids, drugs, ...)
- They differ in (software) build procedures
- However OpenMM unifies them

The CHARMM family

- Originally coupled with the CHARMM software (MacKerell, Karplus), but independent
- Variants of note: C36M
 - [toppar_c36_jul22.tgz](#)
- Based on RTF (templates) and PRM (parameters) files
- Build: **psfgen** -> xxx.psf
 - <https://www.academiccharmm.org/>
 - https://mackerell.umaryland.edu/charmm_ff.shtml

CHARMM system layout



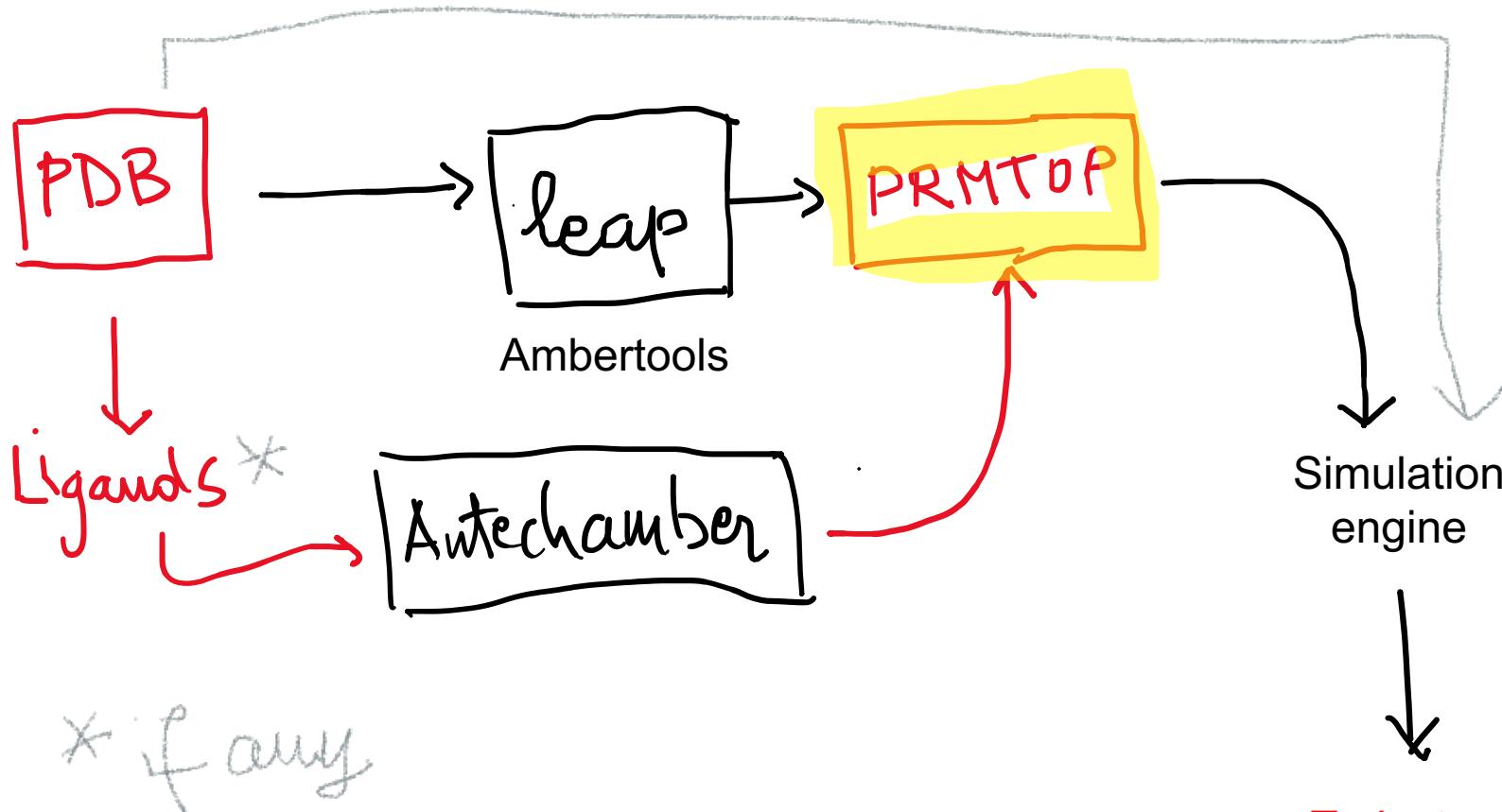
psfgen matches the given structure to FF templates provided in RTF files, completing missing atoms.
The actual parameter values are read at a later stage.

Trajectory

The AMBER family

- Originally coupled with the Amber software (Case, Merz, Kollmann, ...), but independent
- Variants of note: ff19SB
 - <https://ambermd.org/AmberTools.php>
- Based on **tleap** + its database
- Build: **tleap** --> xxx.prmtop

Amber system layout



tleap matches the given structure to FF templates, completing missing atoms, *merging parameters*.

FF choice

Largely equivalent, i.e. subtle differences only appear at late stages.
FF choice is mostly due to the species in the system that one intends to model.

	CHARMM	AMBER
Proteins, peptides	++	++
Water, ions	++	++
Lipids (membranes)	++	+
Small molecules	+ (CGenFF)	++ (GAFF2)
Post-translational modif. *	+/-	+/-
Non-standard charge states	+/-	+/-
DNA	-	-
RNA	--	--
Non-standard AAs *	--	--

Somewhat subjective!

* Check individually

An example for Amber

Molecule/Ion Type	Force Field
protein	ff19SB
DNA	OL21
RNA	OL3
carbohydrates	GLYCAM_06j
lipids	lipids21
organic molecules (usually ligands)	gaff2
ions	<ul style="list-style-type: none">•should be matched to water model; see force fields for ions for further discussion
water model	<ul style="list-style-type: none">•should be matched to atomic ions; common water models include tip3p, spc/e, tip4pew, and OPC

System building

The build procedure

- Know your system
- Build
 - Cleanup
 - Assign forcefield
 - Solvate
 - Ionize
- Minimize
- Equilibrate
- Run

The build procedure

- It used to be somewhat convolved
- Generally
 - take PDB coordinates
 - filter out unwanted species
 - solvate
 - ionize
- Automation was (it still is) challenging
- OpenMM unified the build process (but it is still possible to use the old tools)

High-Throughput Automated Preparation and Simulation of Membrane Proteins with HTMD

Stefan Doerr,^{†,||,ID} Toni Giorgino,^{‡,||,ID} Gerard Martínez-Rosell,^{†,ID} João M. Damas,^{¶,ID} and Gianni De Fabritiis^{*,†,§,ID}

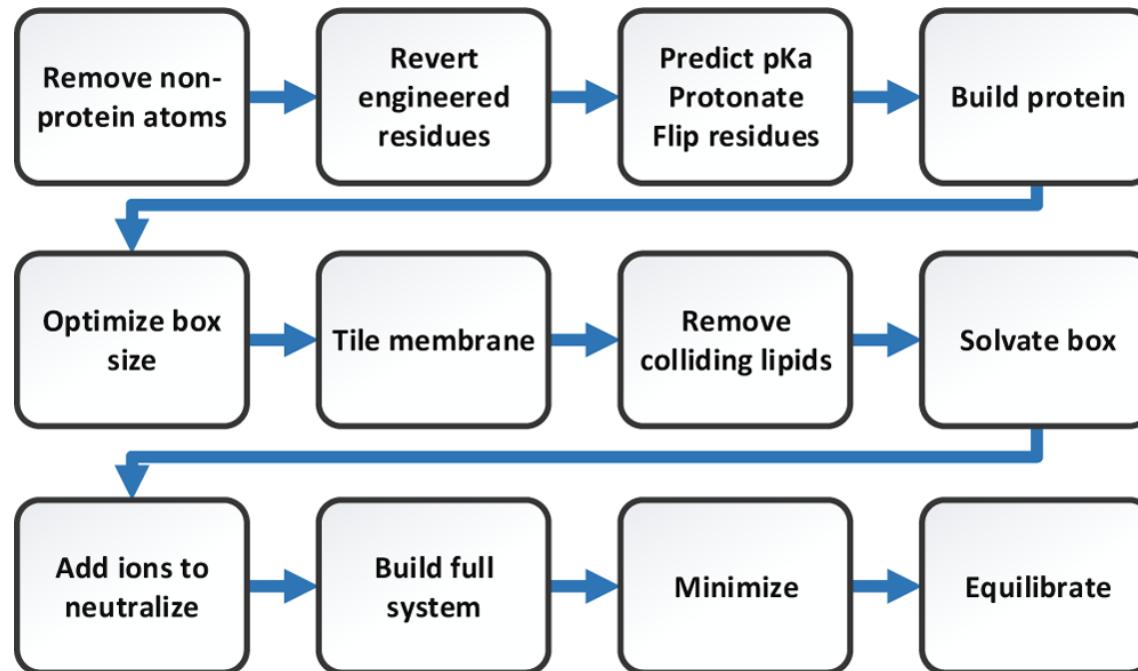


Figure 2. OPM building and simulation protocol workflow starting from a PDB file and ending with an equilibrated system.

Small molecules in the Amber system

Ligands

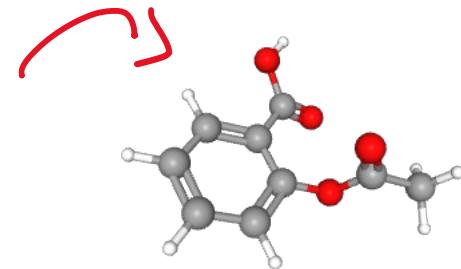
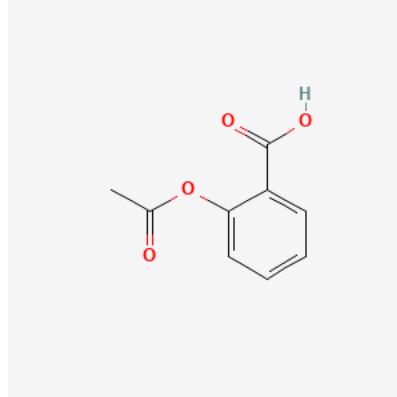
- Anything non-peptide is problematic
- Some common molecules are pre-computed
- Generic small molecules need *parameterization*
- AmberTools has *antechamber* + *GAFF*
 - Semi-automated
 - Do mind stereochemistry, tautomers, protonation!
 - Partial charges are assigned by RESP
 - Other force terms are pattern-matched
- CHARMM has CGenFF, web only

Ligands

CC(=O)OC1=CC=CC=C1C(=O)O

Smiles

(databases
for virtual
screening)



planar
(editing)

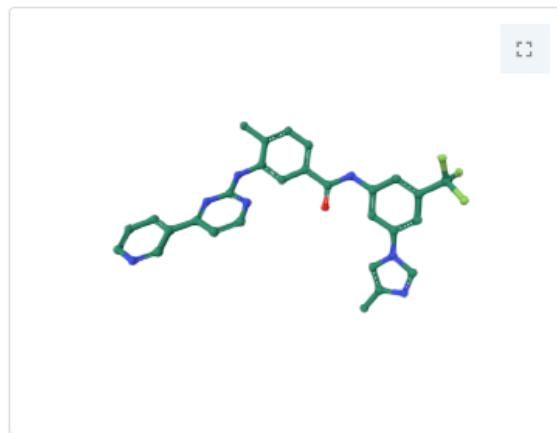
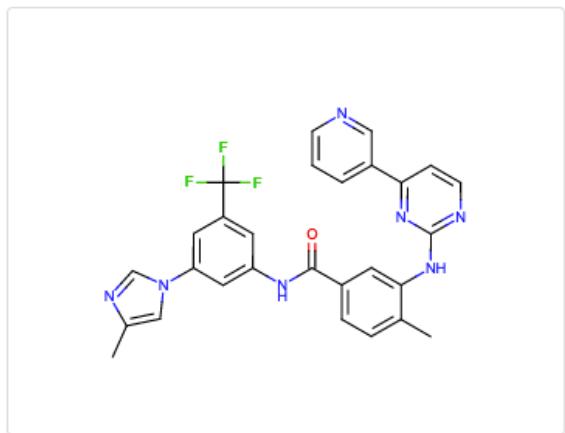
3D
conformers

generated in
crystallized

further
process
(dock, parameterize)

- mol2
- sdf

RCSB PDB “Ligand expo”



[Toggle Hydrogen](#) [Toggle Labels](#)

[Display Files ▾](#) [Download Files ▾](#) [Data API](#)

NIL

Nilotinib

Find entries where: NIL

is present as a standalone ligand in 3 entries
[search](#)

Find related ligands:

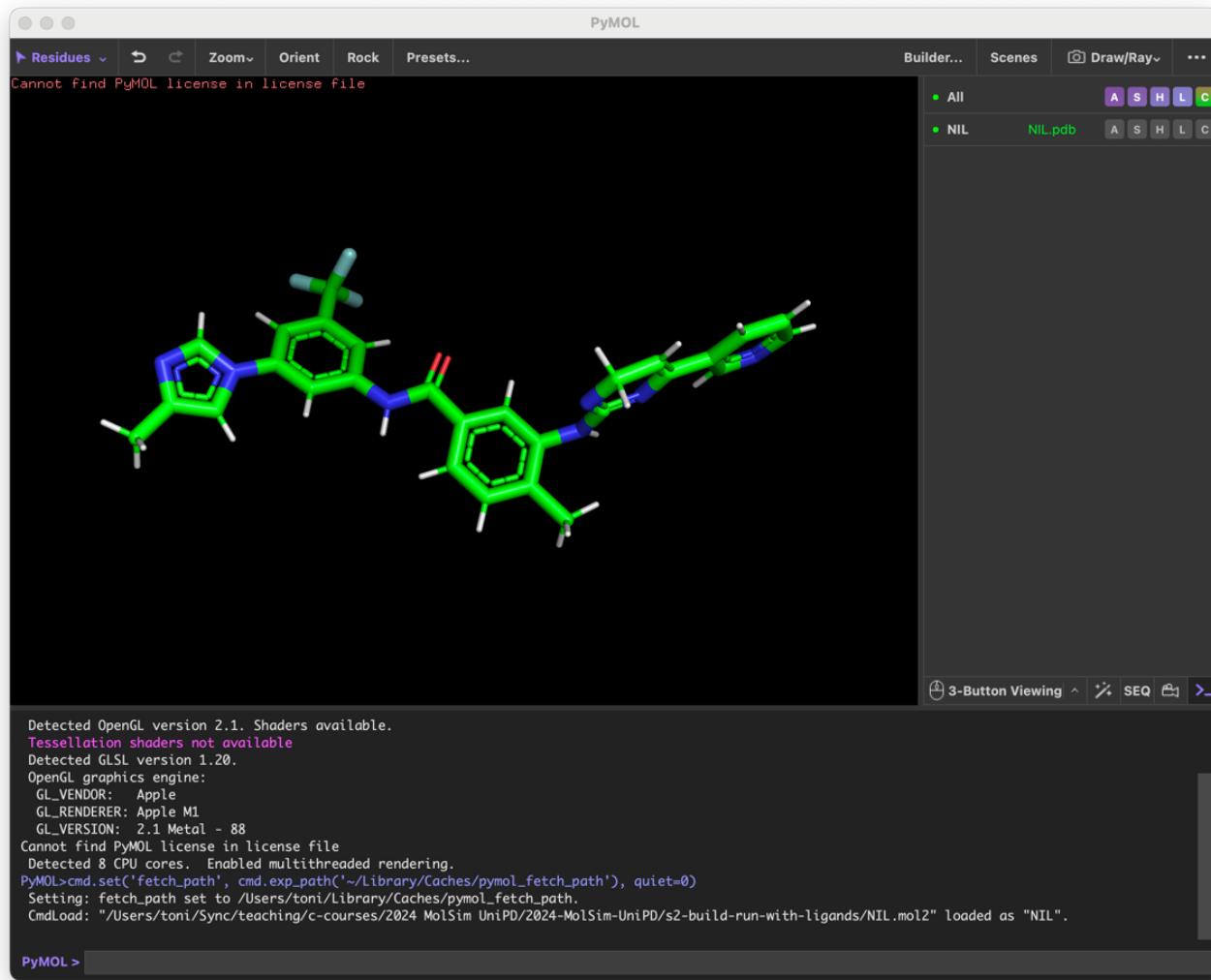
- [Similar Ligands \(Stereospecific\)](#)
- [Similar Ligands \(including Stereoisomers\)](#)
- [Similar Ligands \(Quick Screen\)](#)
- [Similar Ligands \(Substructure Stereospecific\)](#)
- [Similar Ligands \(Substructure including Stereoisomers\)](#)

Chemical Component Summary

Name	Nilotinib
Synonyms	4-methyl-N-[3-(4-methyl-1H-imidazol-1-yl)-5-(trifluoromethyl)phenyl]-3-[(4-pyridin-3-yl)pyrimidin-2-yl]amino]benzamide
Identifiers	4-methyl-N-[3-(4-methylimidazol-1-yl)-5-(trifluoromethyl)phenyl]-3-[(4-pyridin-3-yl)pyrimidin-2-yl]amino]benzamide
Formula	C ₂₈ H ₂₂ F ₃ N ₇ O
Molecular Weight	529.516
Type	NON-POLYMER

Chemical Details

Formal Charge	0
Atom Count	61
Chiral Atom Count	0
Bond Count	65
Aromatic Bond Count	31



Practice, finally

Open these

- www.rcsb.org/structure/3CS9
- github.com/giorginolab/2024-MoLSim-UniPD

2024-MoLSim-UniPD / amber-build-run-with-ligand / 

OpenMM



OpenMM

High performance, customizable molecular simulation.

.org

- OpenMM is a molecular dynamics simulation toolkit that allows for high-performance simulations of biomolecules.
- Allows for simulation of a variety of molecular systems, including proteins, nucleic acids, and small molecules
- OpenMM supports a wide range of force fields and integrators and can run on CPUs and GPUs.
- Open source, written in C++ with Python and other language bindings available

Basic Workflow (object-oriented)

- I. Download, complete and edit the structure:
 - **Topology** (i.e. the identity of atoms, bonds, etc)
 - **Positions** (i.e. the starting coordinates)
2. Create the **system** object.
3. Create the **integrator** object.
4. Create and add custom **forces** to system if needed.
5. Define the **simulation** object.
6. Set the initial positions and velocities.
7. Minimize.
8. Run the simulation.
9. (Analyze the results.)

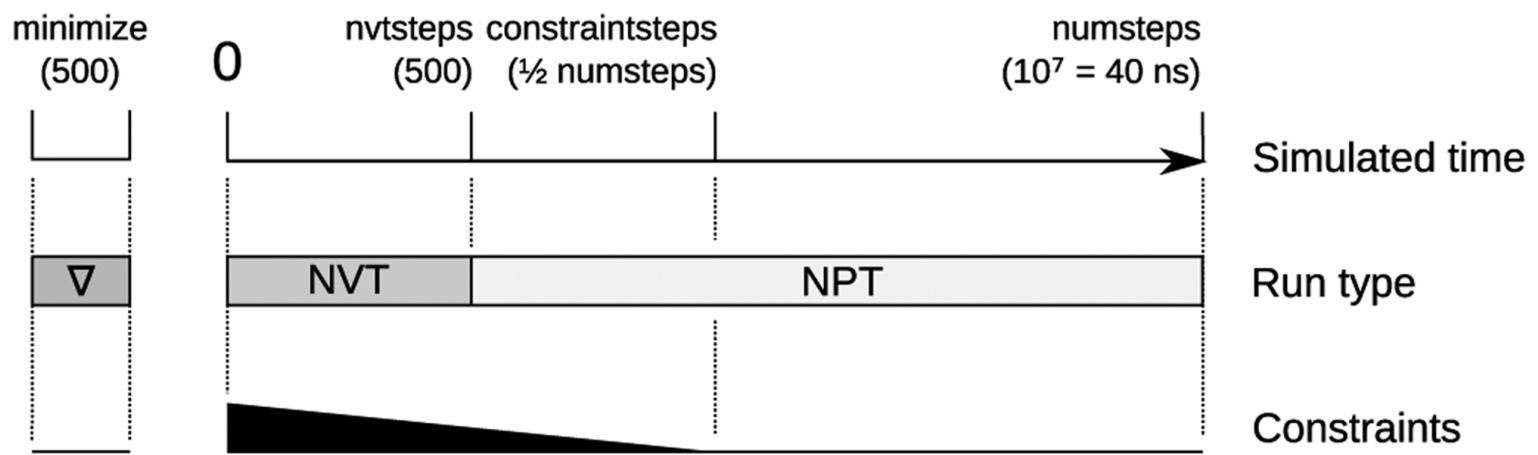
Integrators

- ...are algorithms that solve the equations of motion for a system
- OpenMM includes several integrators, e.g. Langevin dynamics, Verlet integrator, and Monte Carlo barostat
- Different integrators are appropriate for different types of simulations and conditions (e.g.: NPT vs NVT)

Simulating a system

- Once a system has been defined and the force field and integrator selected, it can be simulated
- The simulation (run) involves running a series of steps, where each step involves calculating the forces on each atom, integrating the equations of motion, and updating the system's coordinates
- After the simulation, data analysis can be performed to obtain information about the system's behavior and properties

Equilibration



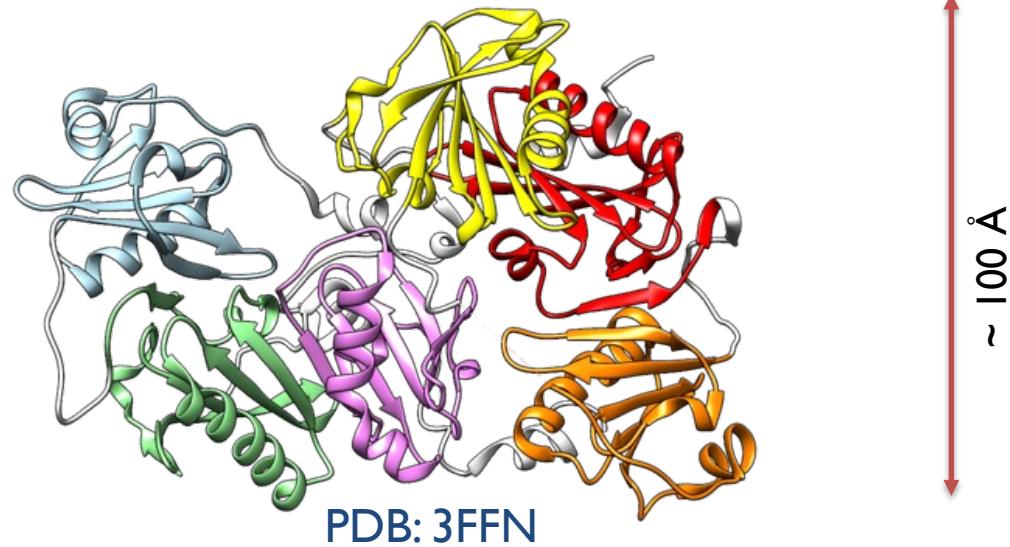
Orders of magnitude

- Integration timestep: 1-4 fs
 - How often forces are calculated, i.e. the finest time granularity. No periods faster than this.
- Logging interval: $O(1 \text{ ps})$
 - How often to print energies etc.
- Frame interval: $O(100 \text{ ps})$
 - How often to save a snapshot
- Total length: $O(10-1000 \text{ ns})$

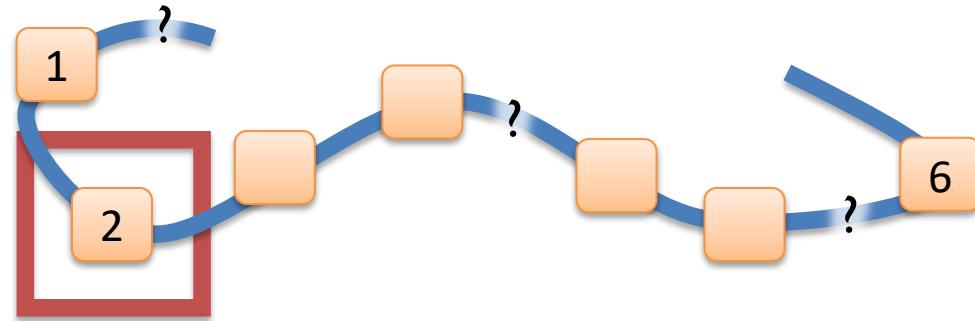
Let's pick a test system



Inactive
No Ca^{2+} .
Closed form.
Resolved (2008).



Active
~ mM Ca^{2+} .
Active form.
Dynamic. Elusive.



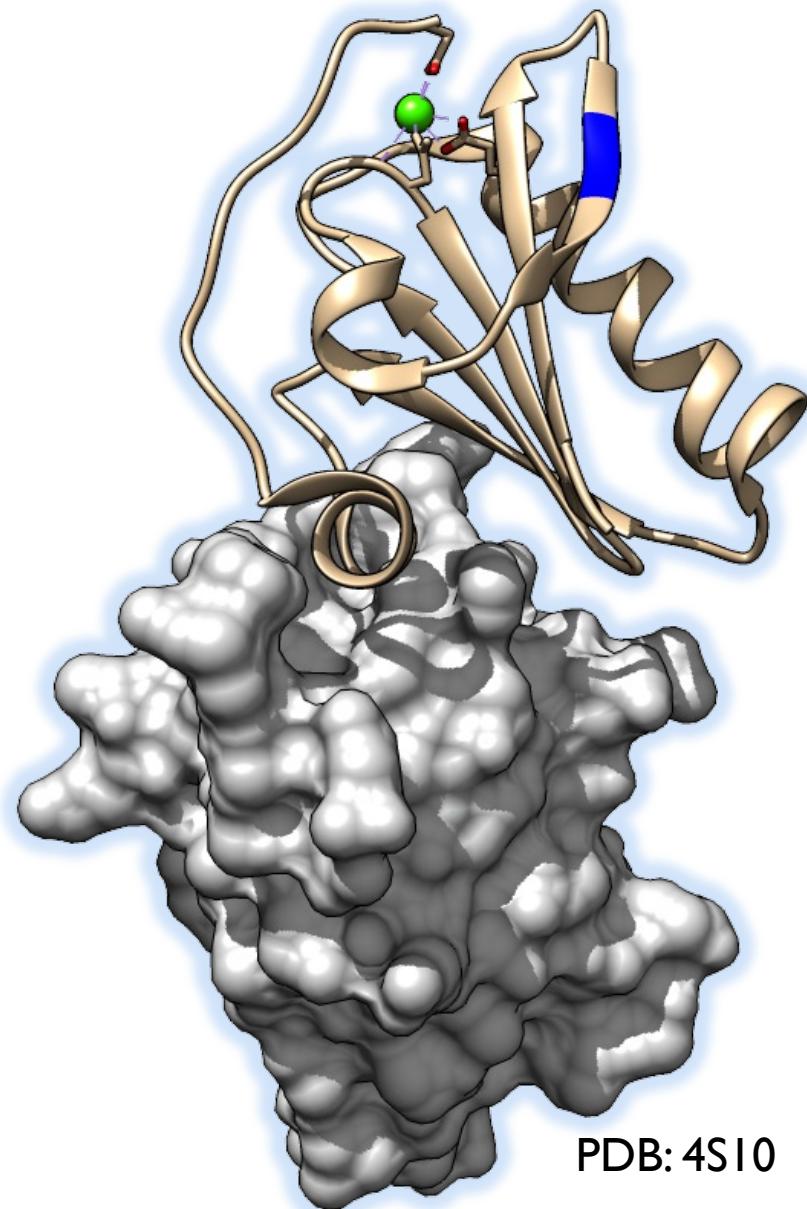
However, crystallisation of Ca^{2+} -bound
isolated domains has been successful.

2015: Gettemans et al. – G2 Nanobodies

- Inoculate llamas with GSN **WT** G2
- Extract nanobodies (NbII)
- Sequence them
- Obtain WT:Nb structure

Idea

- Can the WT-raised NbII re-stabilize D187N enough to allow crystallization?



6H1F: Gelsolin G2+nanobody

[Structure Summary](#)
[3D View](#)
[Annotations](#)
[Experiment](#)
[Sequence](#)
[Genome](#)
[Versions](#)

Biological Assembly 1 



 [3D View: Structure | 1D-3D View](#)

[Electron Density](#) | [Validation Report](#) | [Ligand Interaction](#)

6H1F

 [Display Files](#)  [Download Files](#)

Structure of the nanobody-stabilized gelsolin D187N variant (second domain)

PDB DOI: [10.2211/pdb6H1F/pdb](https://doi.org/10.2211/pdb6H1F/pdb)

Classification: **STRUCTURAL PROTEIN**

Organism(s): *Lama glama*, *Homo sapiens*

Expression System: *Escherichia coli*

Mutation(s): Yes 

Deposited: 2018-07-11 Released: 2019-01-23

Deposition Author(s): [Hassan, A.](#), [Milani, M.](#), [Mastrangelo, E.](#), [de Rosa, M.](#)

Funding Organization(s): Amyloidosis Foundation

Experimental Data Snapshot

Method: X-RAY DIFFRACTION

Resolution: 1.90 Å

R-Value Free: 0.233

R-Value Work: 0.199

R-Value Observed: 0.202

wwPDB Validation

 [3D Report](#)  [Full Report](#)

Metric	Percentile Ranks	Value
Rfree		0.234
Clashscore		6
Ramachandran outliers		0
Sidechain outliers		0
RSRZ outliers		5.2%

Worse  Better 

 Percentile relative to all X-ray structures
 Percentile relative to X-ray structures of similar resolution

This is version 1.0 of the entry. See complete [history](#).

Find Similar Assemblies

Biological assembly 1 assigned by authors and generated by PISA (software)

Biological Assembly Evidence: gel filtration

Macromolecule Content

- Total Structure Weight: 28.49 kDa 
- Atom Count: 1,896 
- Modelled Residue Count: 229 
- Deposited Residue Count: 259 
- Unique protein chains: 2

Literature

[Download Primary Citation](#) 

Nanobody interaction unveils structure, dynamics and proteotoxicity of the Finnish-type amyloidogenic gelsolin variant.

[Giorgino, T.](#), [Matianni, D.](#), [Hassan, A.](#), [Milani, M.](#), [Mastrangelo, E.](#), [Barbiroli, A.](#), [Verhelle, A.](#), [Gettemans, J.](#), [Barzago, M.M.](#), [Diomedede, L.](#), [de Rosa, M.](#)

(2019) *Biochim Biophys Acta Mol Basis Dis* **1865**: 648-660

PubMed: [30625383](https://pubmed.ncbi.nlm.nih.gov/30625383/) [Search on PubMed](#)

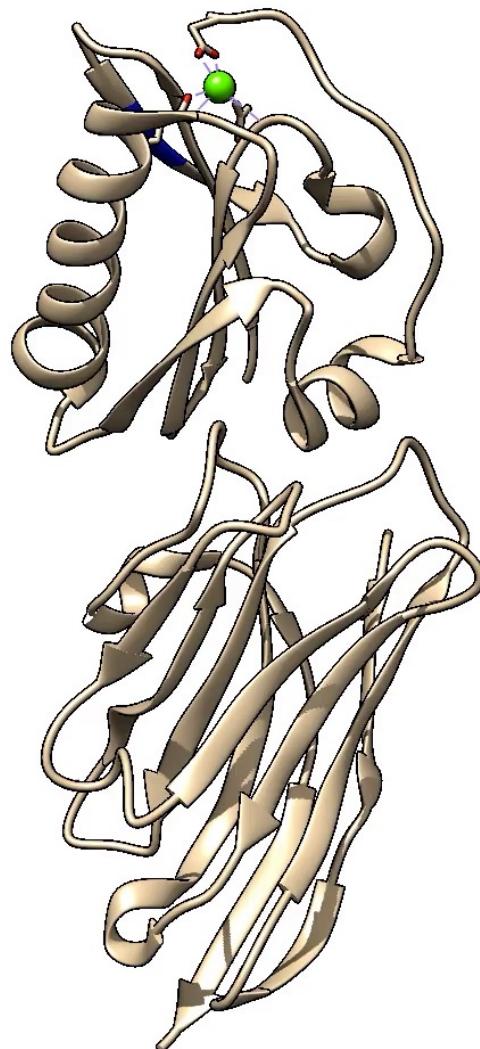
DOI: [10.1016/j.bbadiis.2019.01.010](https://doi.org/10.1016/j.bbadiis.2019.01.010)

Primary Citation of Related Structures:

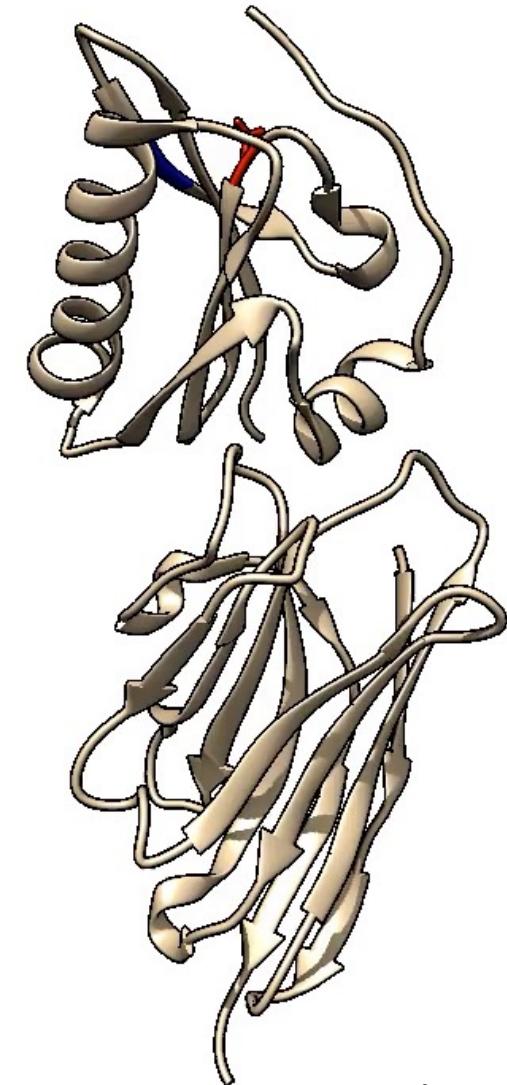
Three puzzles!

WT:NbII complex compared to D187N:NbII.

- I. WT and **D187N** are **virtually identical***: same structure, different function
2. NbII binds far from the furin **cleavage site**...
3. ...and far from the **Ca²⁺** ion



WT: 4S10, 2.6 Å



D187N: **6HIF**, 1.9 Å

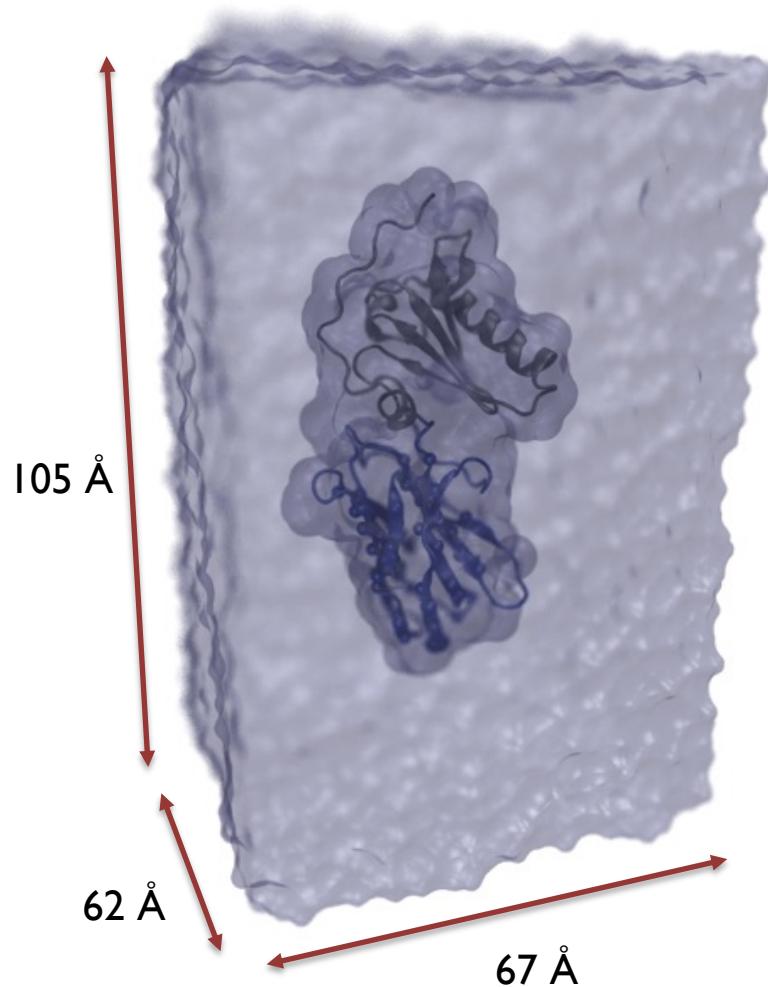
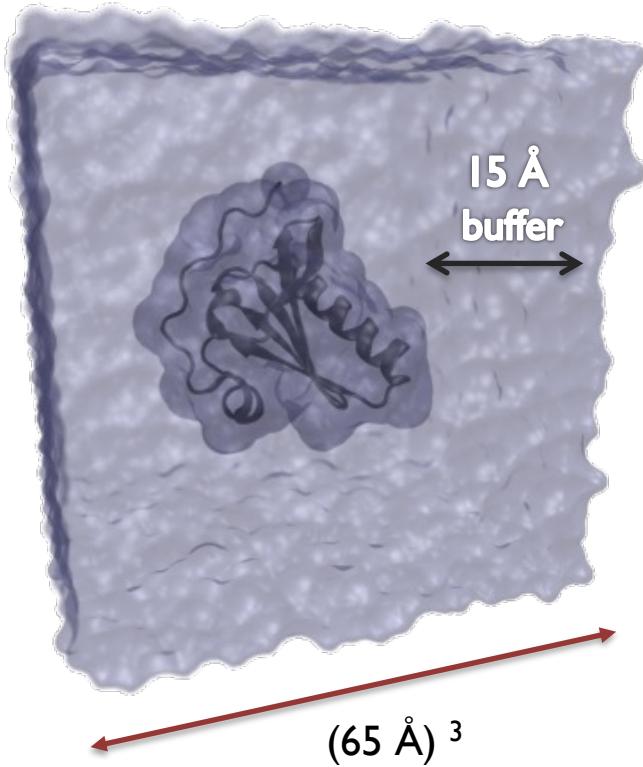
* Except Ca²⁺ binding

GSN \pm NbI MD simulations

- Unbiased sampling @300 °K
- 100 mM NaCl
- Harmonic restraints:
SS NbI @ 0.03 kcal/mol/Å²

CHARMM36

~3 μs tot. ~25k/43k atoms



Open these

- www.rcsb.org/structure/6H1F
- github.com/giorginolab/2024-MoLSim-UniPD

[2024-MoLSim-UniPD](#) / [openmm-python-only](#) / [OpenMM_2024.ipynb](#) 

The code contained in this repo is structure in large part as IPython notebooks. They can be run in three ways:

- On your local machine. That's probably the easiest and most didactic for interpreting the intermediate files.
- On [Google Colab](#). Just open Colab and point it to the file's URL. Simulations run much faster if you request a GPU runtime.
- On Binder. Somewhat slow but comes with a pre-installed environment. [Open in binder](#)

Results

- The results from actual (long) runs are in
github.com/giorginolab/MD-Tutorial-Data

MD-Tutorial-Data / GSN / 

- Next week will do the analysis etc., but we can already visualize and discuss.

Conclusion

Conclusion

- OpenMM is a powerful tool for molecular dynamics simulations
- Good, if fragmented, documentation
- With its customizable force fields and integrators, it can be used to study a wide range of atomistic systems, e.g.
 - “toy” polymers
 - all-atom MD with major FFs
 - ANN potentials

Resources for learning OpenMM

- OpenMM.org website and documentation
- GitHub repository with examples and tutorials
- Community forums and mailing lists for support and discussion
- See also
 - OpenMMtools
 - <https://openforcefield.org/>
 - HTMD, ACEMD
 - <https://github.com/openmm/pdbfixer>
 - Charmm-GUI



RESEARCH ARTICLE

OpenMM 7: Rapid development of high performance algorithms for molecular dynamics

Peter Eastman^{1*}, Jason Swails², John D. Chodera³, Robert T. McGibbon¹, Yutong Zhao¹, Kyle A. Beauchamp^{3a}, Lee-Ping Wang⁴, Andrew C. Simmonett⁵, Matthew P. Harrigan¹, Chaya D. Stern^{3,6}, Rafal P. Wiewiora^{3,6}, Bernard R. Brooks⁵, Vijay S. Pande^{1,7}

¹ Department of Chemistry, Stanford University, Stanford, California, United States of America,

² Department of Chemistry and Chemical Biology and BioMaPS Institute, Rutgers University, Piscataway,

End