

# Resolving Transitions with Markov State Models – and 1000x Faster with AI

Erik Lindahl



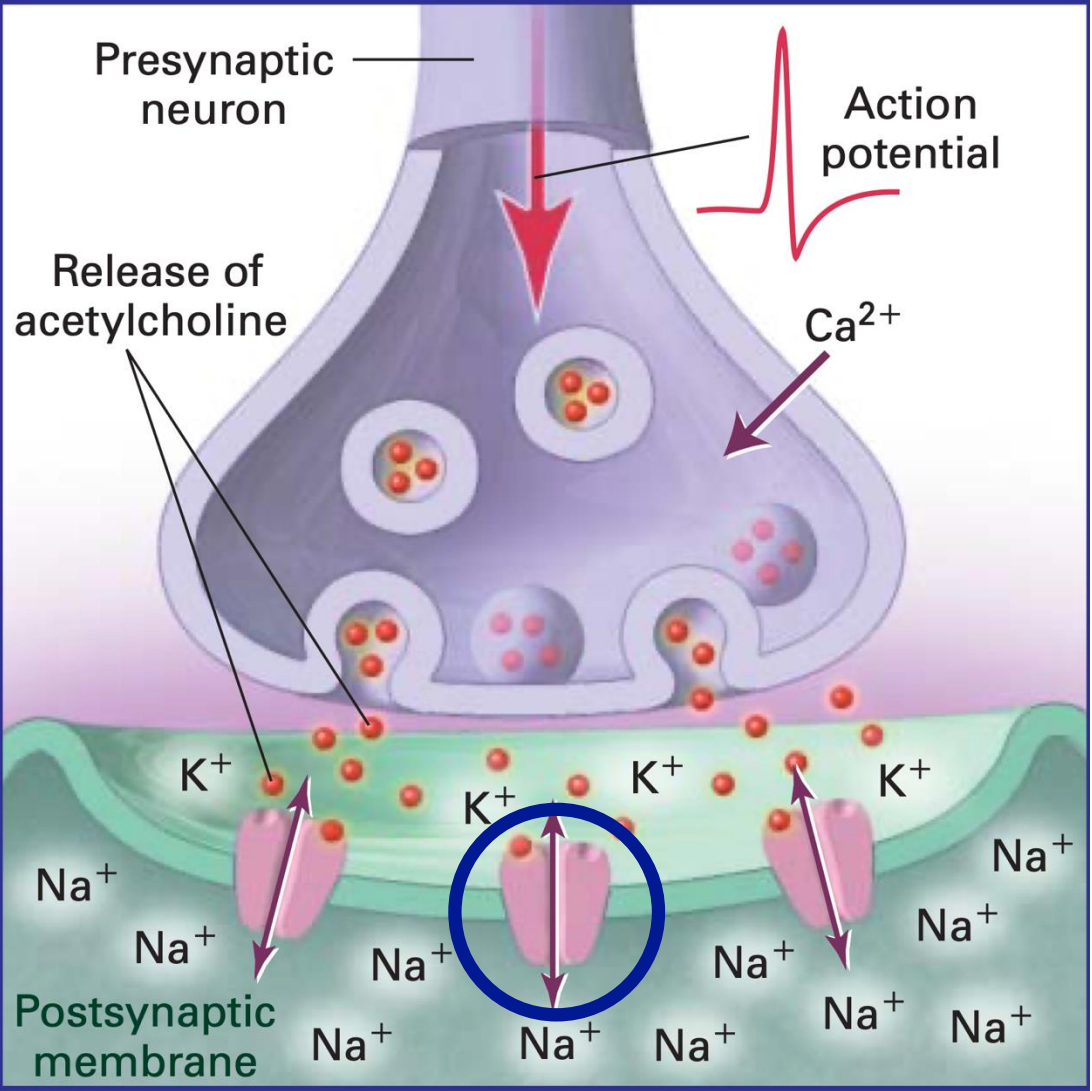
“Suffering so great as I underwent cannot be expressed in words . . . but the blank whirlwind of emotion, the horror of great darkness, and the sense of desertion by God and man, which swept through my mind, and overwhelmed me, I can never forget”



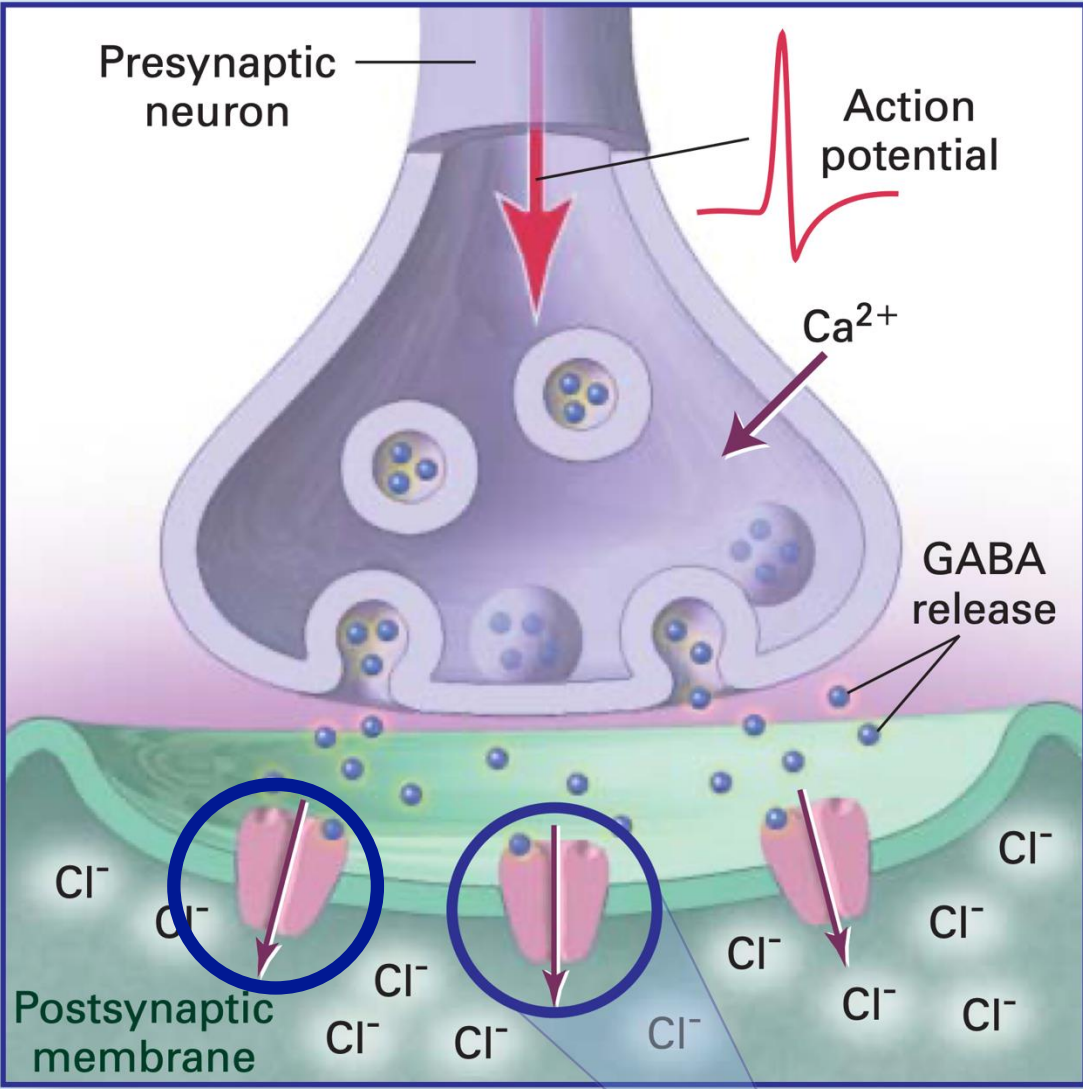
Morton Auditorium, Massachusetts General Hospital  
October 16, 1846



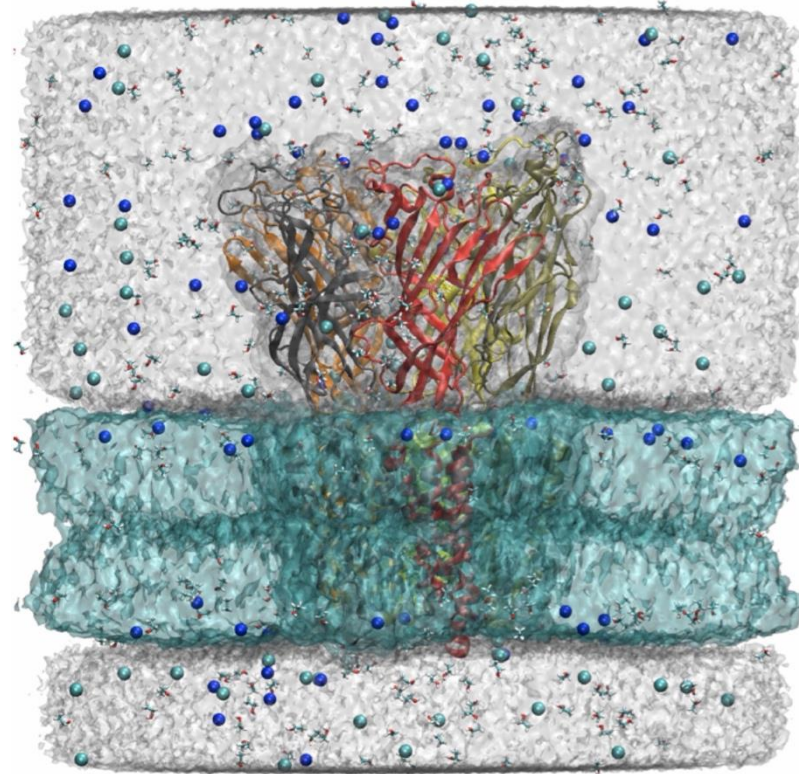
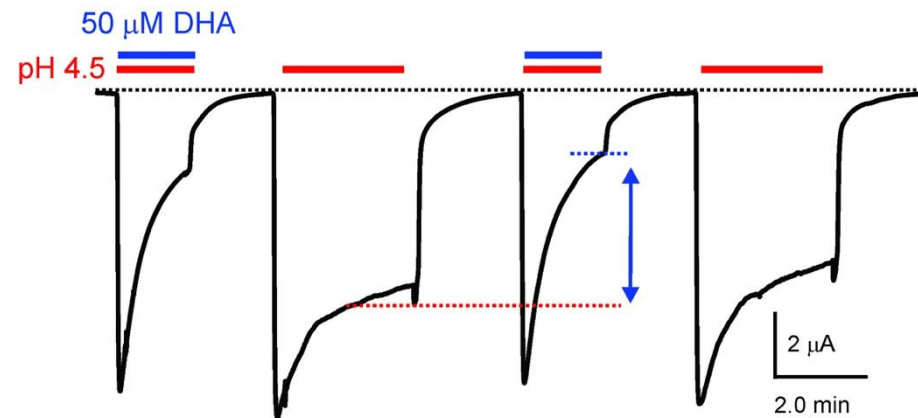
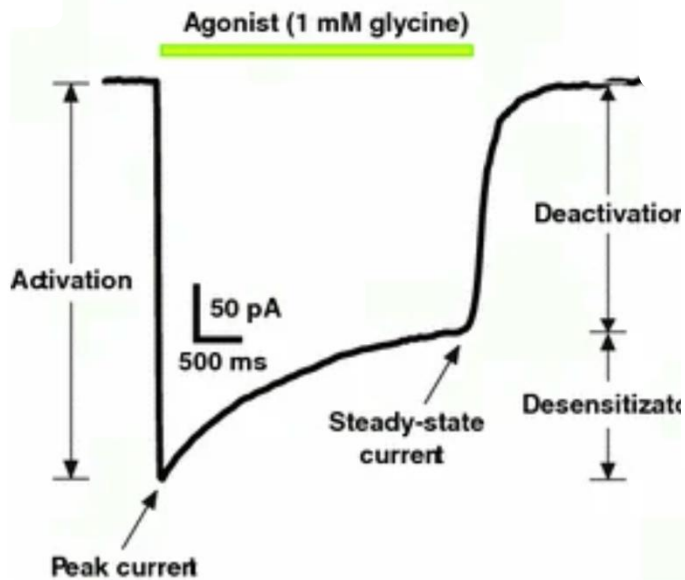
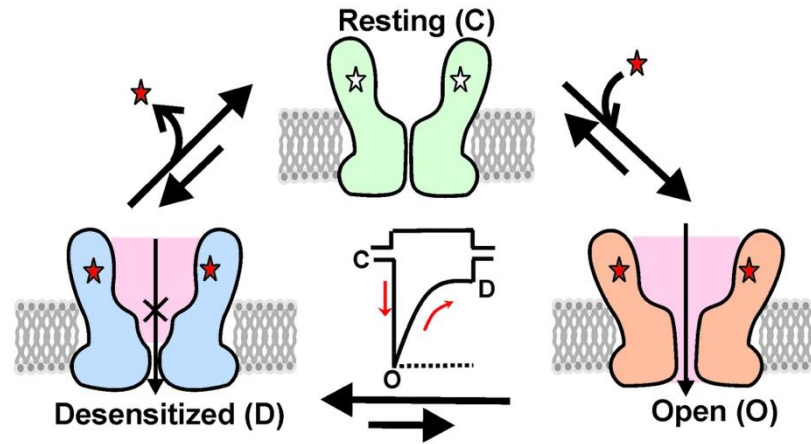
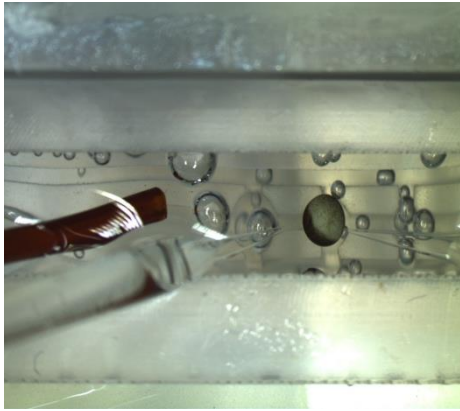
# Neuroexcitation



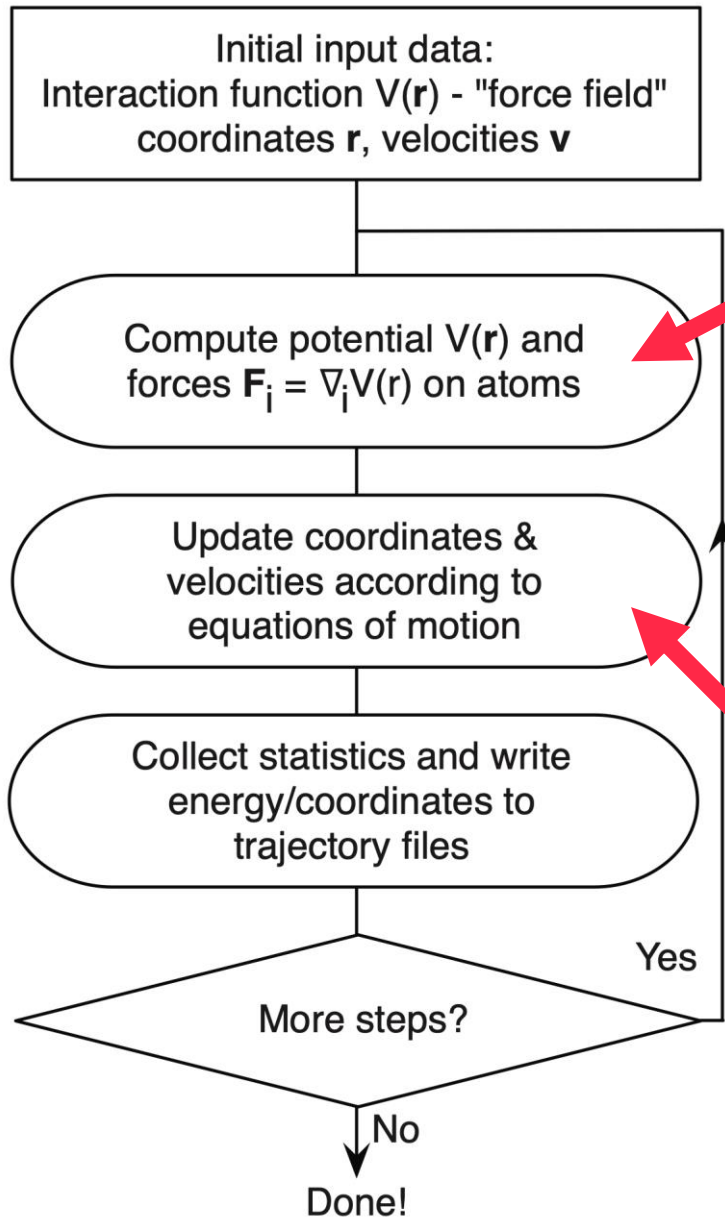
# Neuroinhibition



# MEASURING GATING



# MD is intrinsically parallel, but also sequential



$$\begin{aligned}
 V(r) = & \sum_{bonds} \frac{1}{2} k_{ij}^b (r_{ij} - r_{ij}^0)^2 \\
 & + \sum_{angles} \frac{1}{2} k_{ijk}^\theta (\theta_{ijk} - \theta_{ijk}^0)^2 \\
 & + \sum_{torsions} \left\{ \sum_n k_\theta [1 + \cos(n\phi - \phi_0)] \right\} \\
 & + \sum_{impropers} k_\xi (\xi_{ijkl} - \xi_{ijkl}^0) \\
 & + \sum_{i,j} \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}} \\
 & + \sum_{i,j} \left[ \frac{C_{12}}{r_{ij}^{12}} - \frac{C_6}{r_{ij}^6} \right]
 \end{aligned}$$

Costly, because these terms involve all pairs

$$\begin{aligned}
 m_i \frac{\partial^2 r_i}{\partial t^2} &= F_i \quad i = 1..N \\
 F_i &= - \frac{\partial V(r)}{\partial r_i}
 \end{aligned}$$

With  $\Delta t \sim 1\text{fs}$  and  $\mu\text{s}$  to  $\text{s}$  timescales of interest, we need  $10^9$ - $10^{15}$  steps.

# The floating-point performance generation of HPC scientists

```
for(k=nj0; (k<nj1); k++)
{
  /* Get j neighbor index, and coordinate index */
  jnr = ijnr[k];
  j3 = 3*jnr;

  /* load j atom coordinates */
  jx1 = pos[j3+0];
  jy1 = pos[j3+1];
  jz1 = pos[j3+2];

  /* Calculate distance */
  dx11 = ix1 - jx1;
  dy11 = iy1 - jy1;
  dz11 = iz1 - jz1;
  rsq11 = dx11*dx11+dy11*dy11+dz11*dz11;

  /* Calculate 1/r and 1/r2 */
  rin11 = 1.0/sqrt(rsq11);

  /* Load parameters for j atom */
  qq = iq*charge[jnr];
  tj = nti+2*type[jnr];
  c6 = vdwparam[tj];
  c12 = vdwparam[tj+1];
  rinvsq = rin11*rin11;

  /* Coulomb interaction */
  vcoul = qq*rin11;
  vctot = vctot+vcoul;

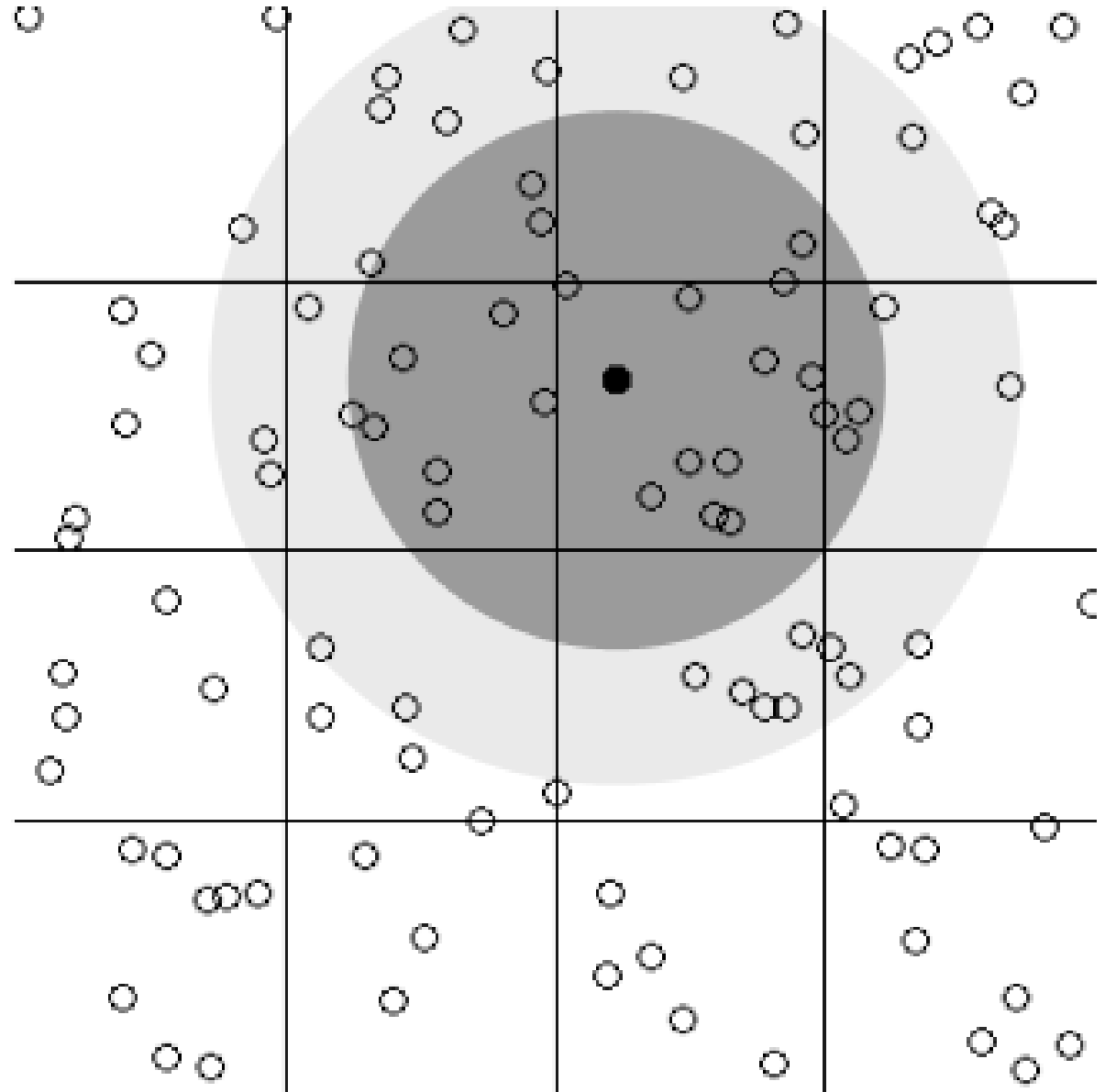
  /* Lennard-Jones interaction */
  rinvsix = rinvsq*rinvsq*rinvsq;
  Vvdw6 = c6*rinvsix;
  Vvdw12 = c12*rinvsix*rinvsix;
  Vvdwtot = Vvdwtot+Vvdw12-Vvdw6;
  fscal = (vcoul+12.0*Vvdw12-6.0*Vvdw6)*rinvsq;

  /* Calculate temporary vectorial force */
  tx = fscal*dx11;
  ty = fscal*dy11;
  tz = fscal*dz11;

  /* Increment i atom force */
  fix1 = fix1 + tx;
  fiy1 = fiy1 + ty;
  fiz1 = fiz1 + tz;

  /* Decrement j atom force */
  faction[j3+0] = faction[j3+0] - tx;
  faction[j3+1] = faction[j3+1] - ty;
  faction[j3+2] = faction[j3+2] - tz;

  /* Inner loop uses 38 flops/iteration */
}
```





ATI!  
PN7120009000  
0414C

ATI RADEON

PN 109-426100-0

# STREAM COMPUTING ON GRAPHICS HARDWARE

A DISSERTATION

SUBMITTED TO THE DEPARTMENT OF COMPUTER SCIENCE

AND THE COMMITTEE ON GRADUATE STUDIES

OF STANFORD UNIVERSITY

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

DOCTOR OF PHILOSOPHY

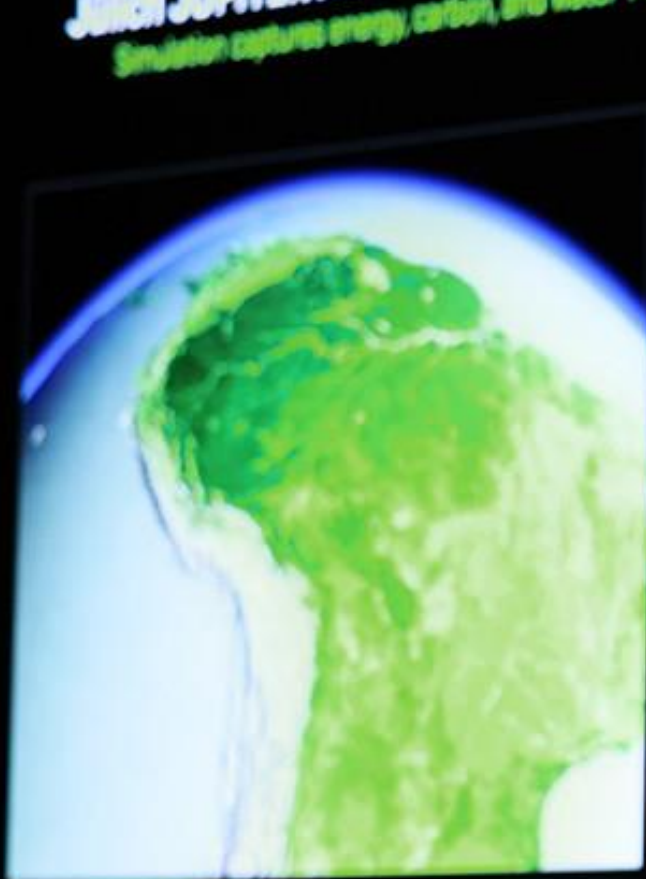
Ian Buck

September 2006

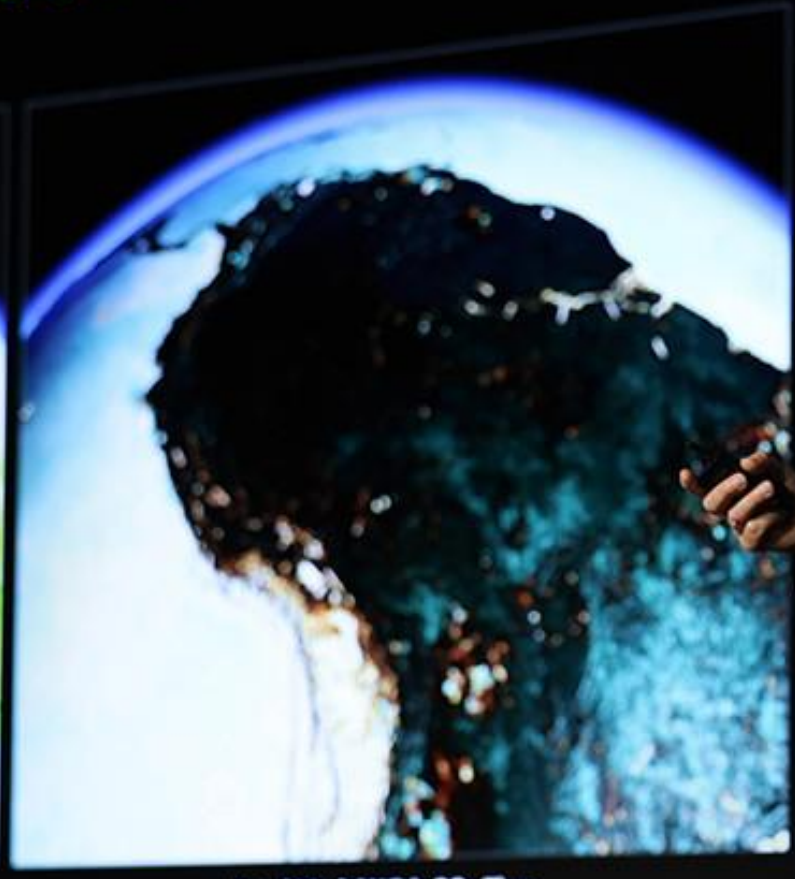
Ian Buck, vice president and general manager of accelerated computing at NVIDIA, delivers a special address at SC25

# Jülich JUPITER Achieves Largest and Most Complex Climate Simulation

Simulation captures energy, carbon, and water 146 days per day on 23,480 Grace Hopper Superchips



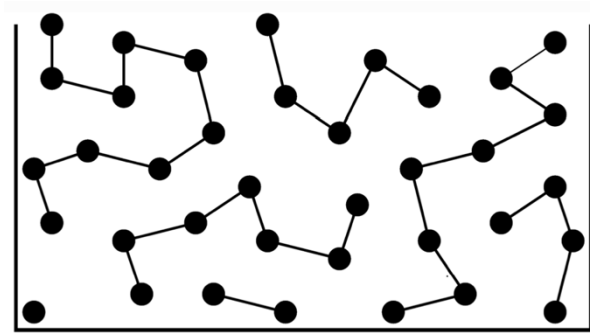
Land Model With Vegetation



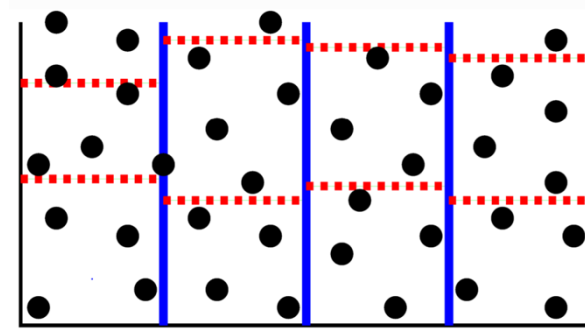
Land Model With CO<sub>2</sub> Flux



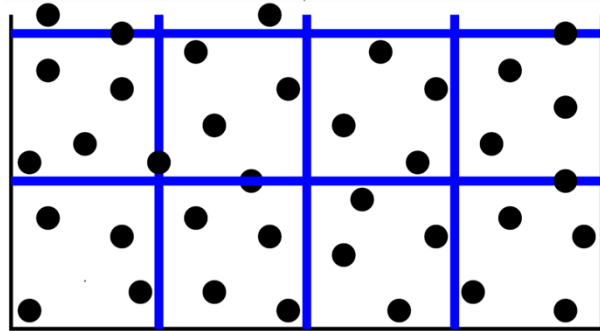
# Don't optimize your **code** for GPUs - optimize the *algorithm*



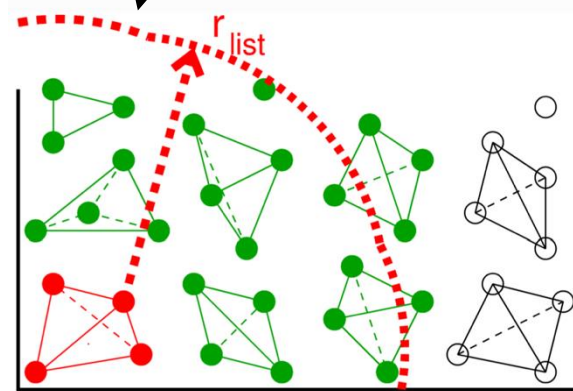
x,y grid  
z sort  
z bin



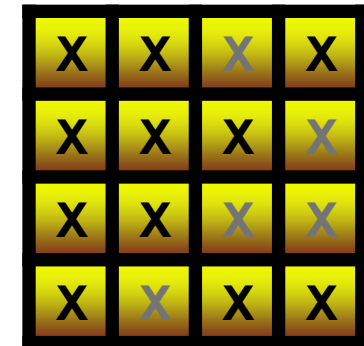
x,y,z  
gridding



Cluster pairlist



Organize  
as tiles with  
all-vs-all  
interactions:



- $i=3$ : 5 6 9 12 15 17 18 25 32 ...
- $i=4$ : 7 8 9 11 12 15 17 25 32 43 54 ...
- ... 8 9 10 11 12 13 19 20 ...

Tile interaction algorithms:  
Load  $N$  atoms, compute  $N^2$  forces

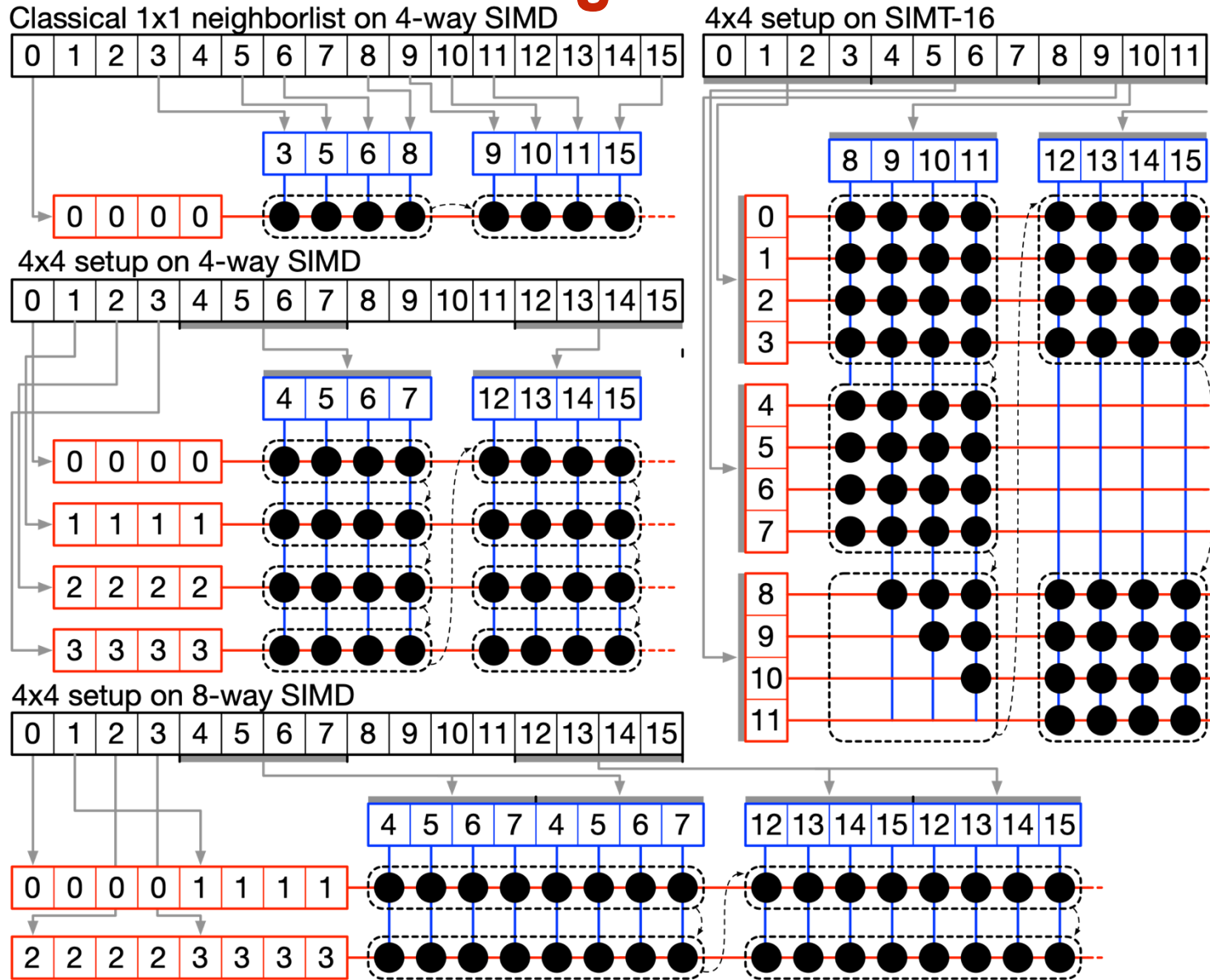
The Link-cell algorithm: Load 1 atom, calculate 1 interaction  
Verlet, Phys Rev 159, 98-103 (1967)

**TO FIRST APPROXIMATION, GPUS PROVIDE  
AN INFINITE AMOUNT OF FLOPS**

**OUR FIRST JOB IS TO HANDLE  
MEMORY BANDWIDTH**

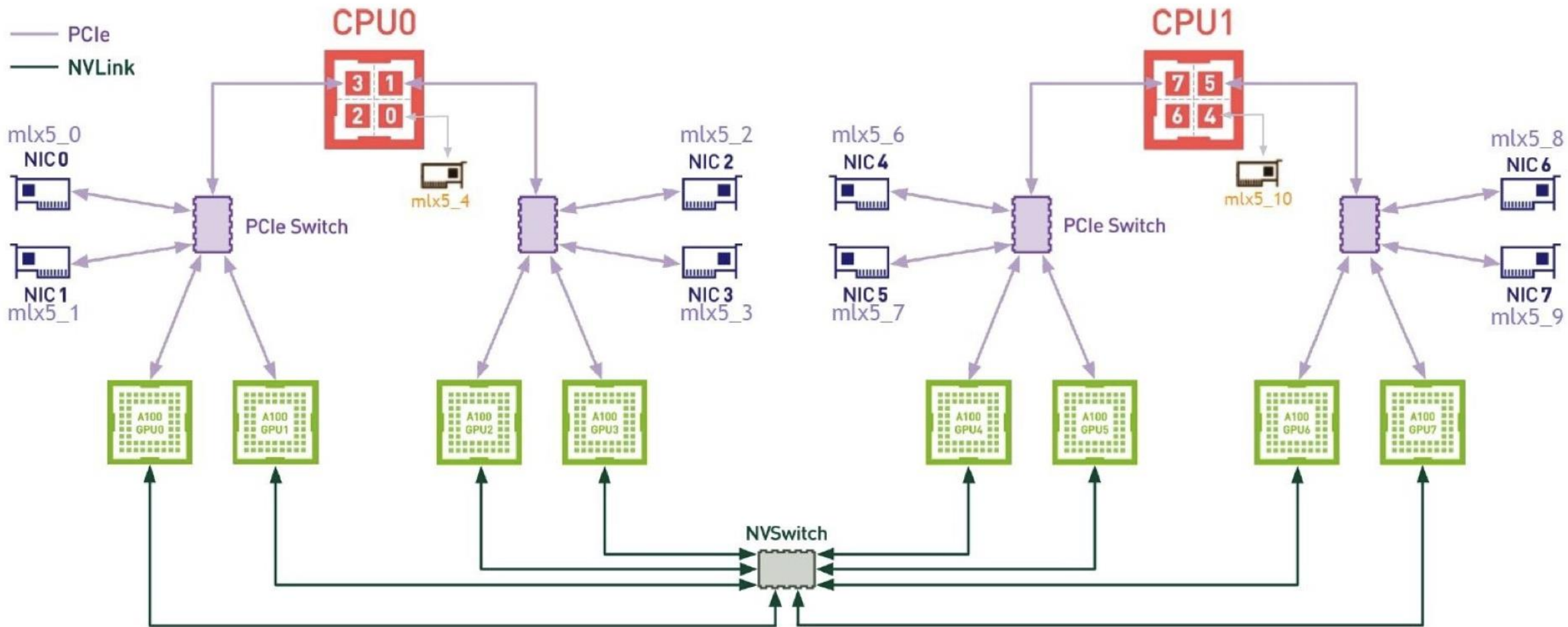
# We change the algorithms to improve compute/memory ratio

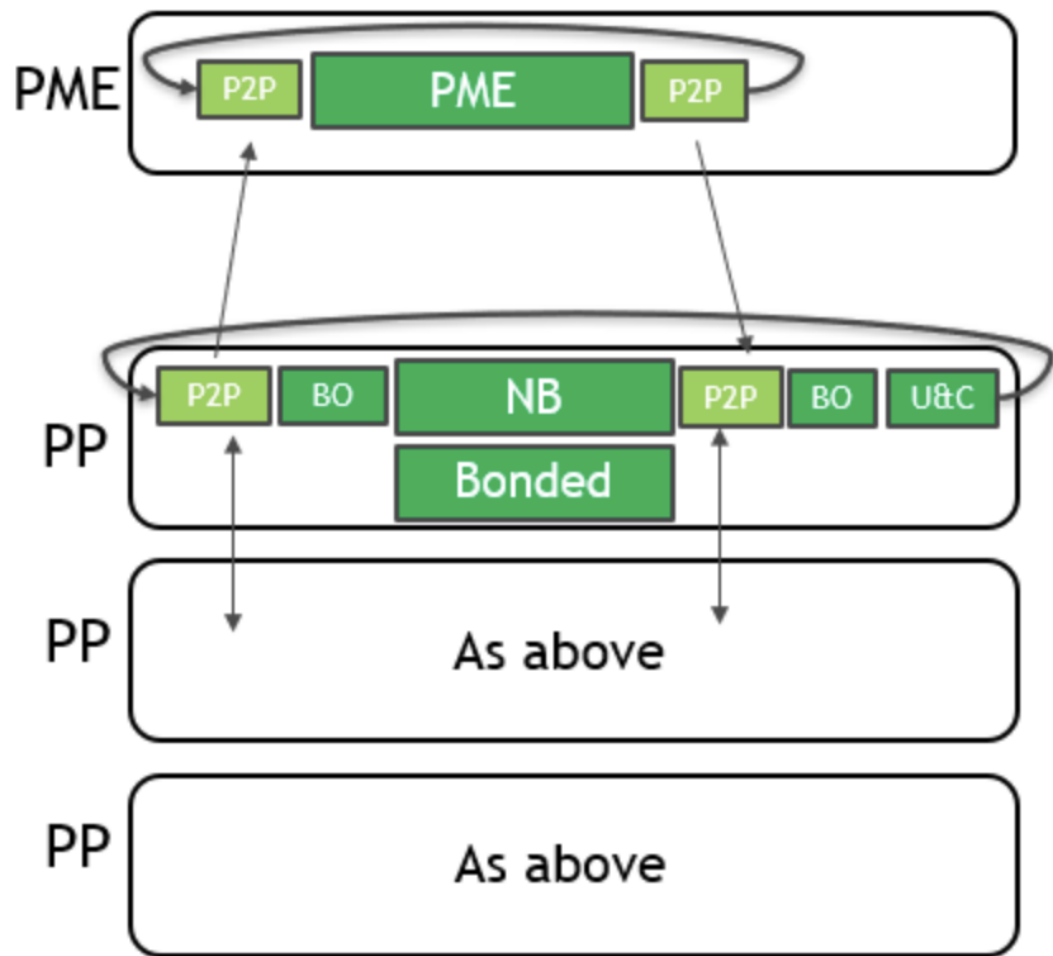
## This is great for modern CPUs too



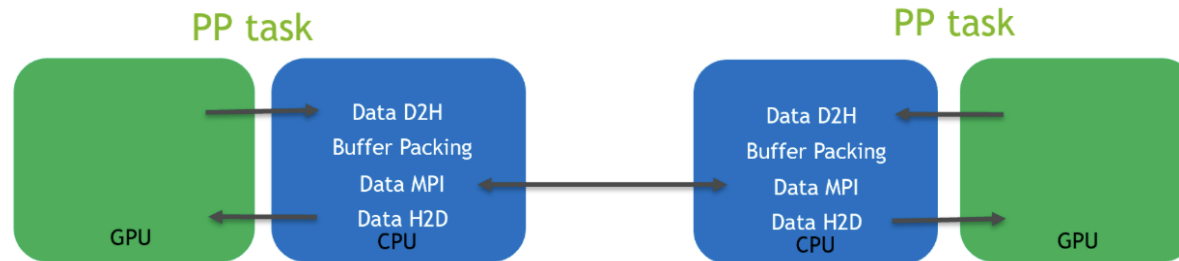
- CUDA
- OpenCL
- Intel MIC
- x86 SSE2
- x86 SSE4.1
- x86 AVX
- x86 AVX-128-FMA
- x86 AVX2
- x86 AVX2\_128
- x86 AVX-512F
- x86 AVX-512ER
- Arm Neon
- Arm64 Asimd
- IBM QPX
- IBM VMX
- IBM VSX
- Fujitsu HPC-ACE
- ARM SVE & SVE2
- SYCL
- HIP







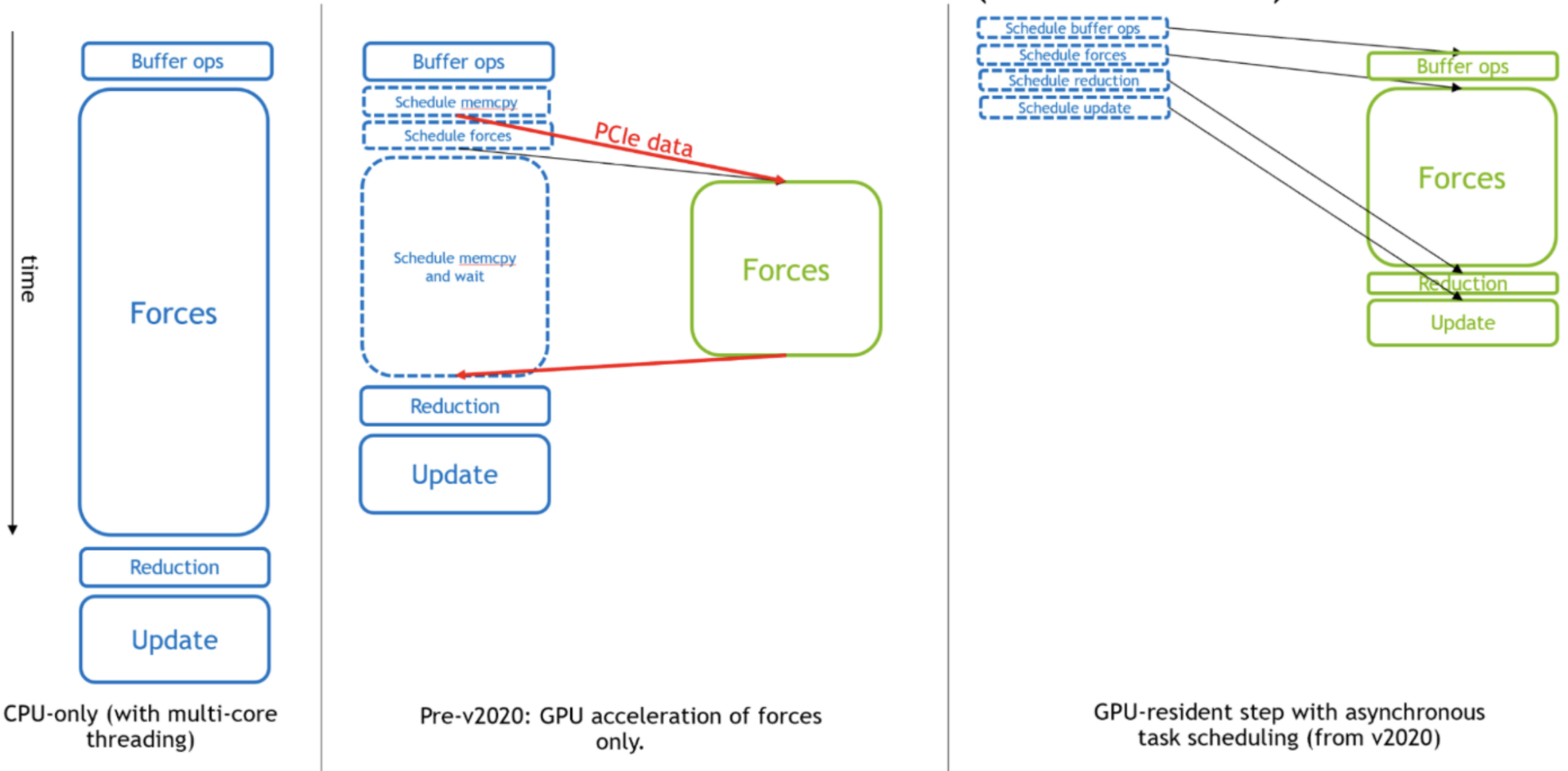
### MPI staged through CPUs



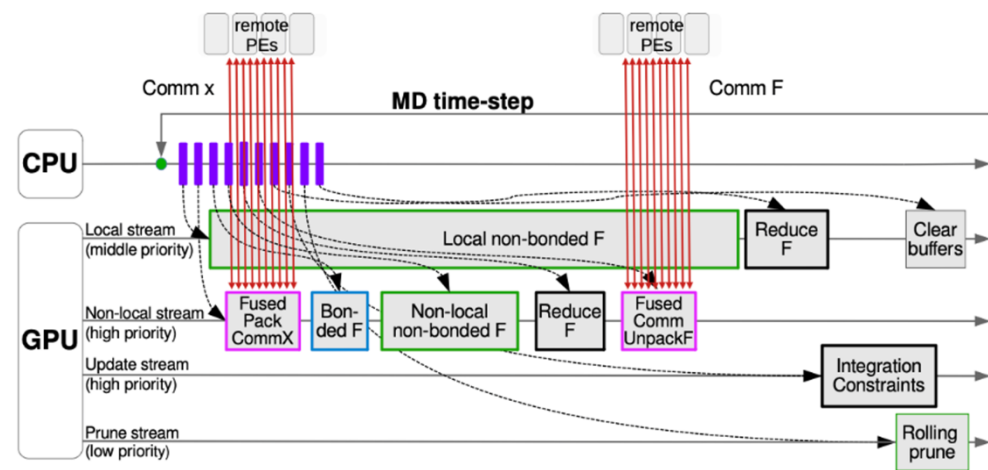
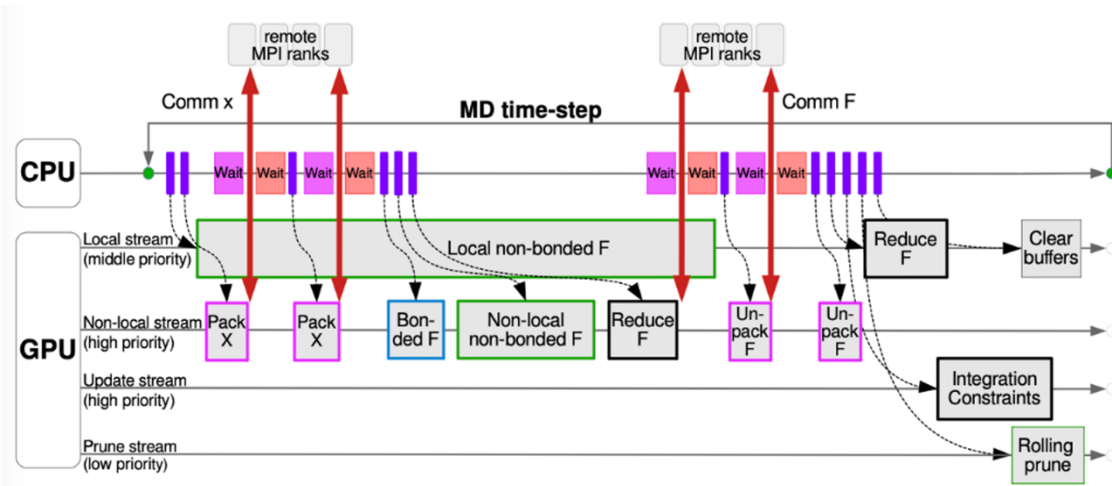
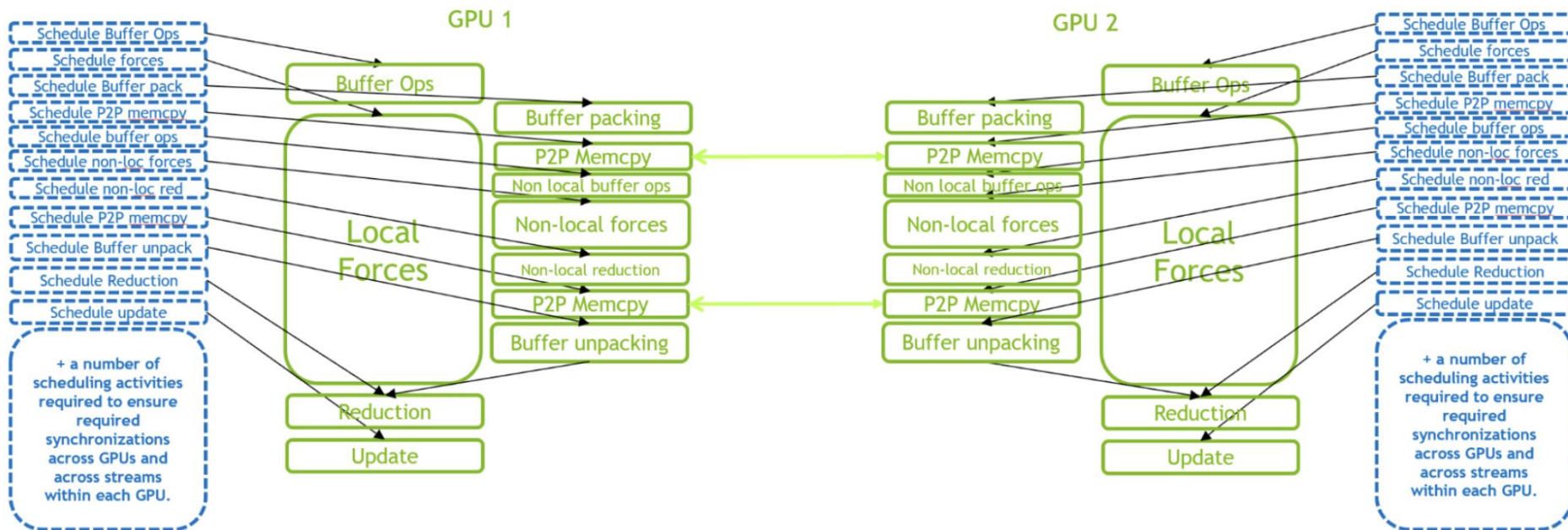
### GPU-direct halo communication



# From CPU to offloading to GPU-resident CPU-scheduled kernels

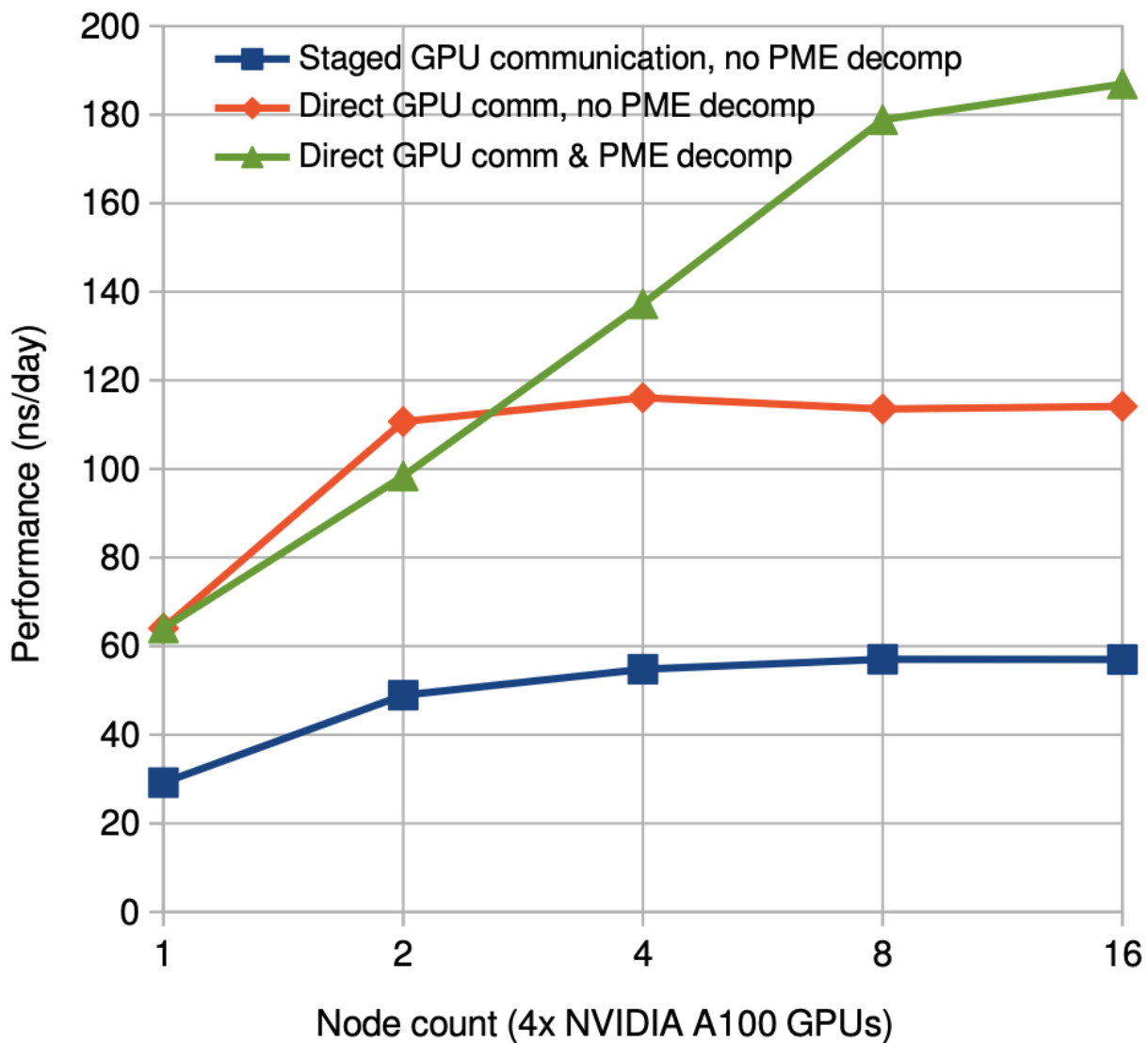


# Multiple GPUs & Multiple nodes

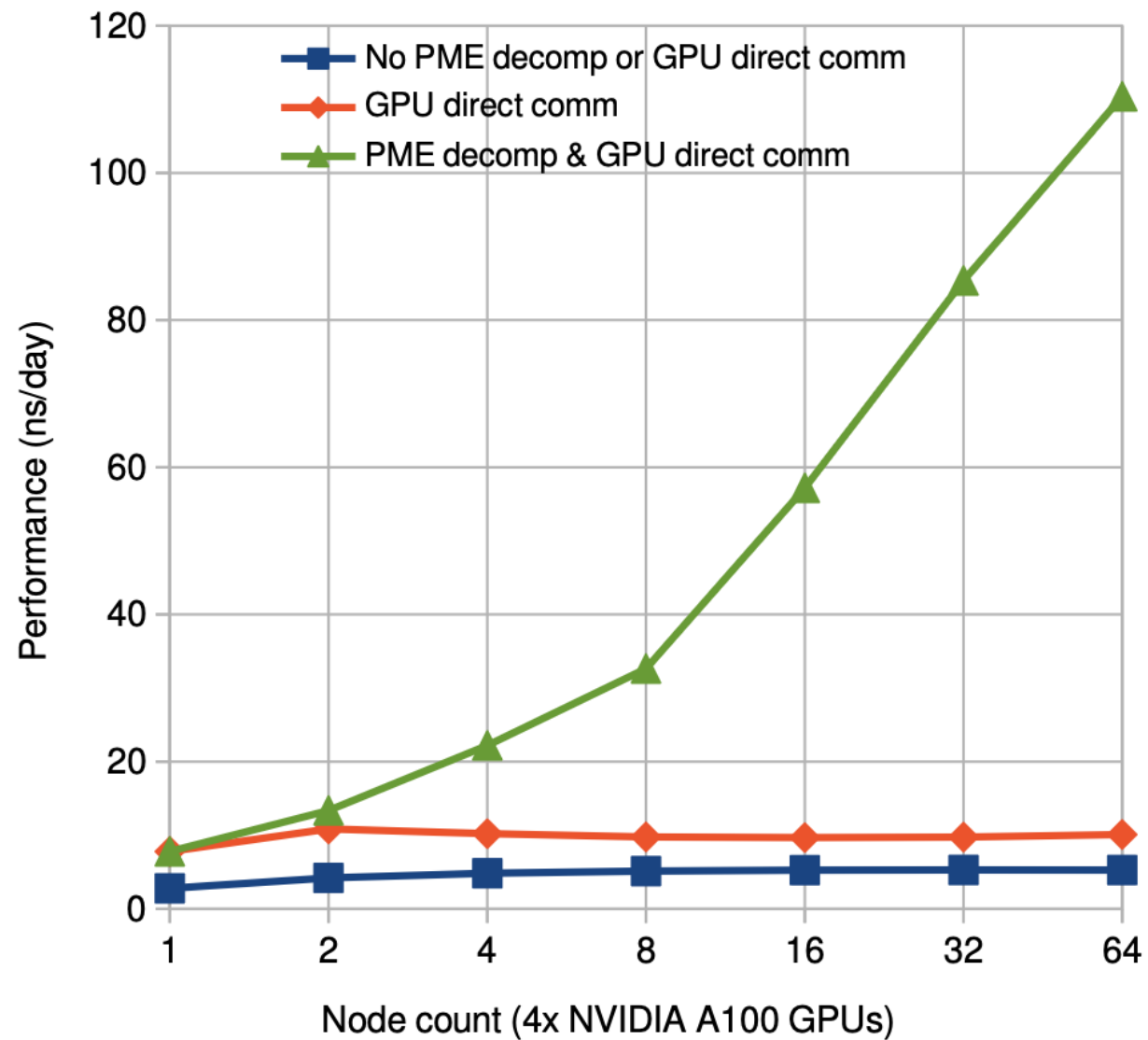


# Efficient strong scaling down to ~25,000 atoms per GPU

## STMV 1M atoms

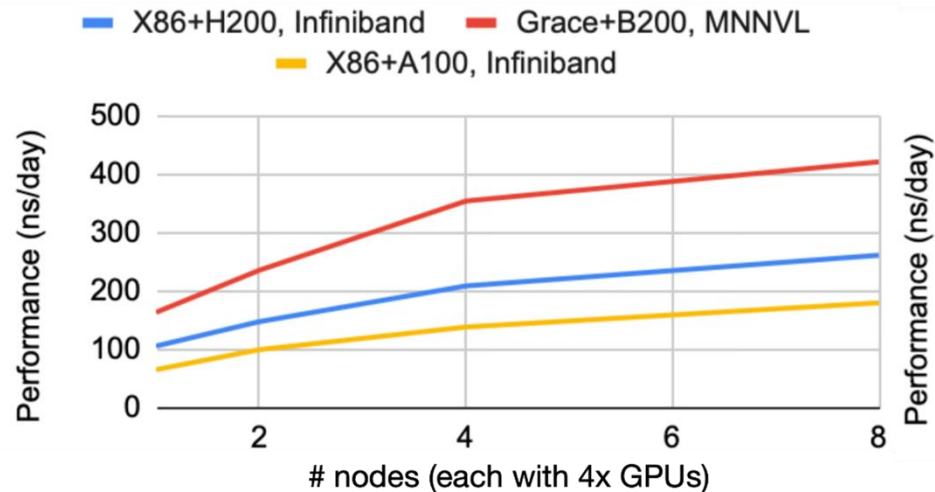


## benchPEP-h 12M atoms

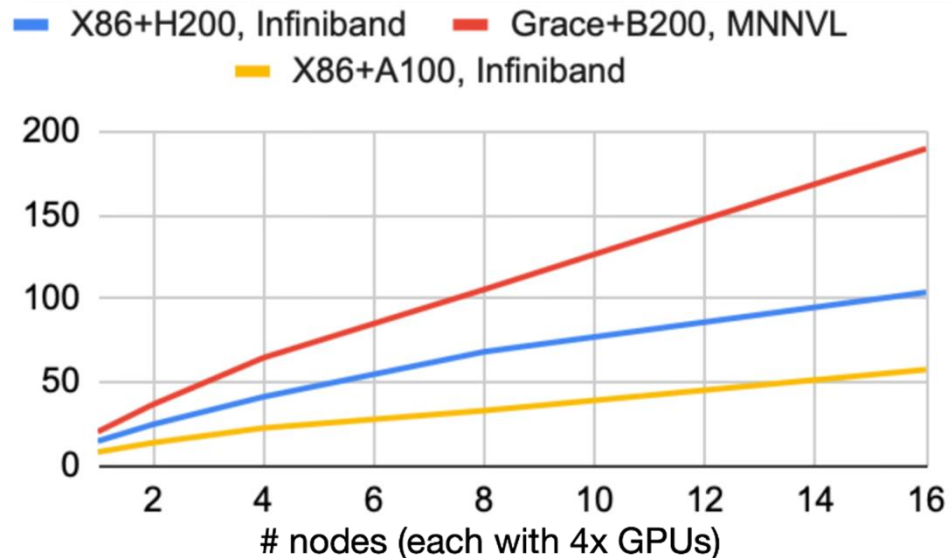


# Tightly connected expensive nodes are not only good for AI!

## STMV (1M atoms)

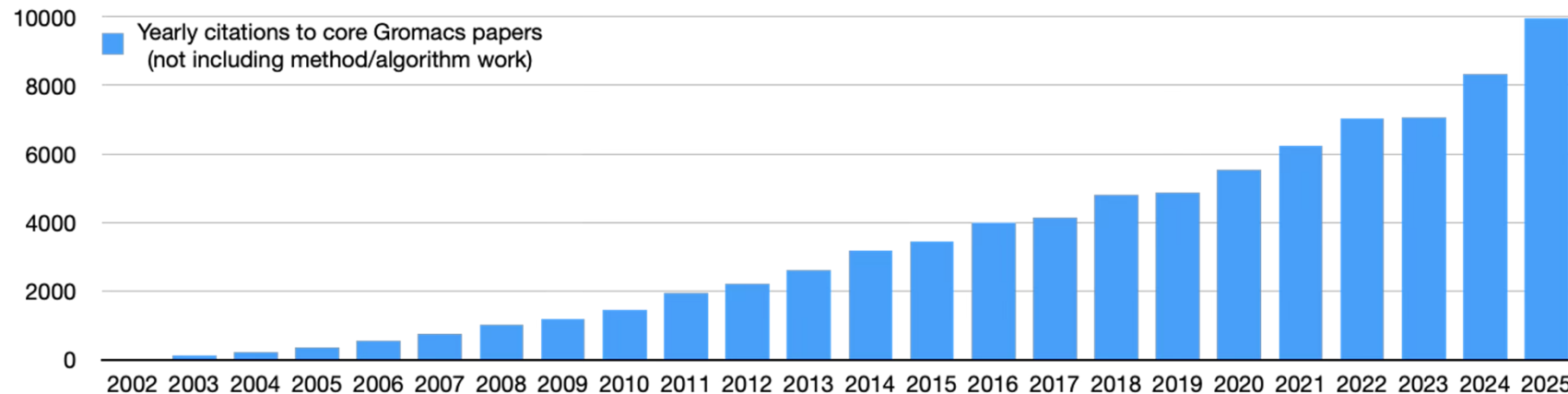
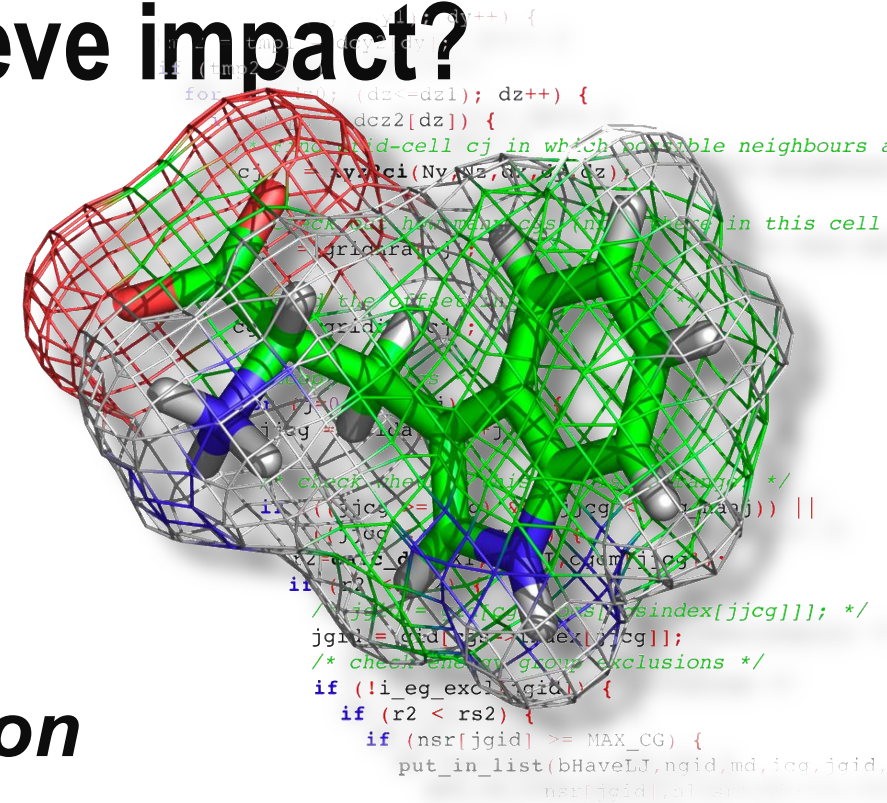


## BenchPep (12M atoms)



# In retrospect: How does software achieve impact?

- Usability, quality, performance, portability
- You must build a broad *external* user base
- Support your users: blood, sweat & tears
- Just as for research, you first do the work, invest time & effort, show promising early impact on others - and then get funding.
- **"Life's most persistent and urgent question is, 'What are you doing for others?'"**

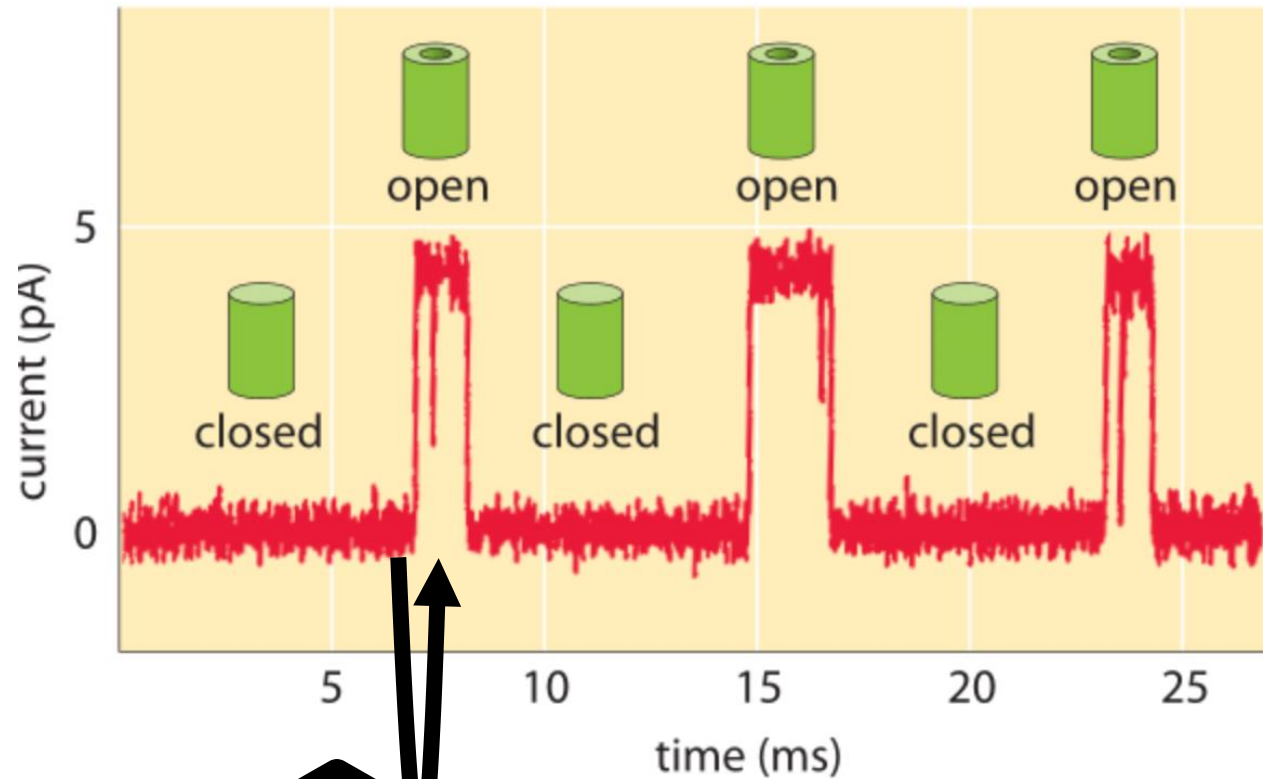


**26,000 FPS**

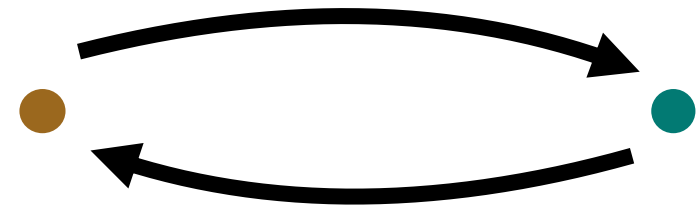


# **MODELING CONFORMATIONAL CHANGES WITH (MANY) SIMULATIONS**

# QUANTIFYING SLOW GATING IN SIMULATIONS: MARKOV STATE MODELS



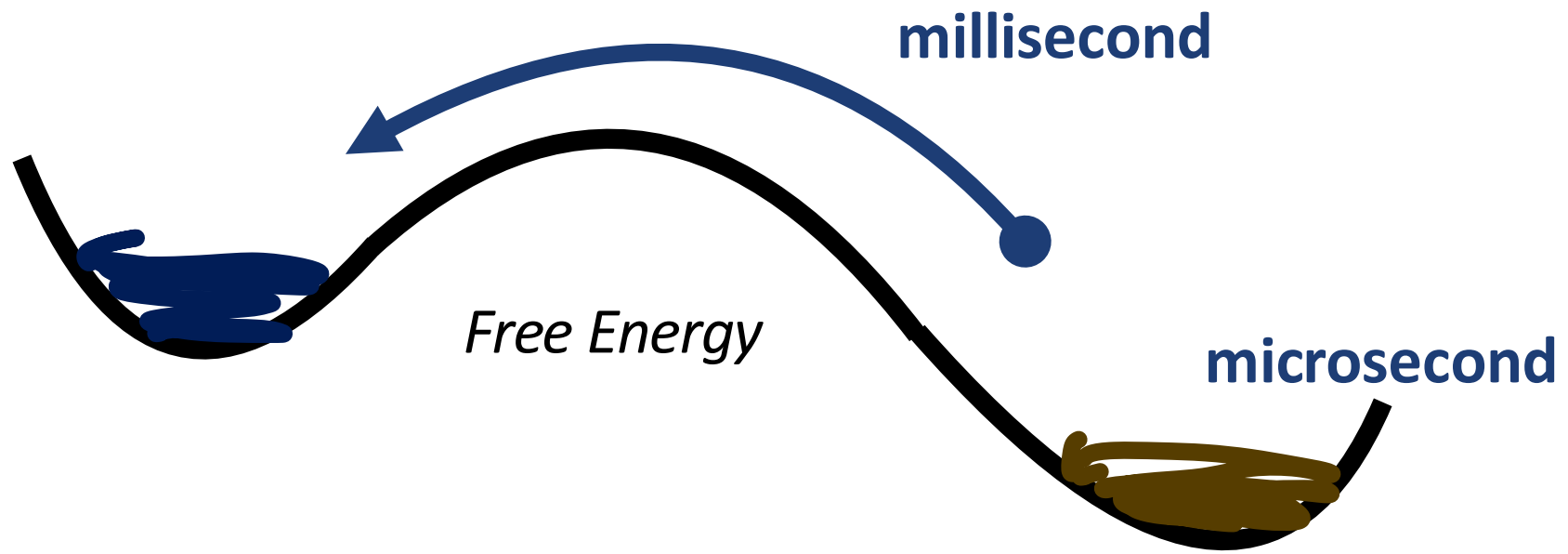
Kinetic model



Open



Closed

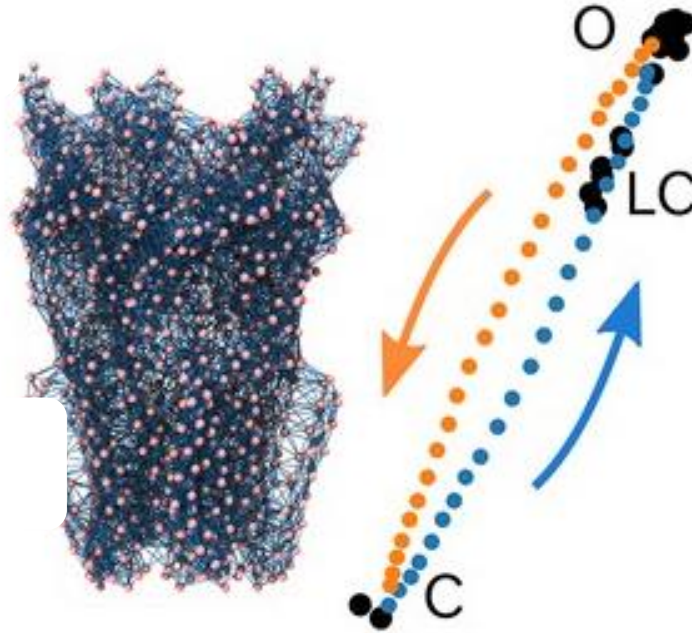
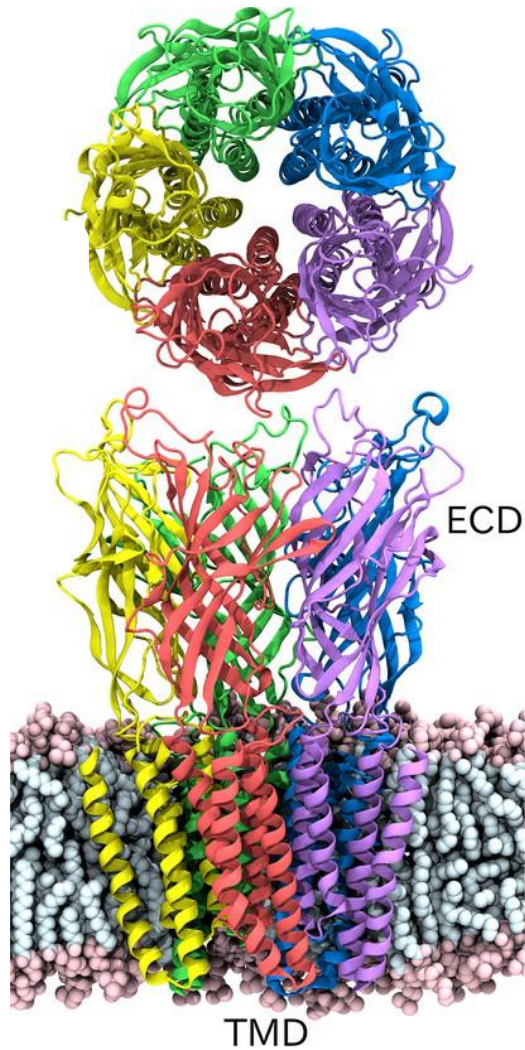


millisecond

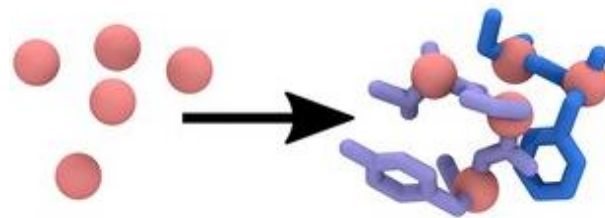
*Free Energy*

microsecond

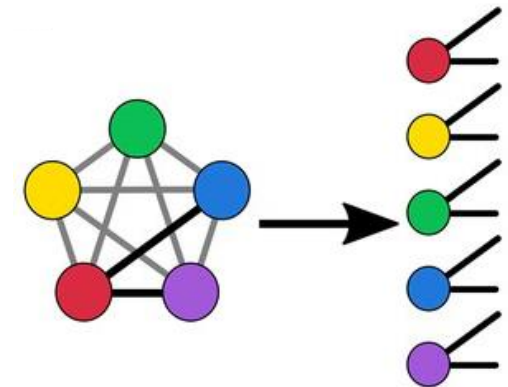
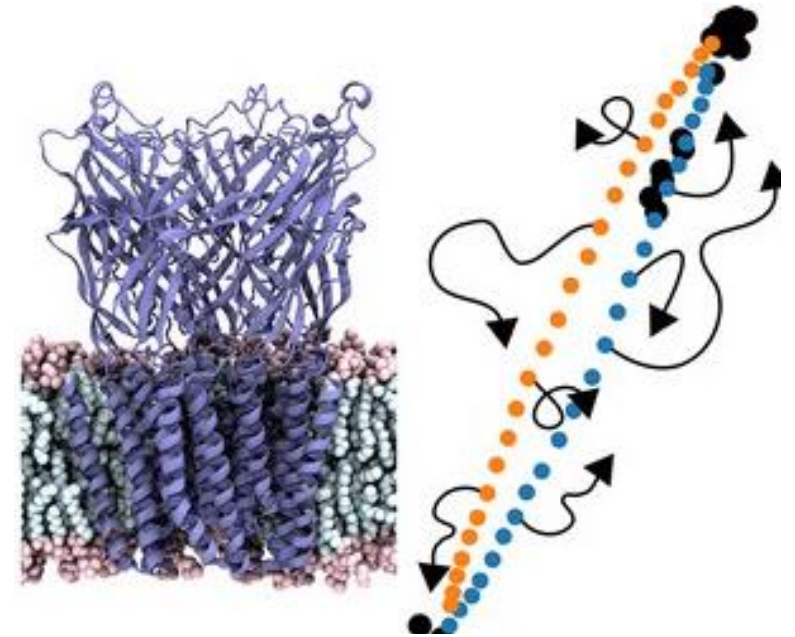
# QUANTIFYING THE GATING PROCESS WITH MARKOV STATE MODELS



<https://ebdims.biophysics.se/>

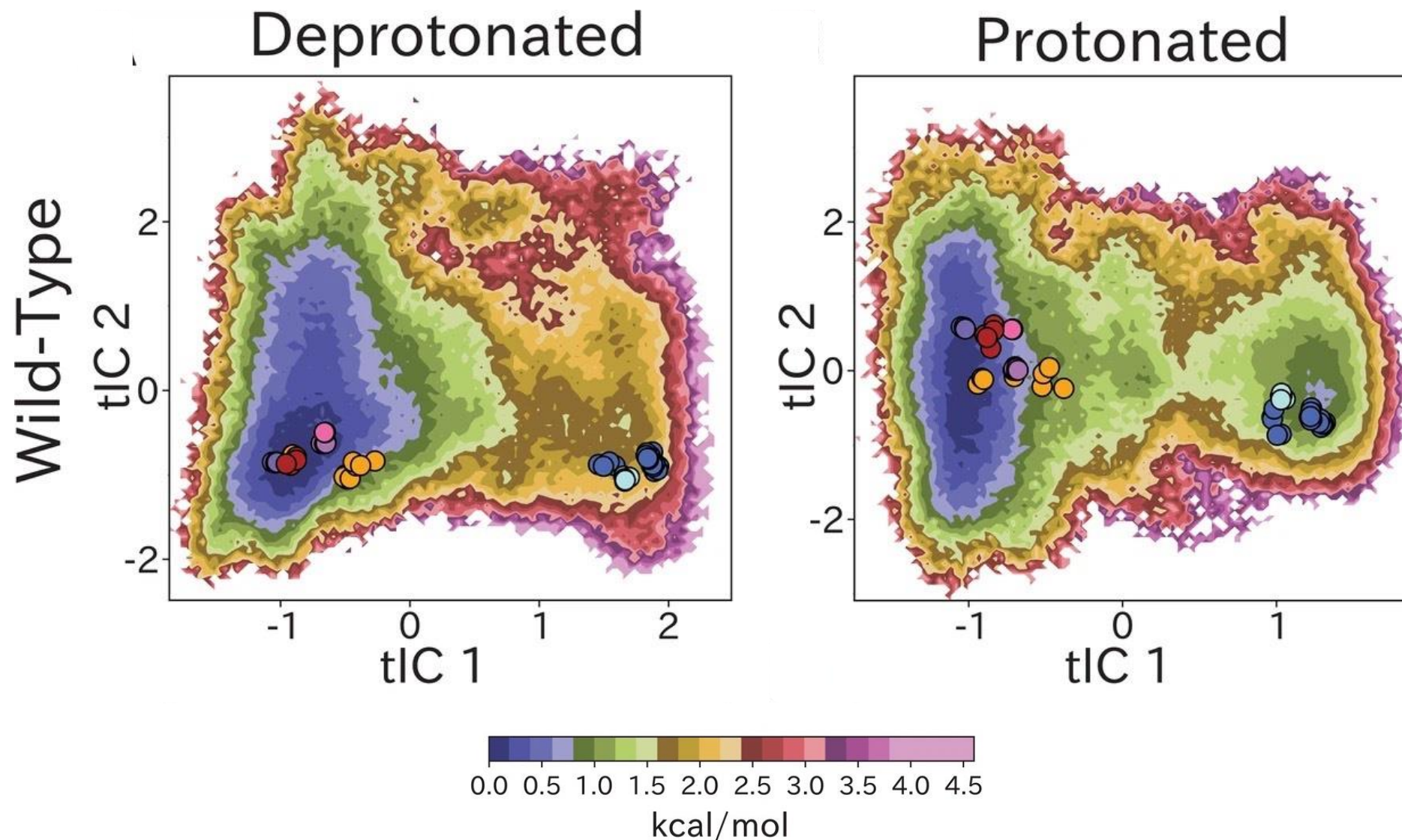


Laura Orellana, Nature Comm. 7, 12575 (2016)

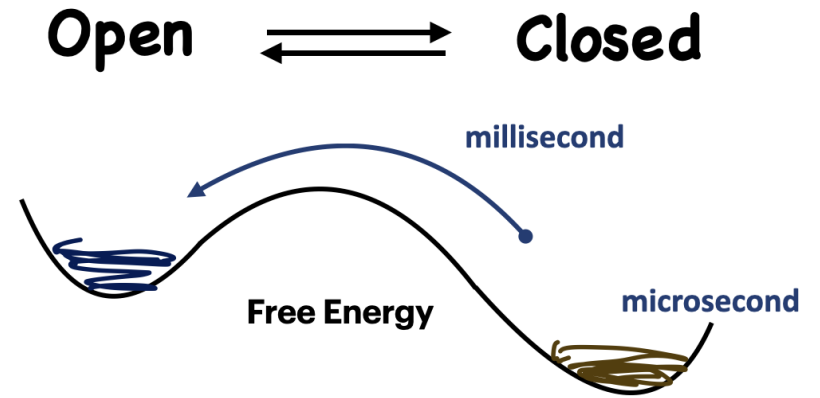
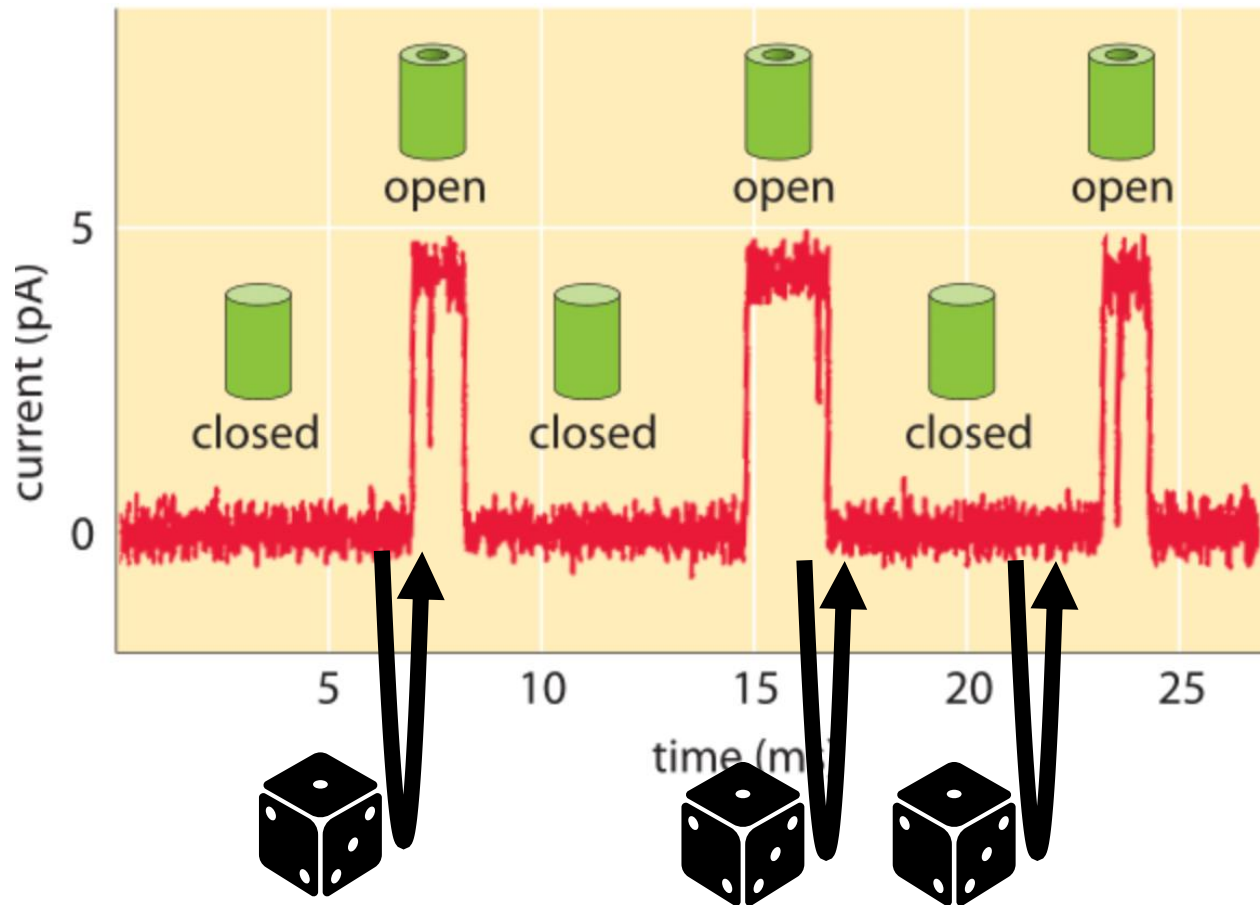


Cathrine Bergh, eLife 10:e68369 (2021)

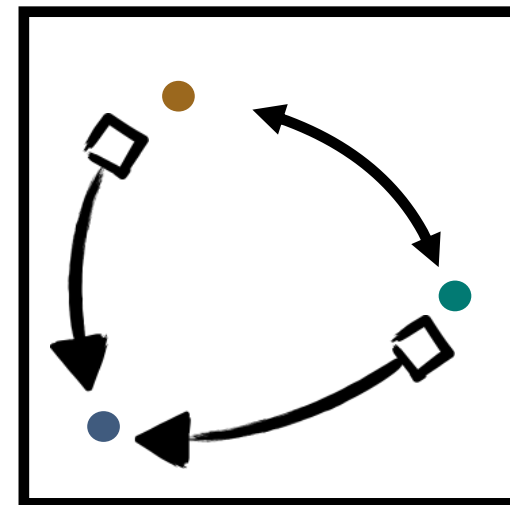
# MSM FREE ENERGY LANDSCAPES INDICATE LOW OPEN PROBABILITIES



# MODELING A MORE REALISTIC CONFORMATIONAL CYCLE

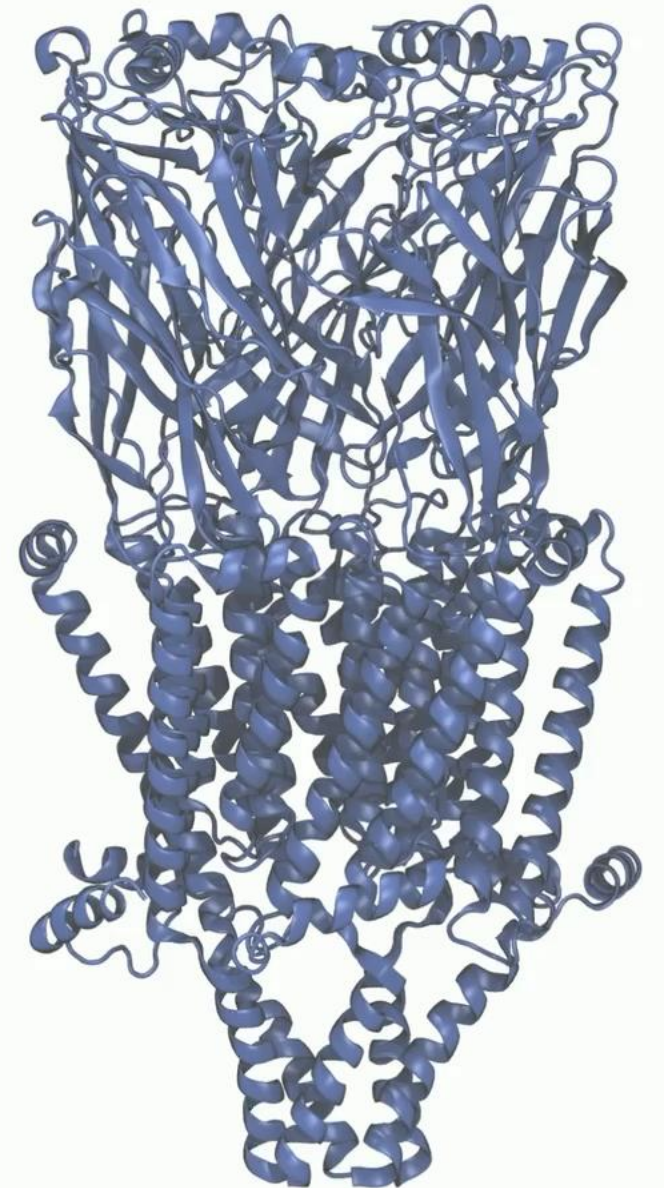
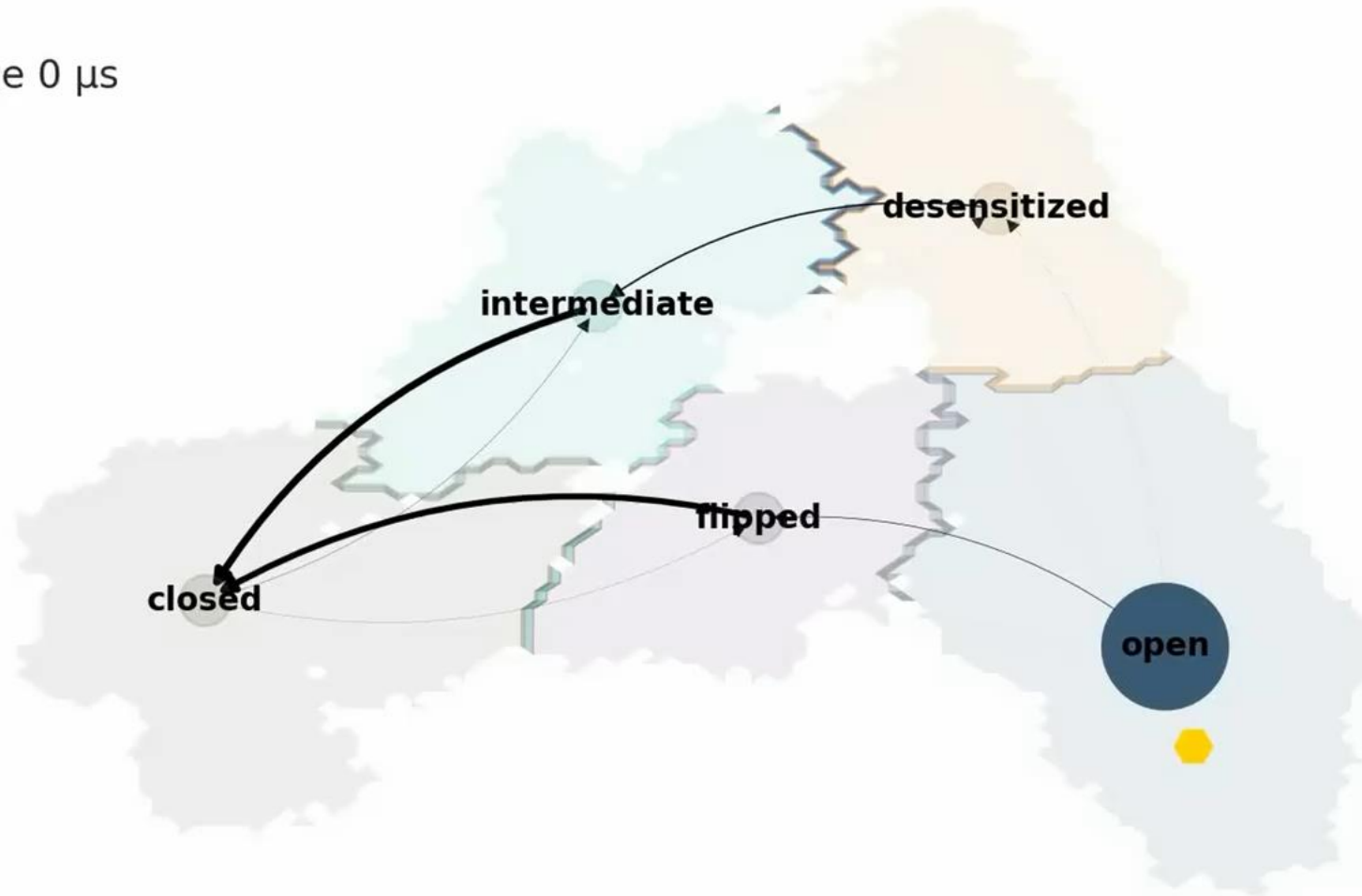


*Kinetic model  
including desensitization*



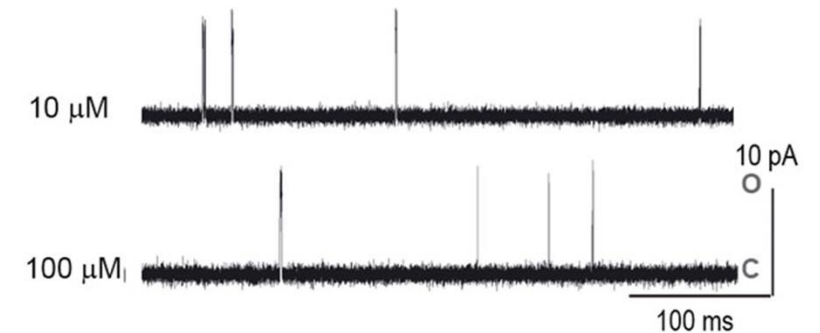
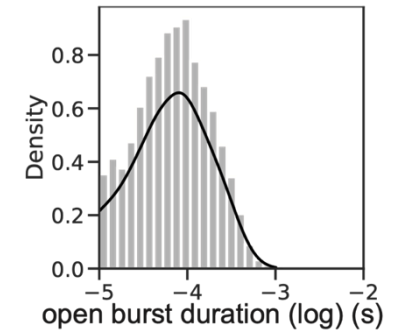
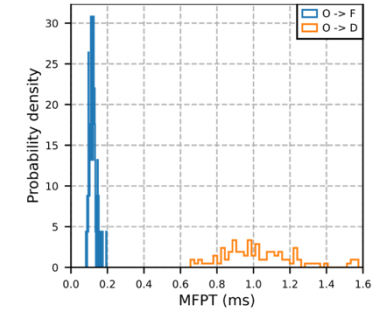
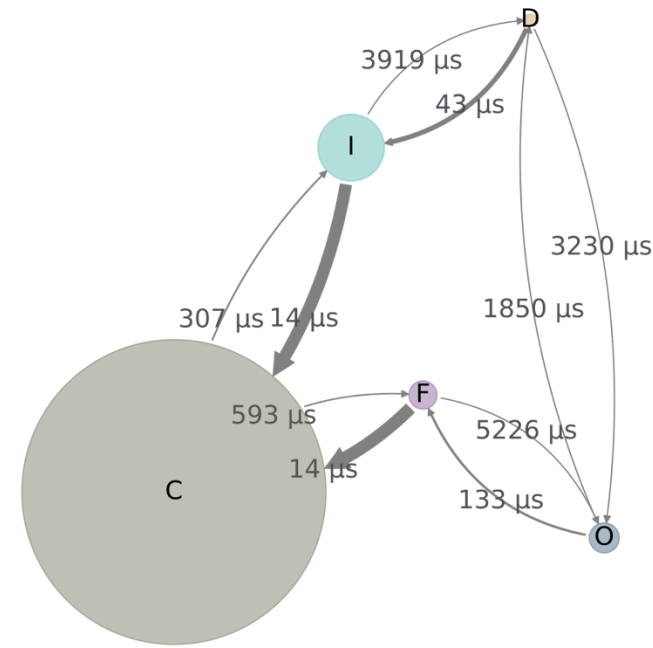
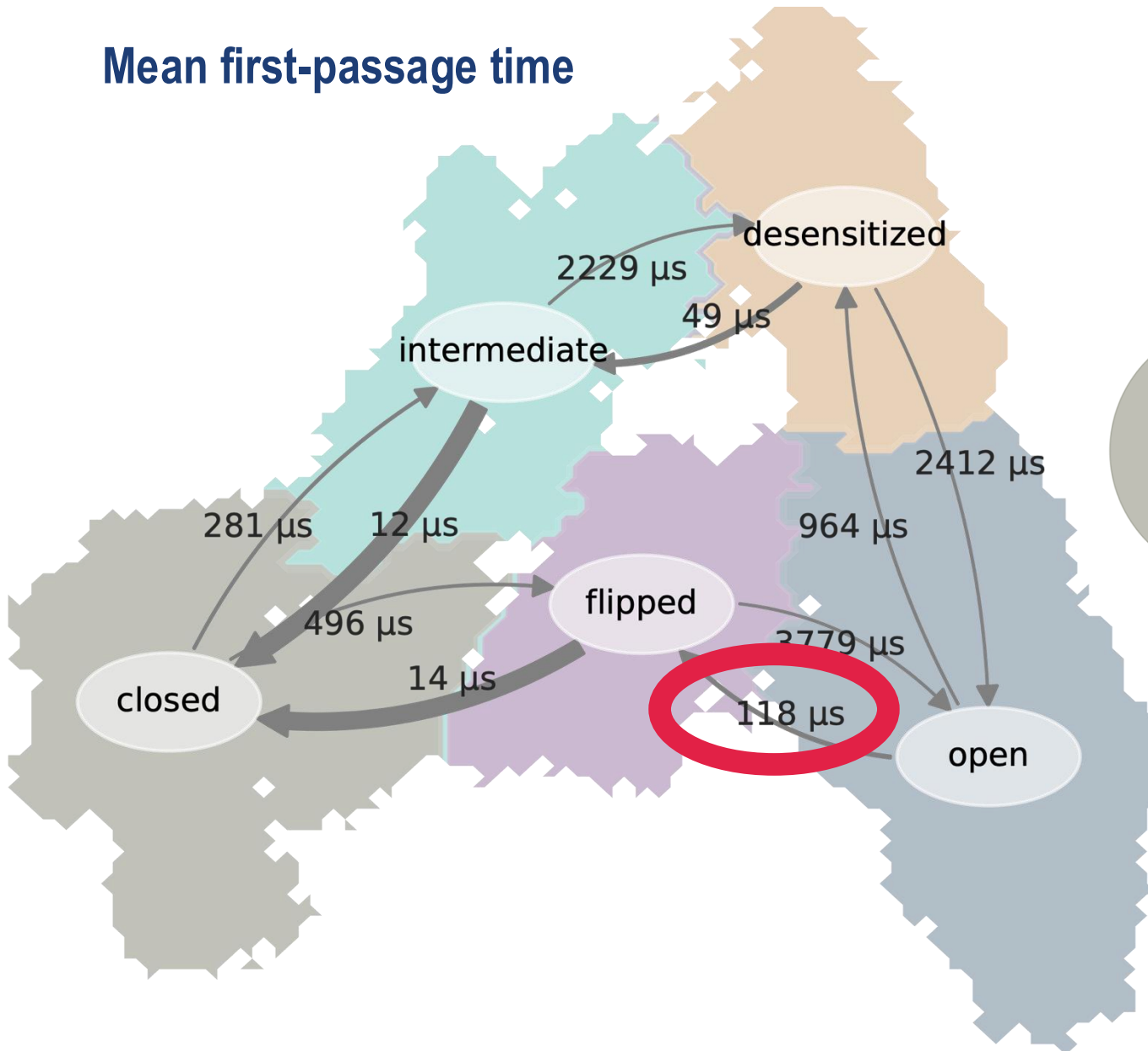
# VISUALIZING 1MS OF NACHR GATING FROM MSM SAMPLING

Time 0  $\mu$ s



# TIMESCALES OF nAChR GATING & DESENSITIZATION FROM MSMs

Mean first-passage time

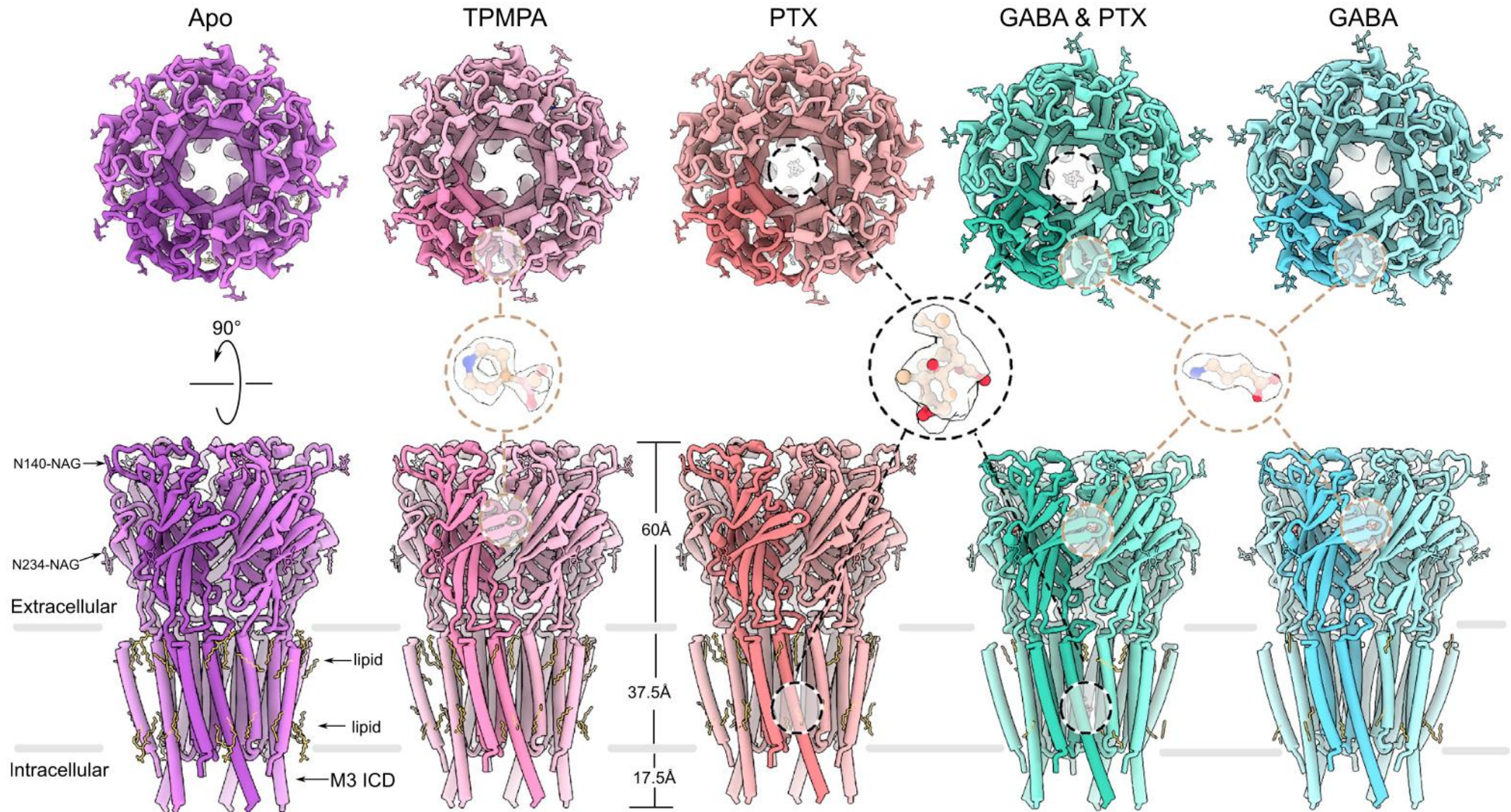


E-phys: “Opening exists as 100  $\mu\text{s}$  isolated events”

Bouzat, 2008, 2018

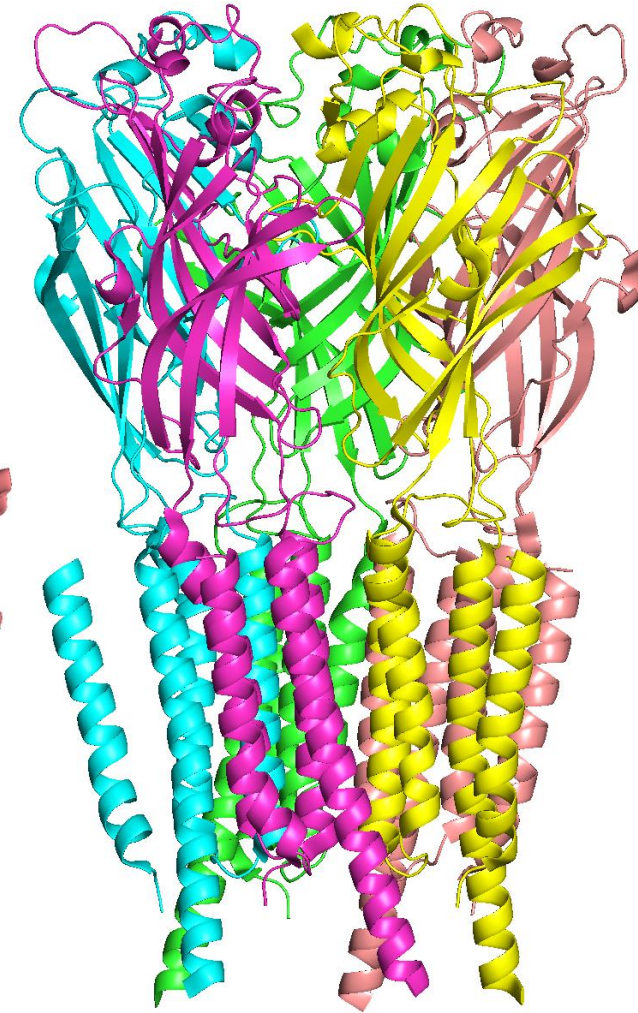
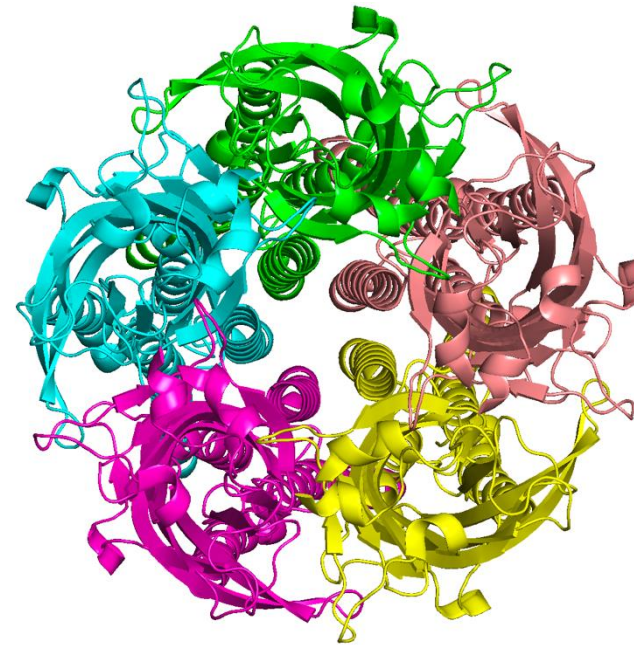
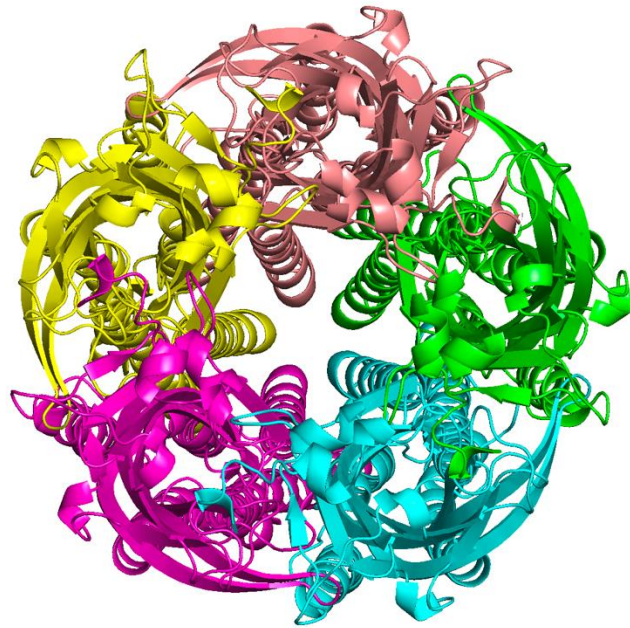
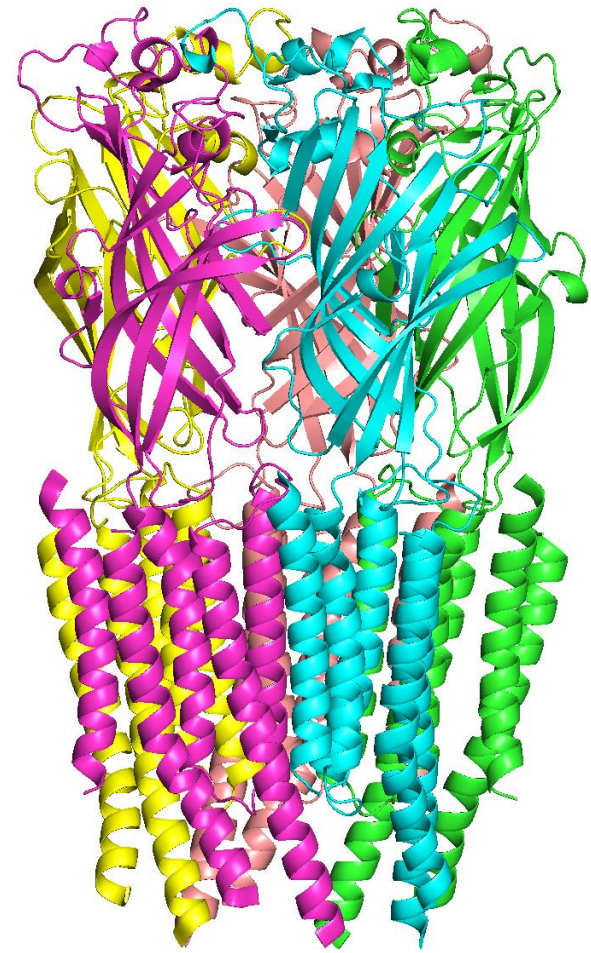
# REVISITING THE SAME CHALLENGE WITH AI

# THE GABA<sub>A</sub>- $\rho$ 1 RECEPTOR

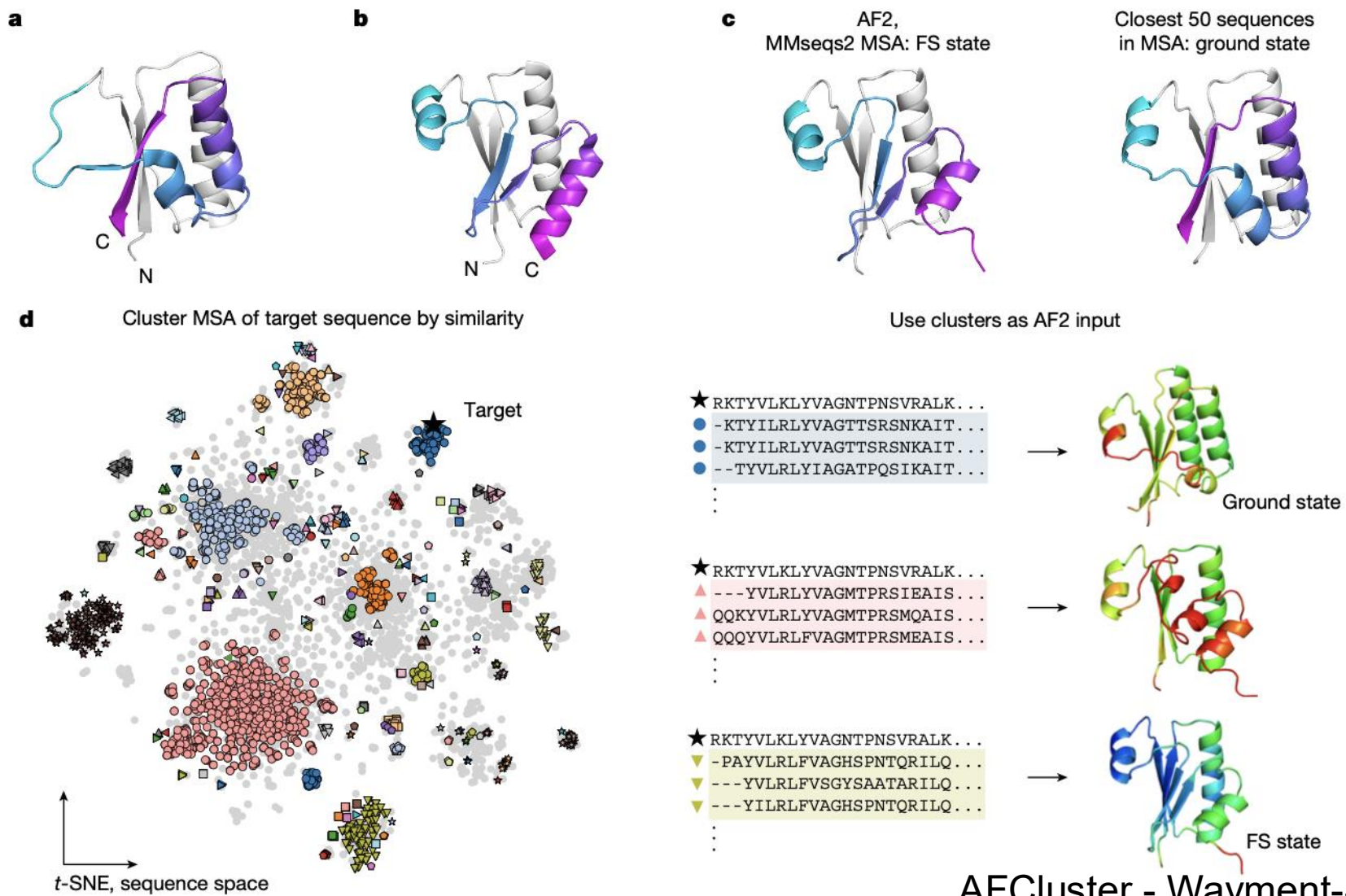


# ALPHAFOLD 3

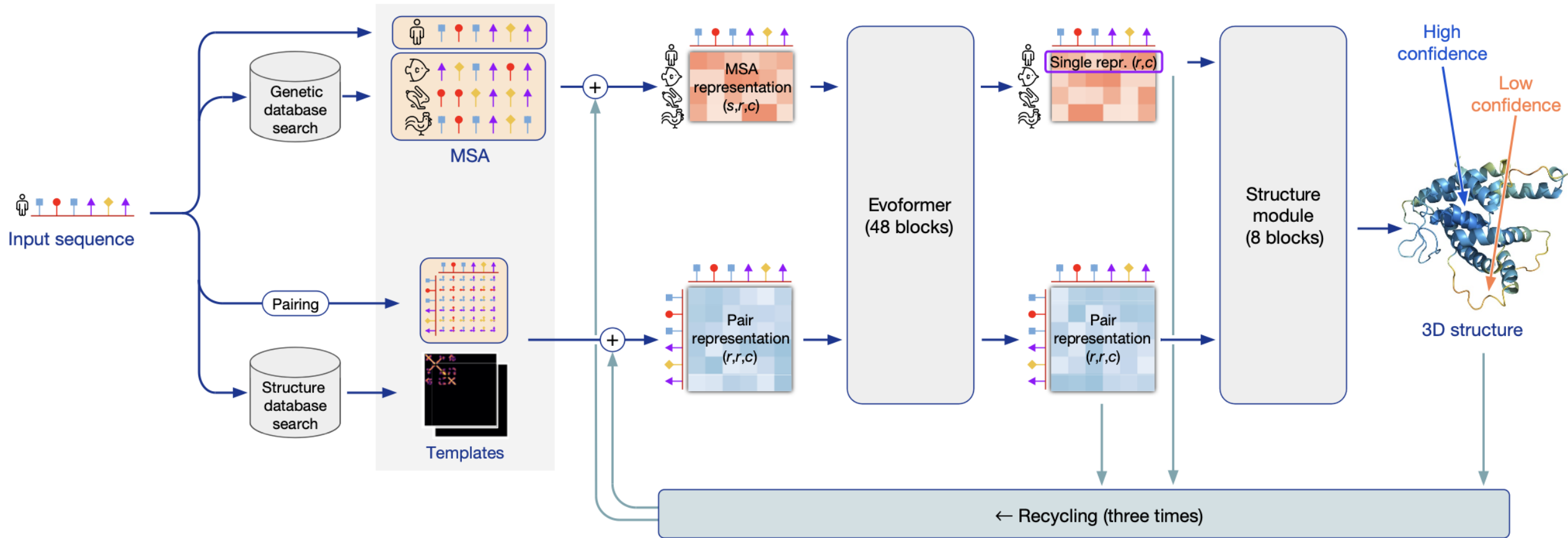
# CRYO-EM



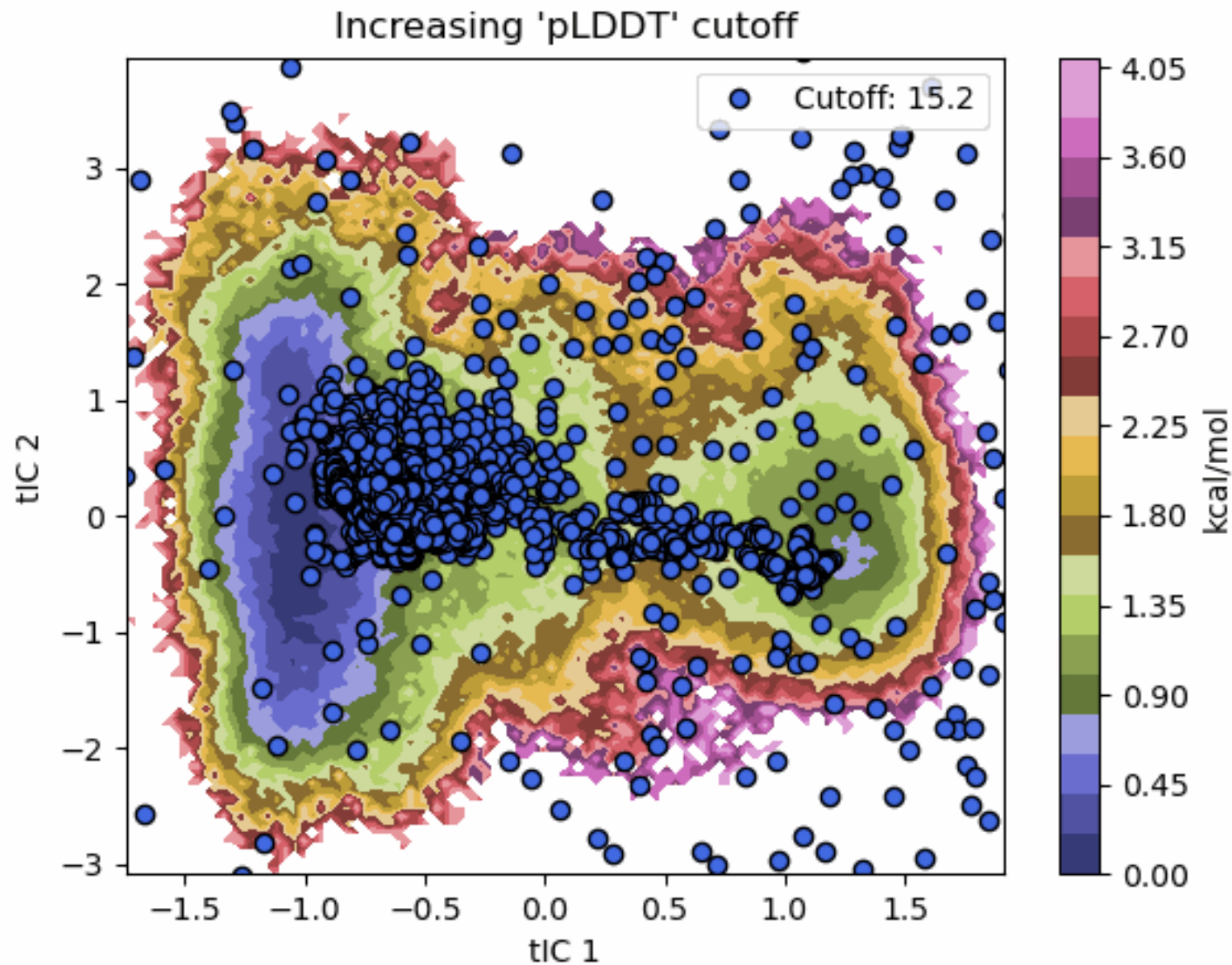
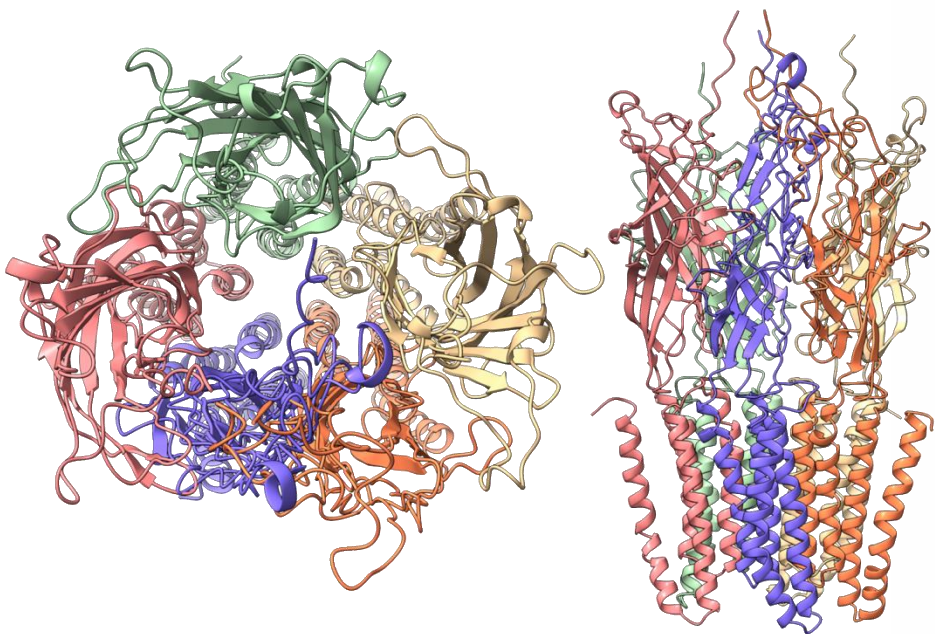
# Can structure prediction methods resolve multiple different states?



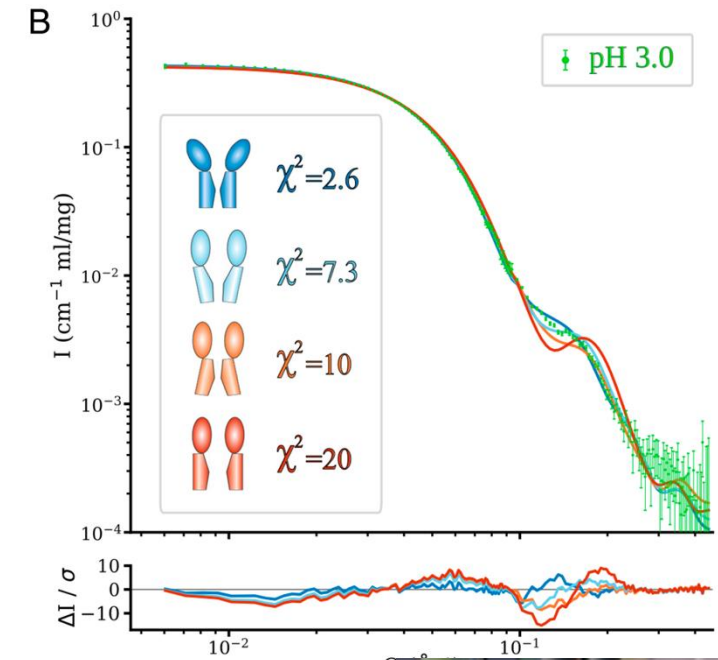
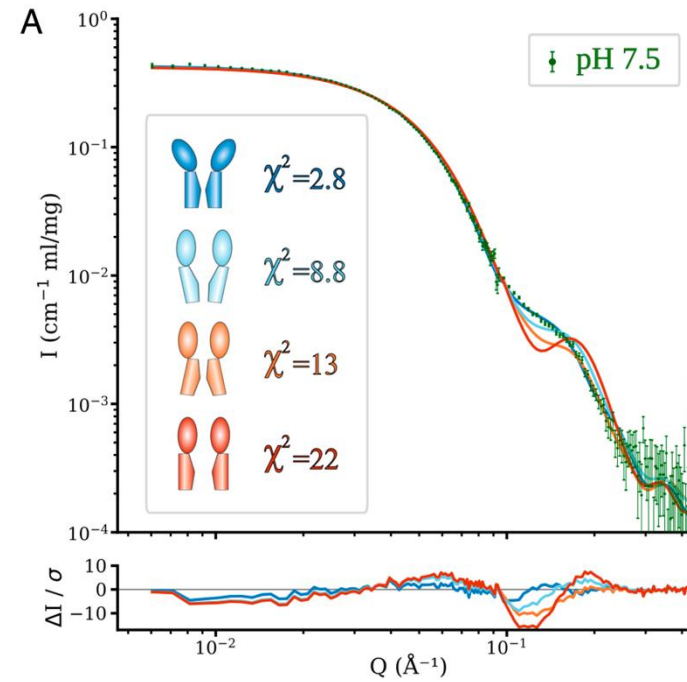
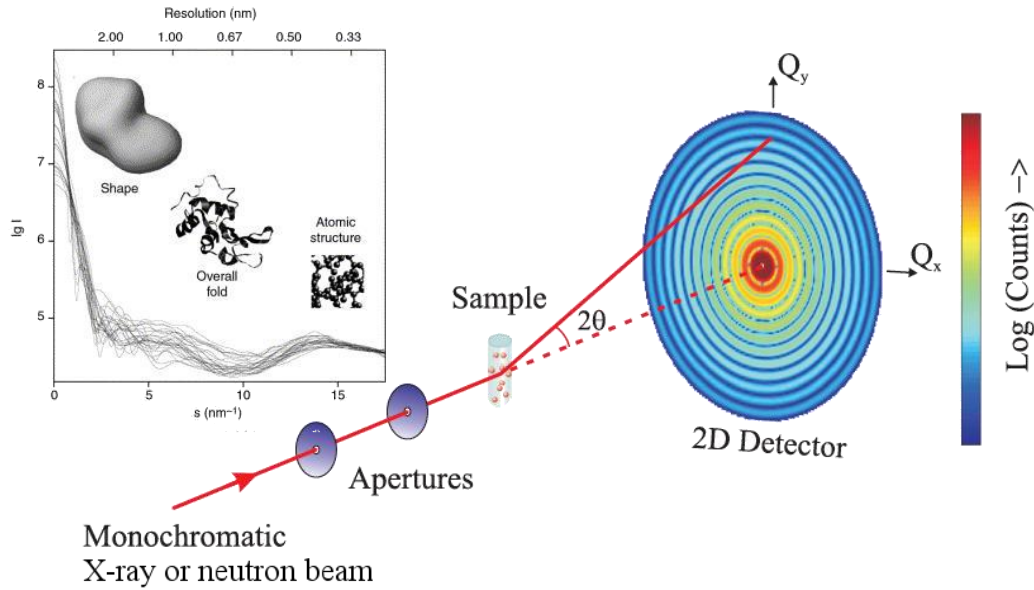
# The MSA stage is critical to what AlphaFold generates



# STOCHASTIC MSA SUBSAMPLING & FILTERING



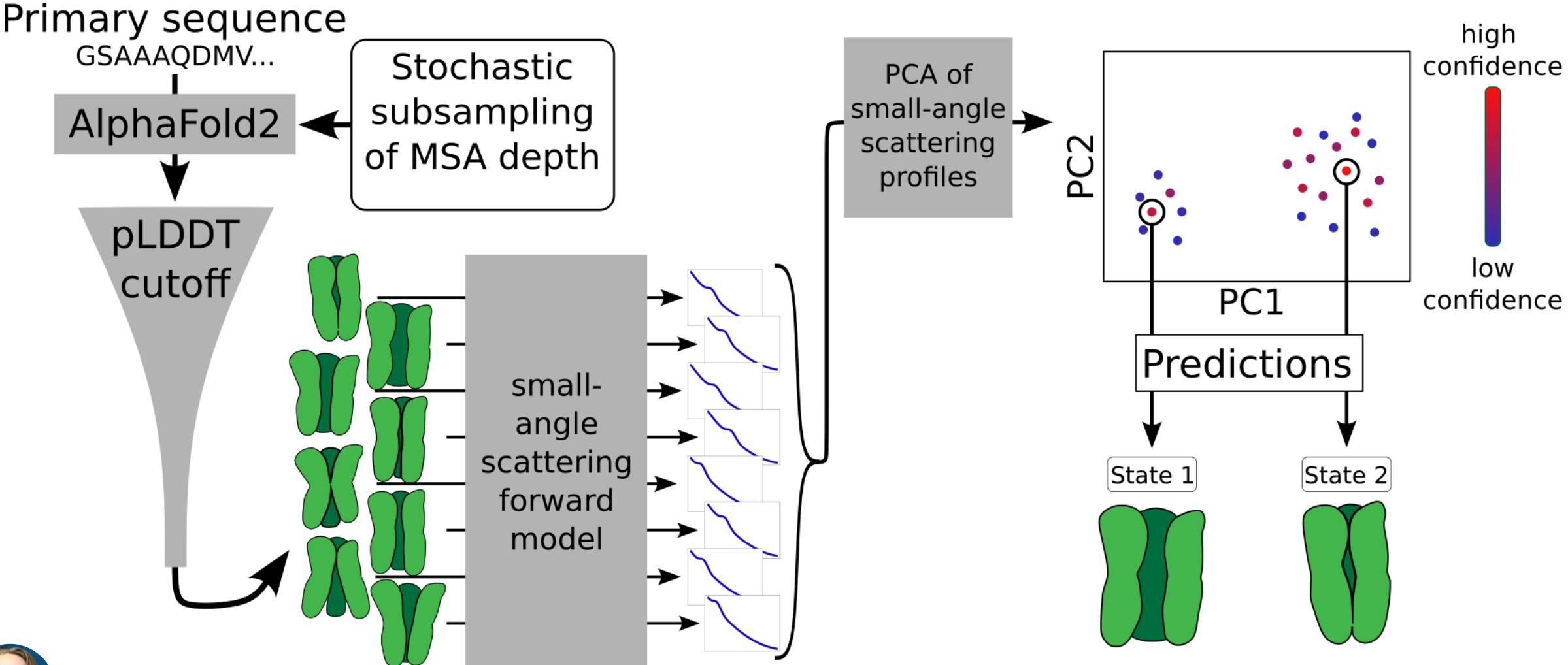
# Use low-resolution data (SANS) to discriminate resting/active state



Marie Lycksell

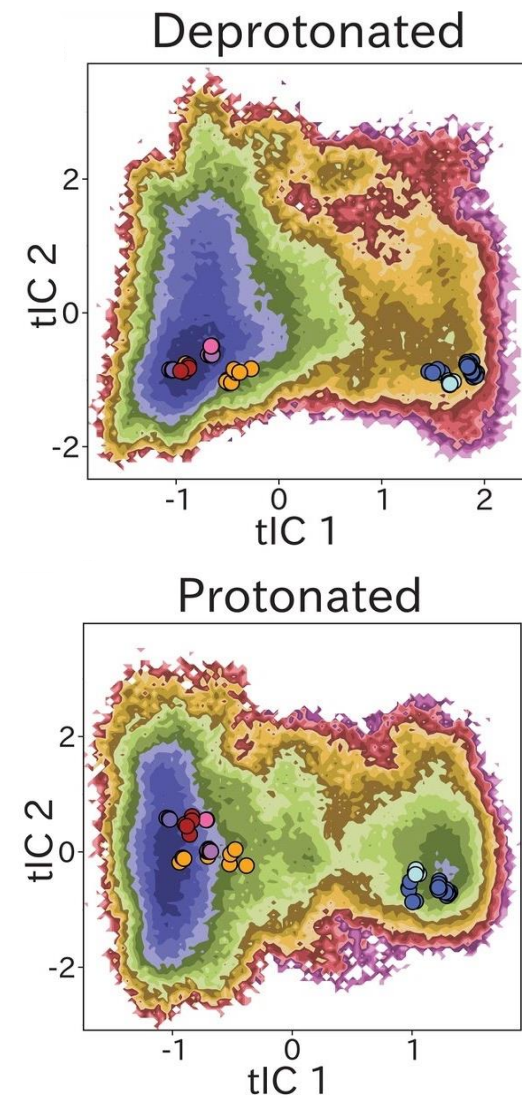
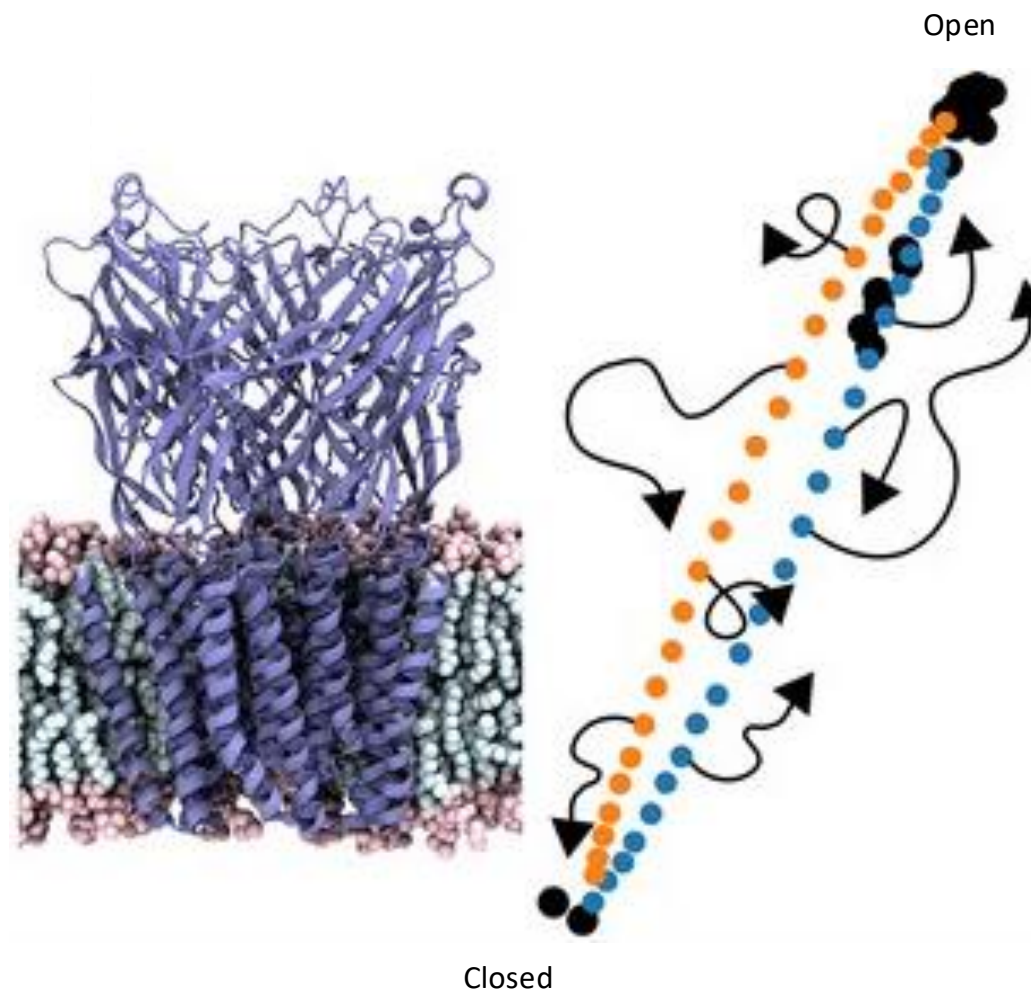
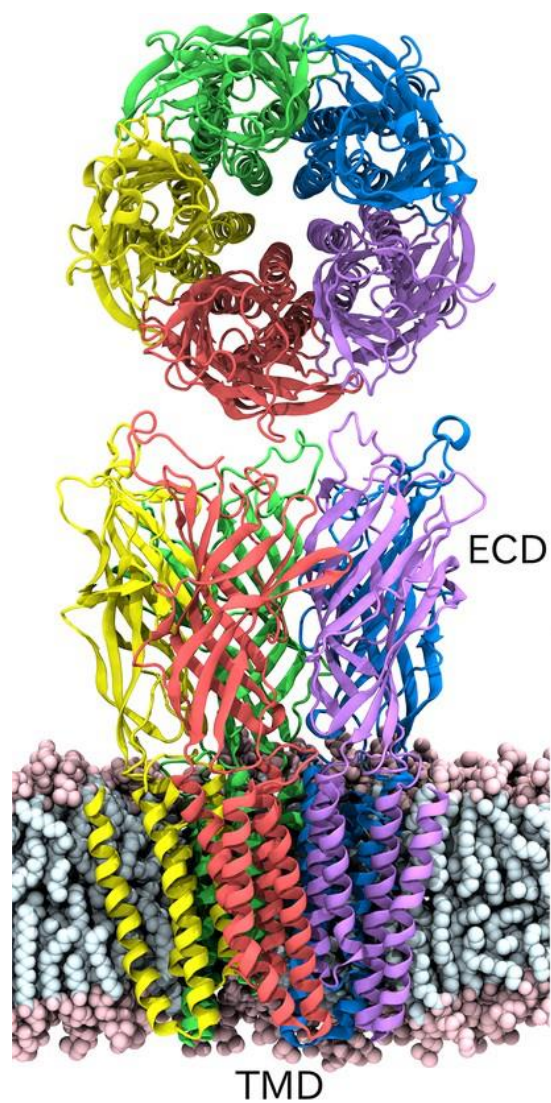


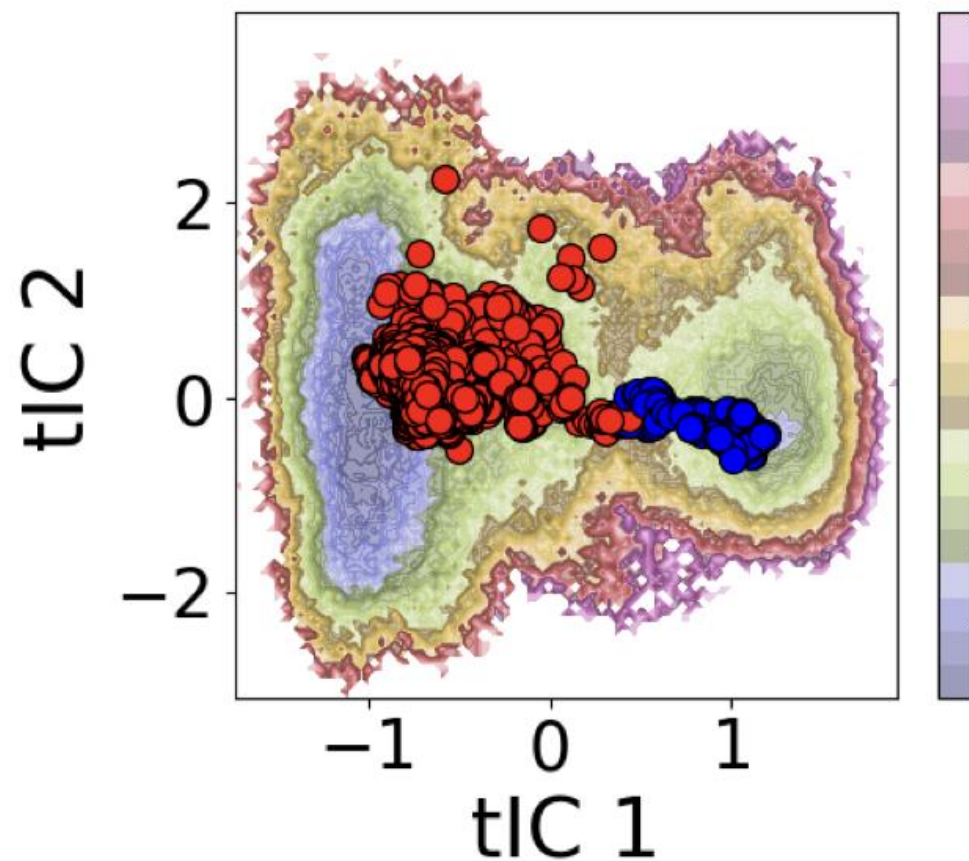
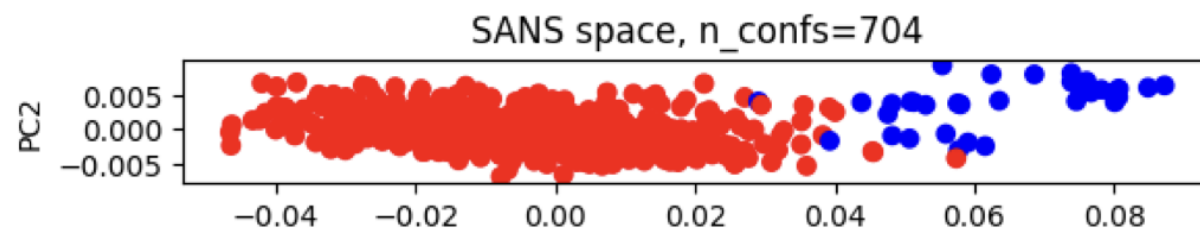
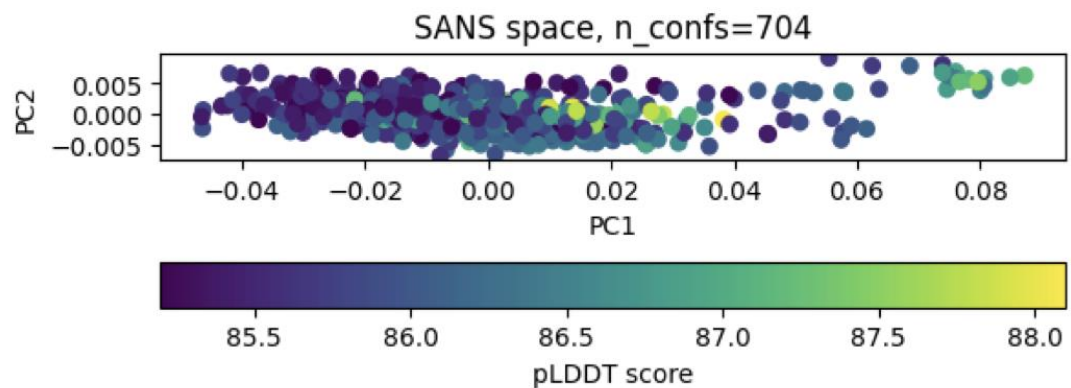
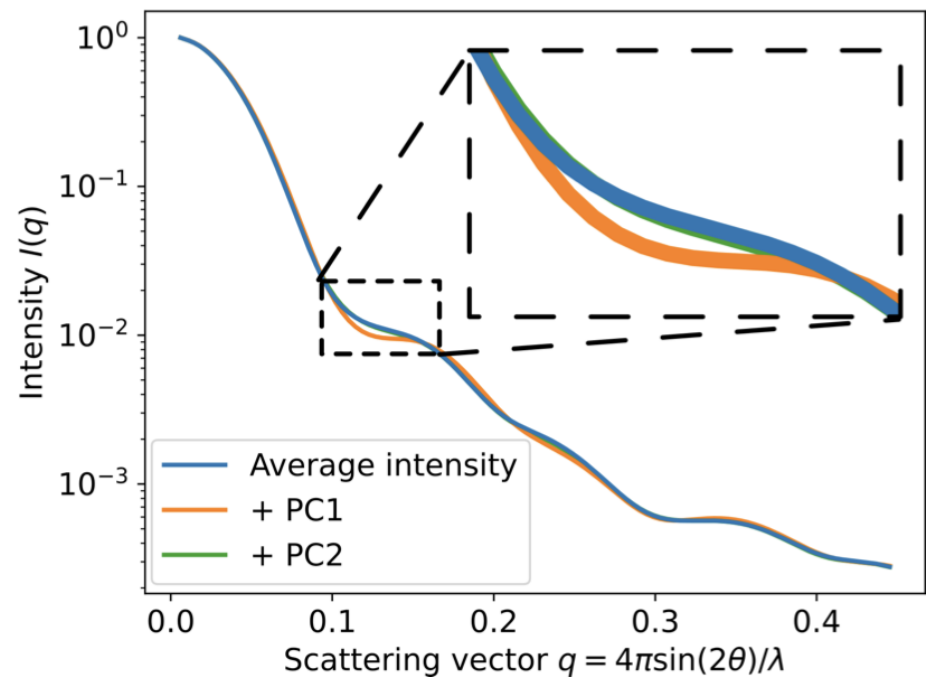
# AlphaFold generates candidates while SANS discriminates



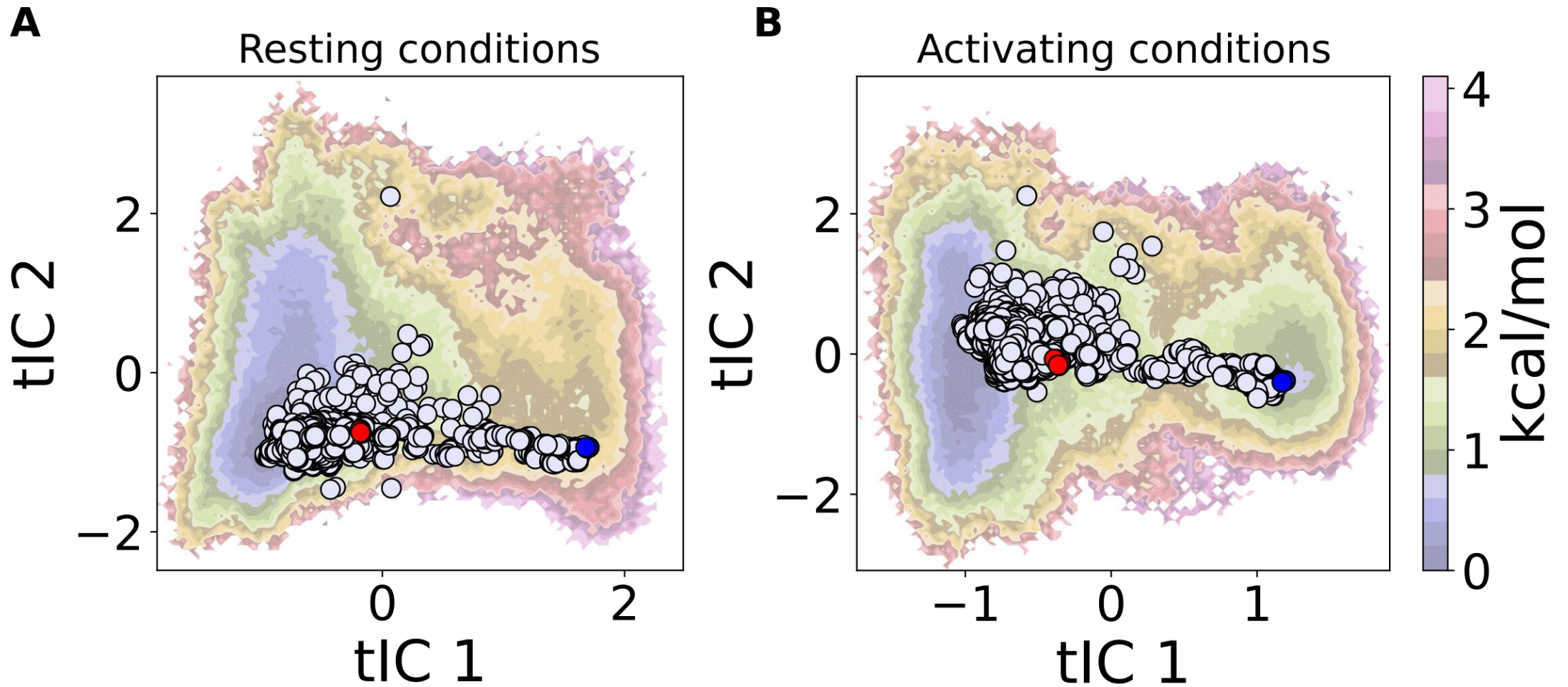
Samuel Eriksson Lidbrink

# Assess against full free energy landscapes from MSM simulations





# Structure ensembles reproduce minima and paths in energy landscapes

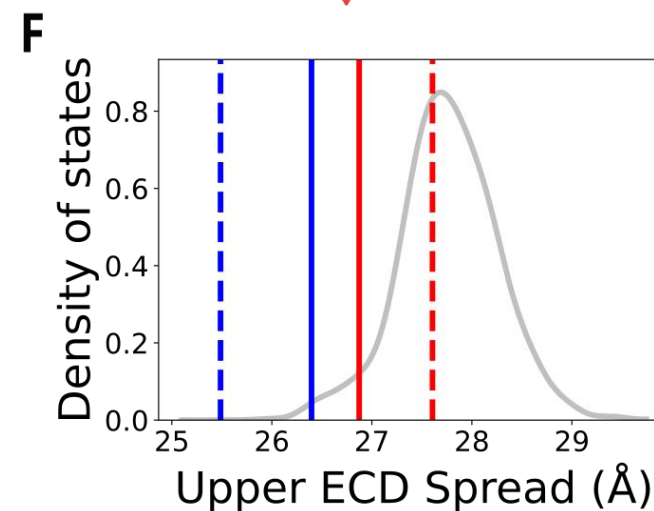
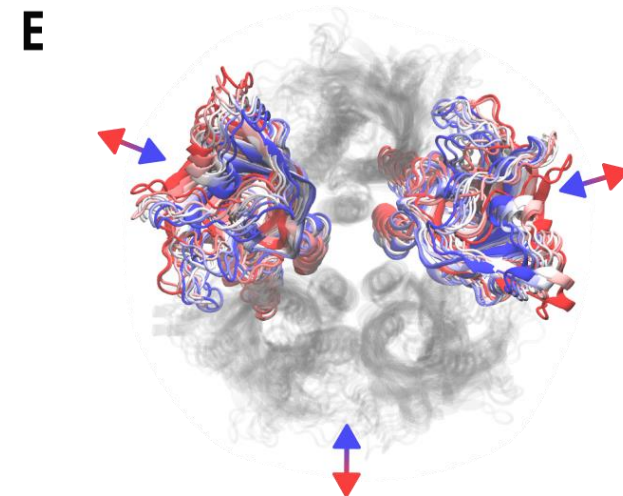
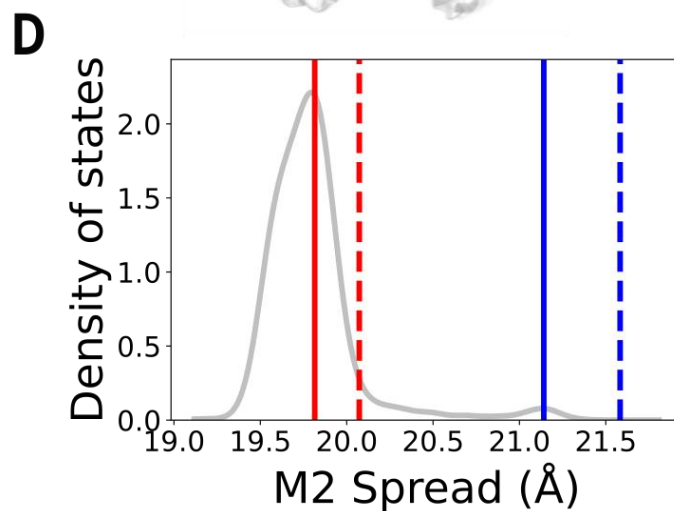
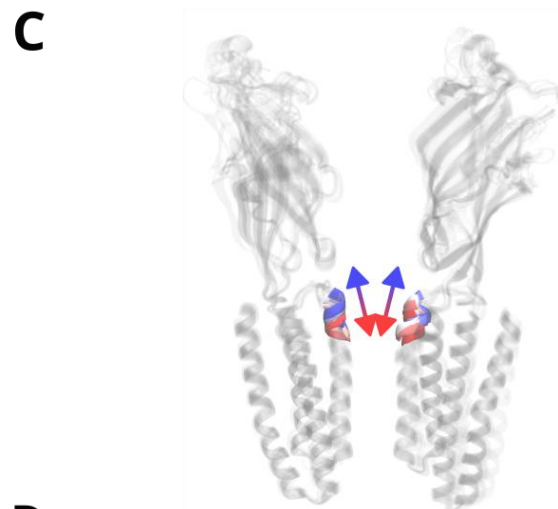
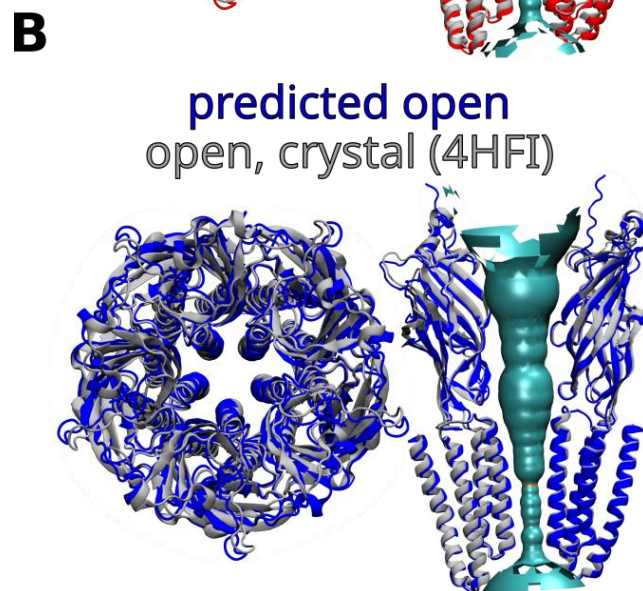
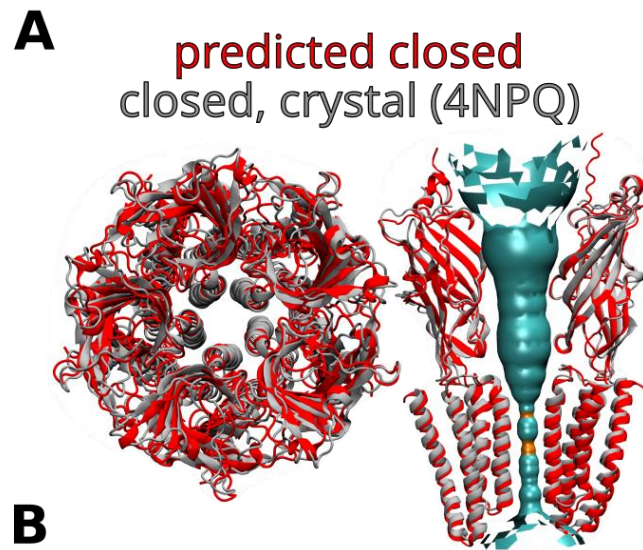


○ AF2, pLDDT  $\geq 75$

● Predicted closed

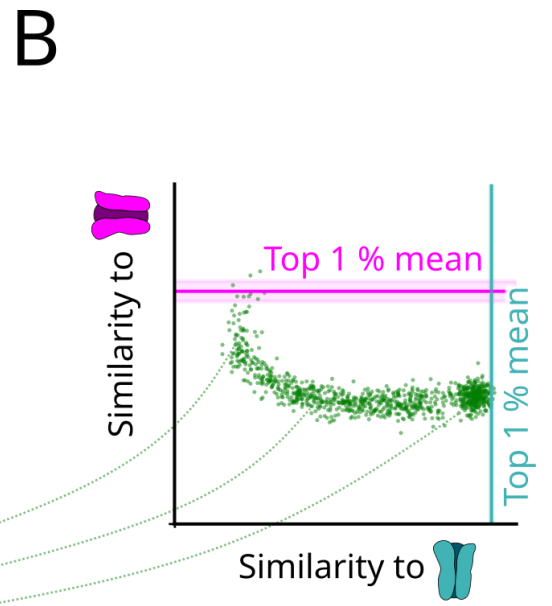
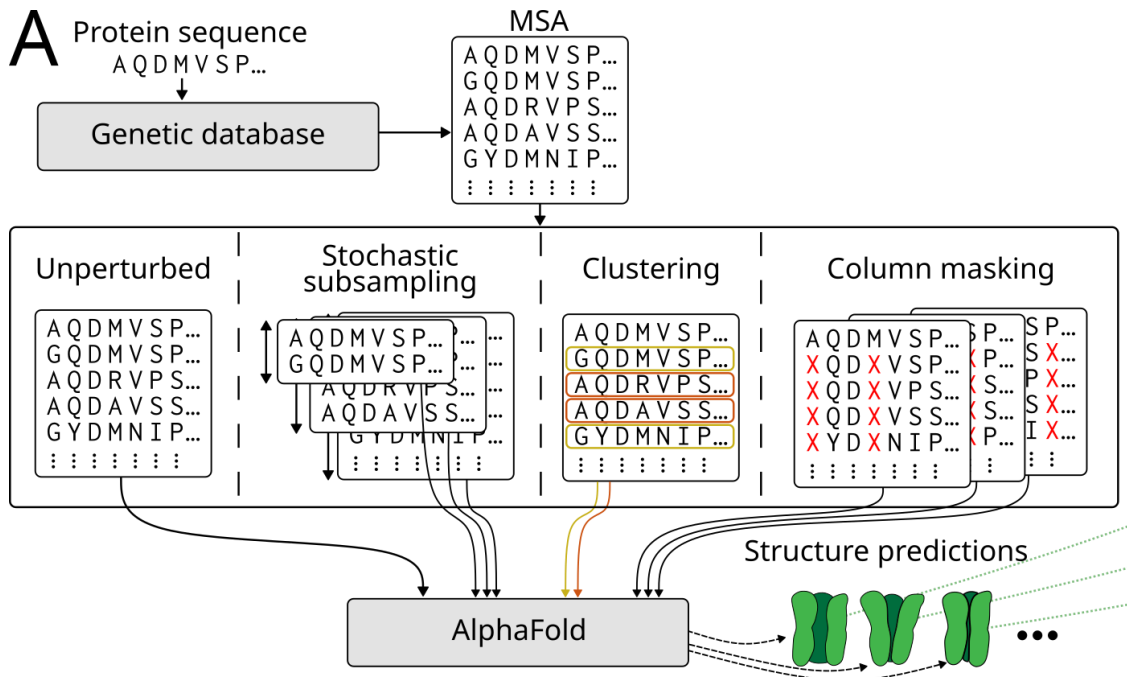
● Predicted open

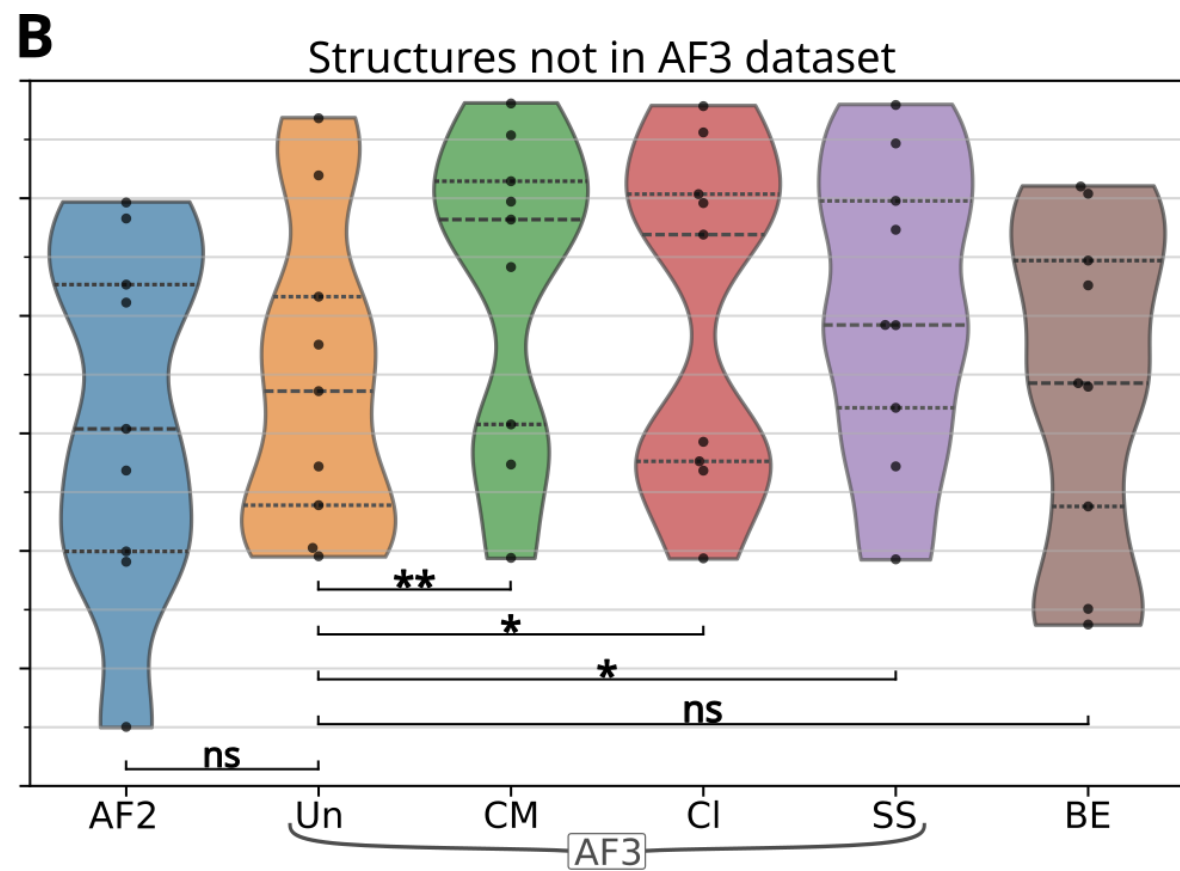
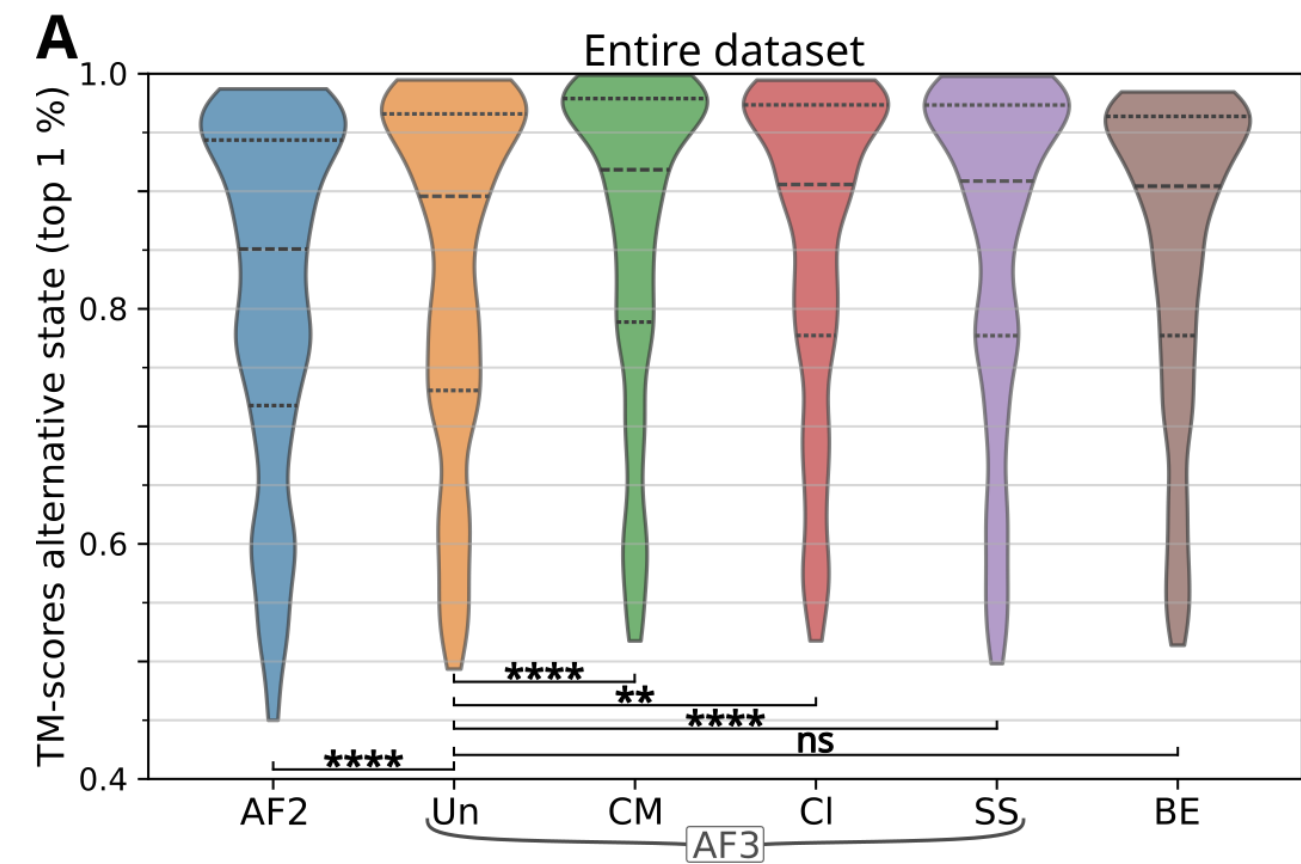
# Selected predictions are close to X-ray structures of closed/open states



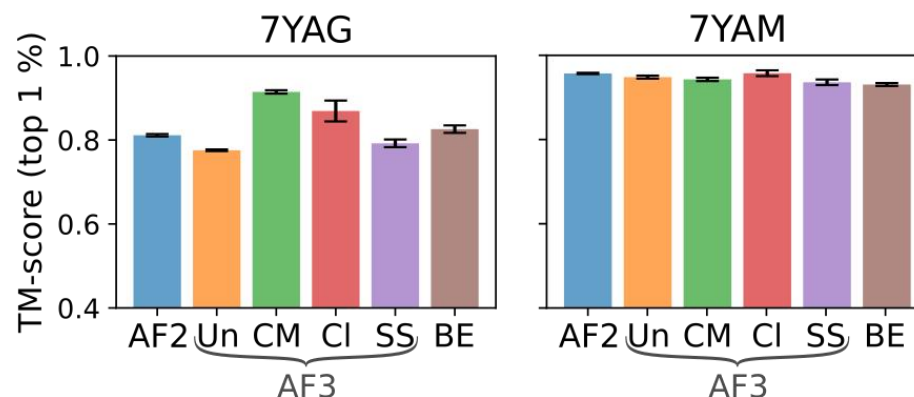
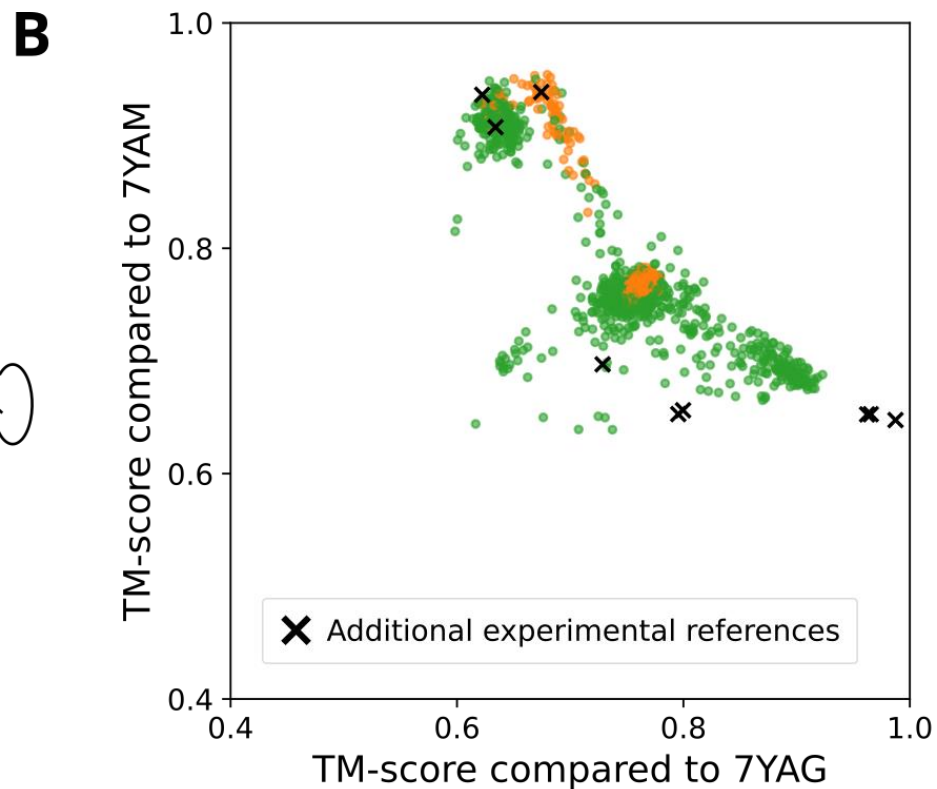
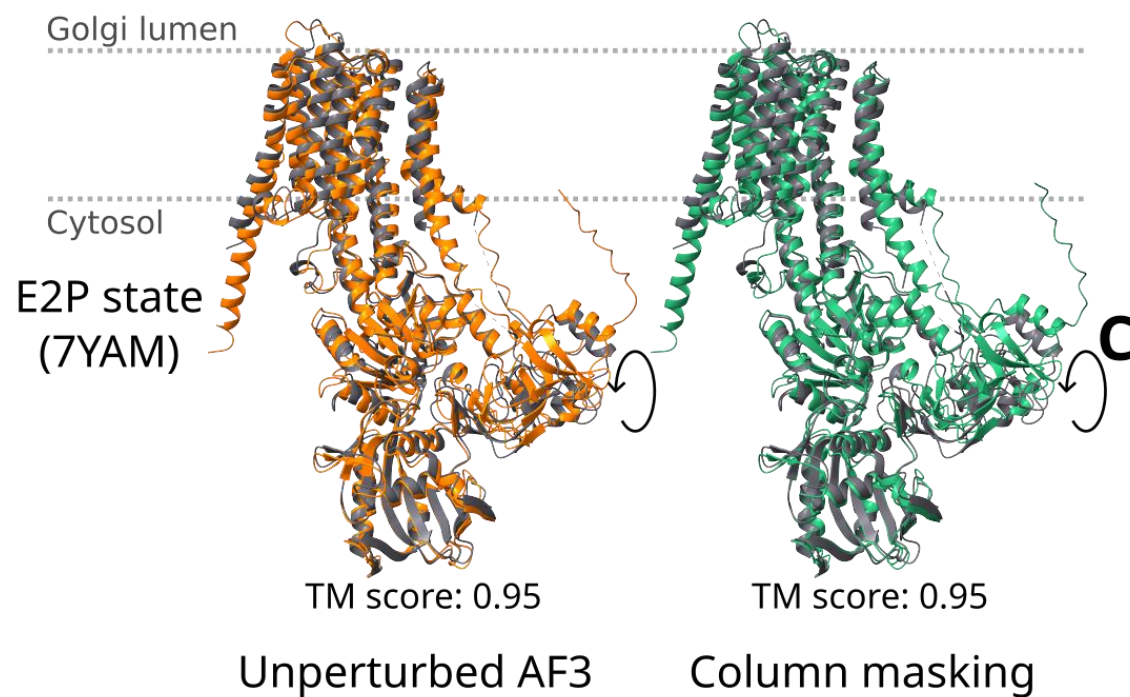
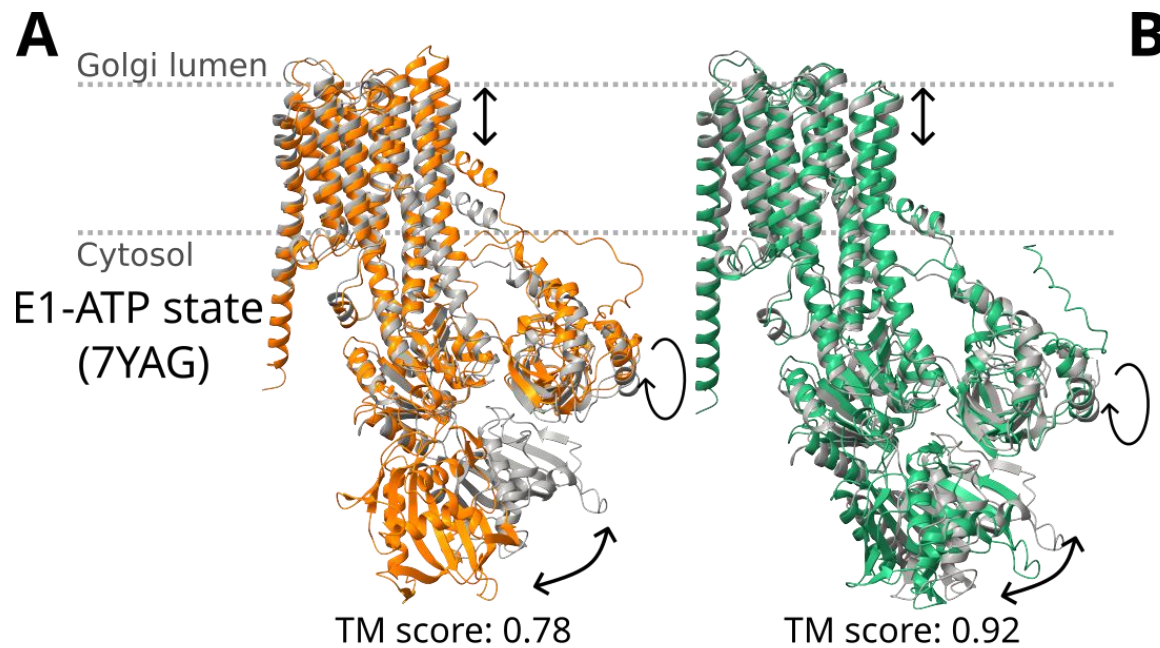
--- Closed, crystal      --- Open, crystal      — AF2, pLDDT $\geq$ 75  
— Predicted closed      — Predicted open

**CAN AI METHODS PREDICT  
CONFORMATIONAL LANDSCAPES IN  
GENERAL?**





● AF2 ● Unperturbed AF3 (Un) ● Column masking (CM) ● Clustering (CI) ● Stochastic subsampling (SS) ● BioEmu (BE)

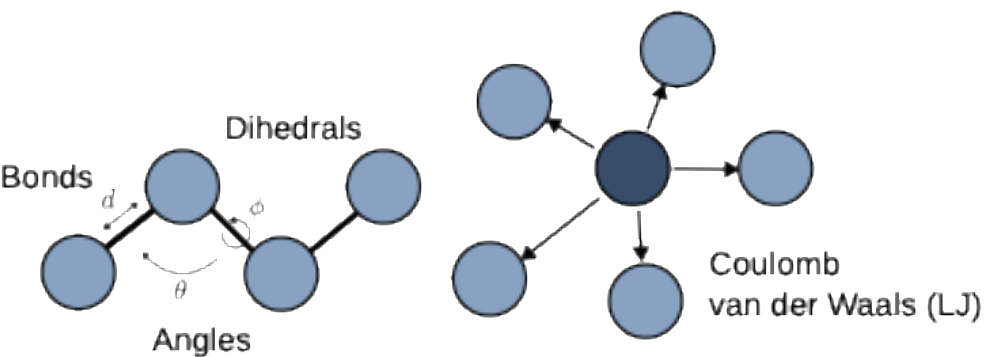


# USING AI FOR YOUR SIMULATIONS: NEURAL NETWORK POTENTIALS

# The Efficiency-Accuracy Tradeoff

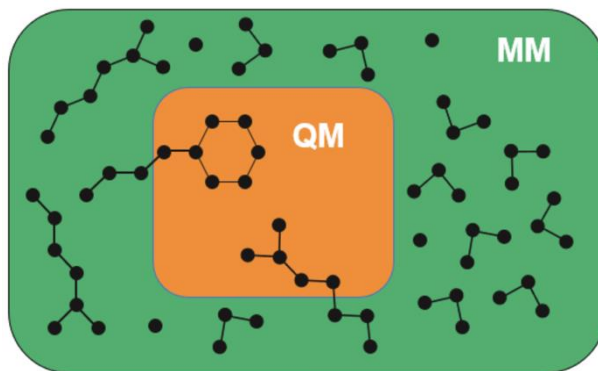
## Molecular Mechanics (MM)

- + computationally efficient
- rigid topology
- non-standard systems need careful validation



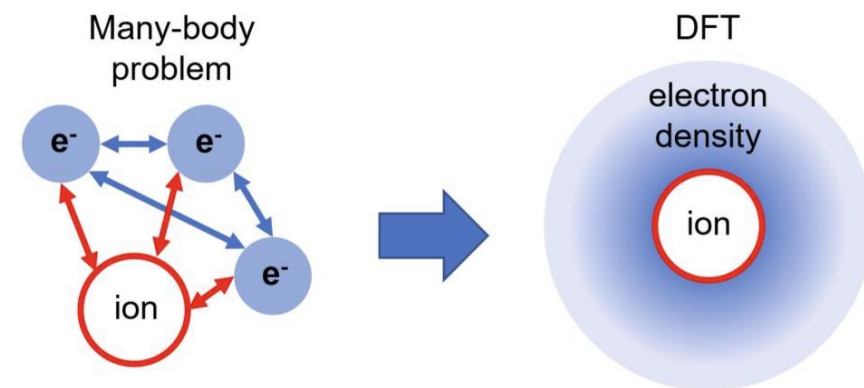
## Hybrid QM/MM

- + best of both worlds
- ... still too expensive!



## Quantum Mechanics (QM)

- + first principles, QM-accurate
- + flexible
- + can handle bond breaking
- Computationally very expensive



$$U(\mathbf{r}^N) = U_{\text{bonded}} + U_{\text{Lennard-Jones}} + U_{\text{Coulomb}}$$

$$H = H_{\text{QM}} + H_{\text{MM}} + H_{\text{QM-MM}}$$

$$i\hbar \frac{d}{dt} |\Psi(t)\rangle = \hat{H} |\Psi(t)\rangle$$

Figure credit:  
GROMACS Website / Dmitry Morozov  
Sergei Posysaev

Can we do better?

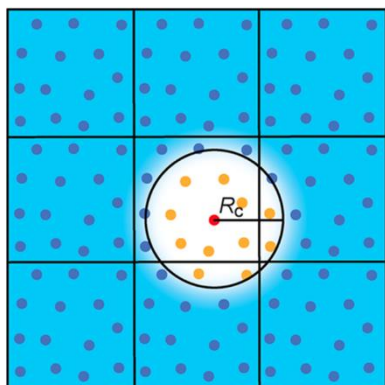
# Neural network potentials in GROMACS 2026

Quantum accuracy at (almost) force-field performance

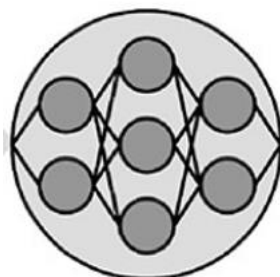
Local chemical environments

Cartesian coordinates  $\mathbf{r}_i$

Nuclear charges  $Z_i$



NNP  
Parameters  $\theta$



Potential Energy, Forces  
(via backprop.)

$$E(\mathbf{r}^N), F(\mathbf{r}^N)$$

MD  
Simulation



Minimize loss function

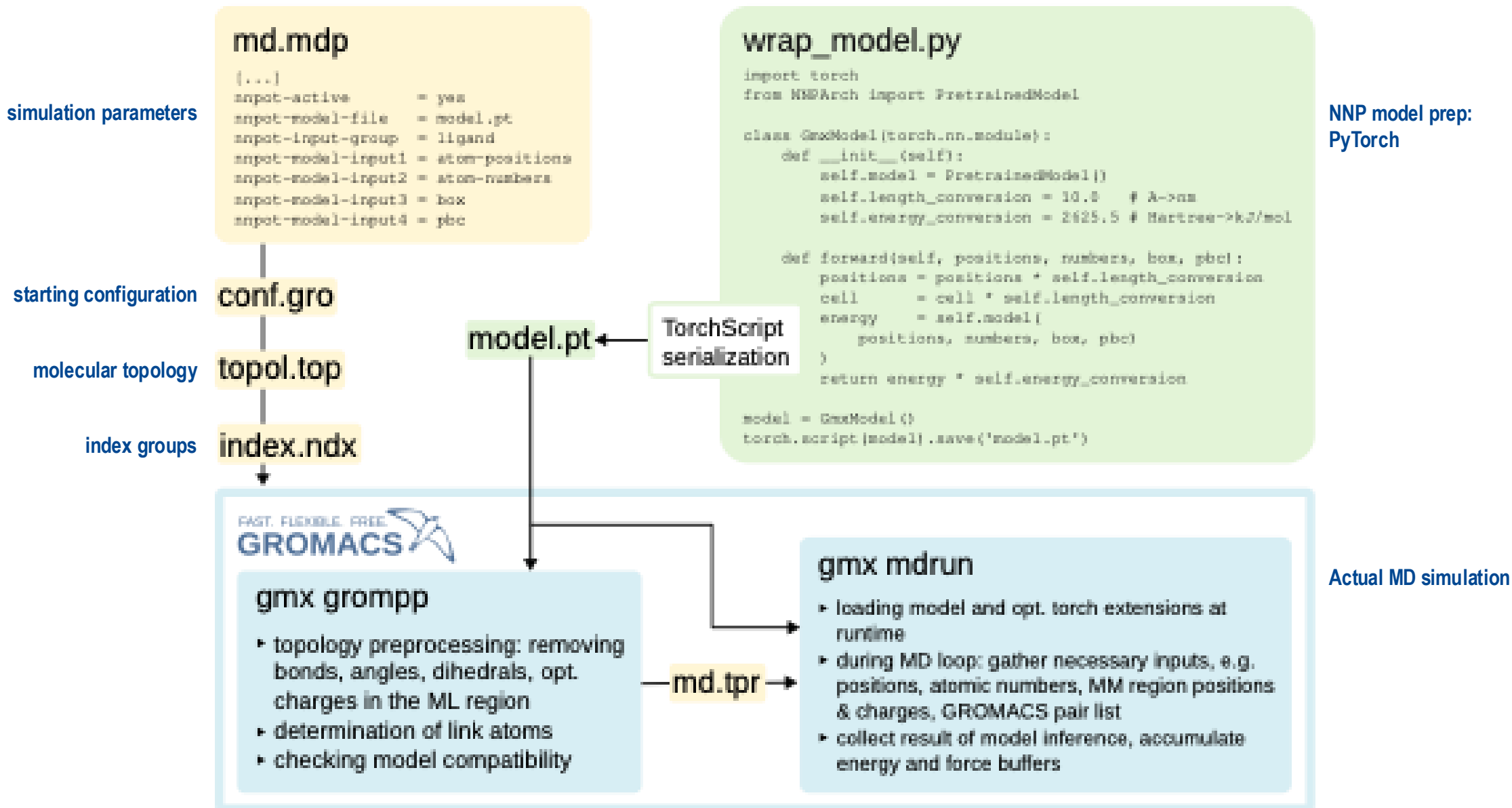
$$\mathcal{L} = \frac{1}{N} \sum_i (E_i - \hat{E}_i)^2$$

Electronic structure data  
DFT/CCSD(T)

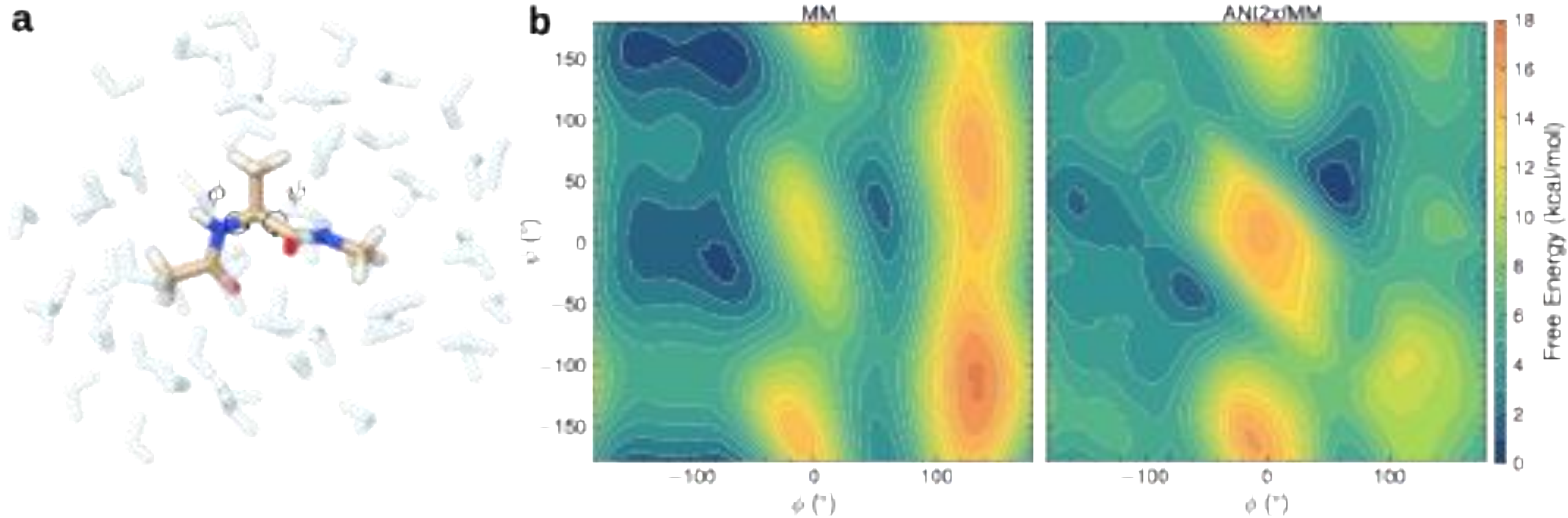
Popular NNP architectures:

- ANI-1x/2x [1,2]
- SchNet [3]
- MACE [4]
- EMLE [5]

# The nnpot interface in GROMACS 2026



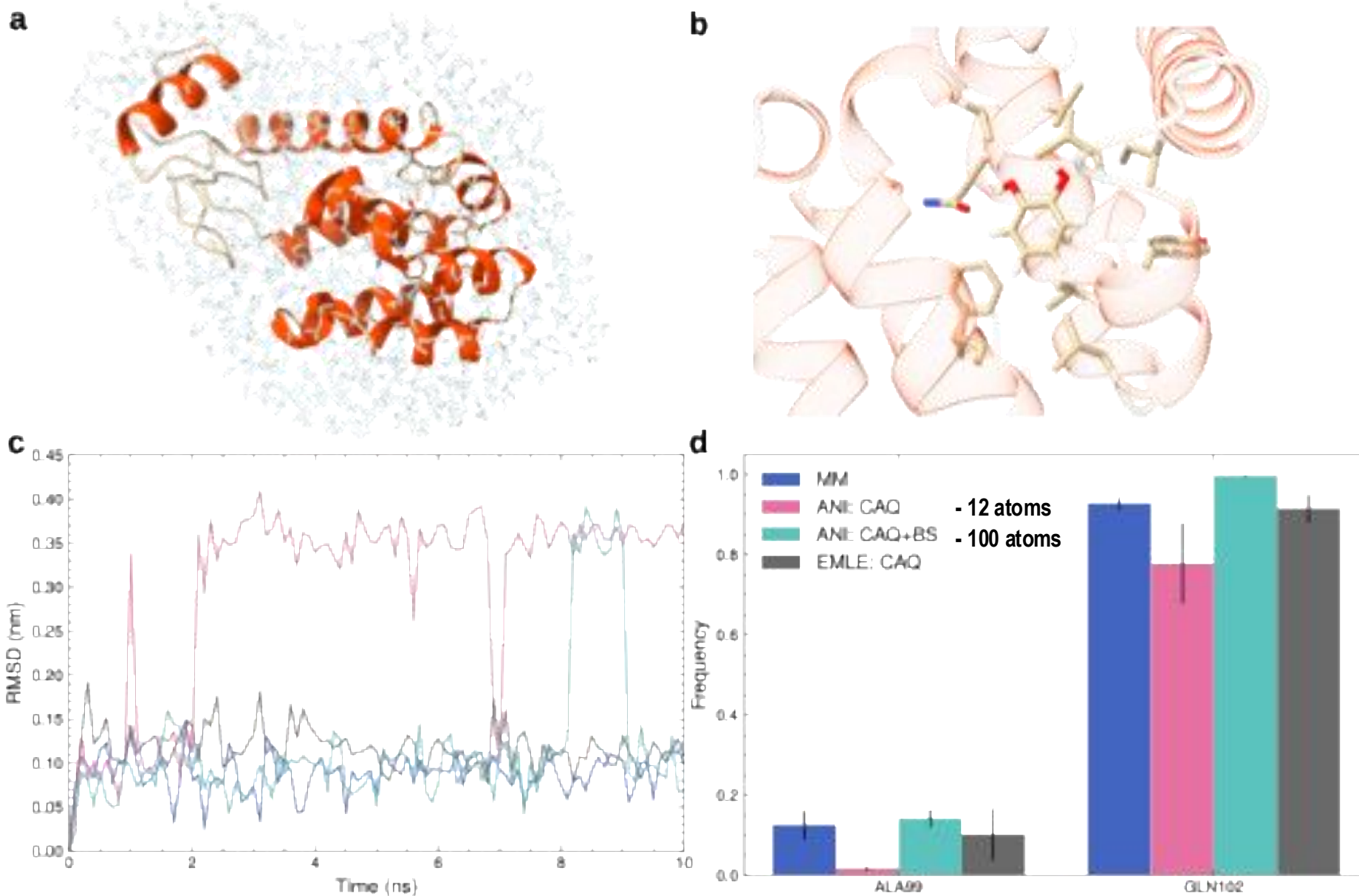
# ML/MM Enhanced Sampling with GROMACS AWH



[Unpublished data]

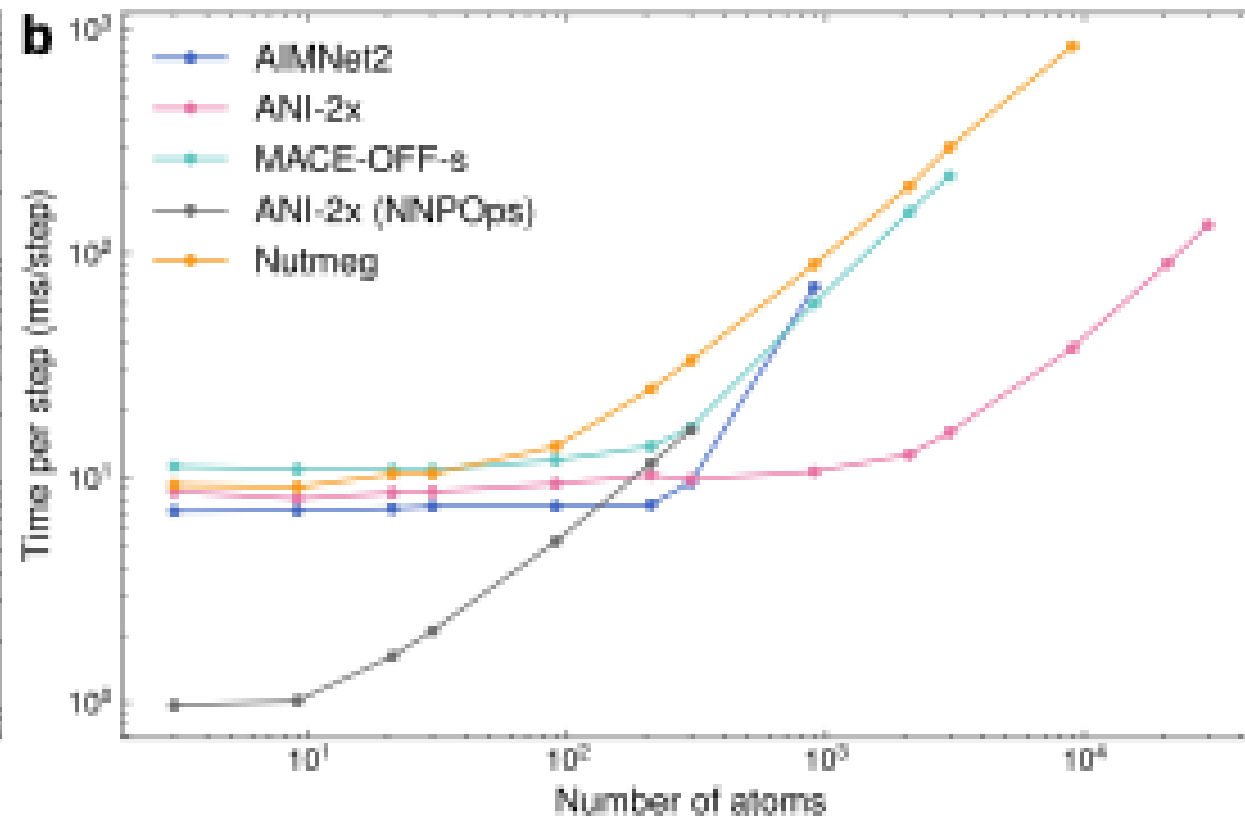
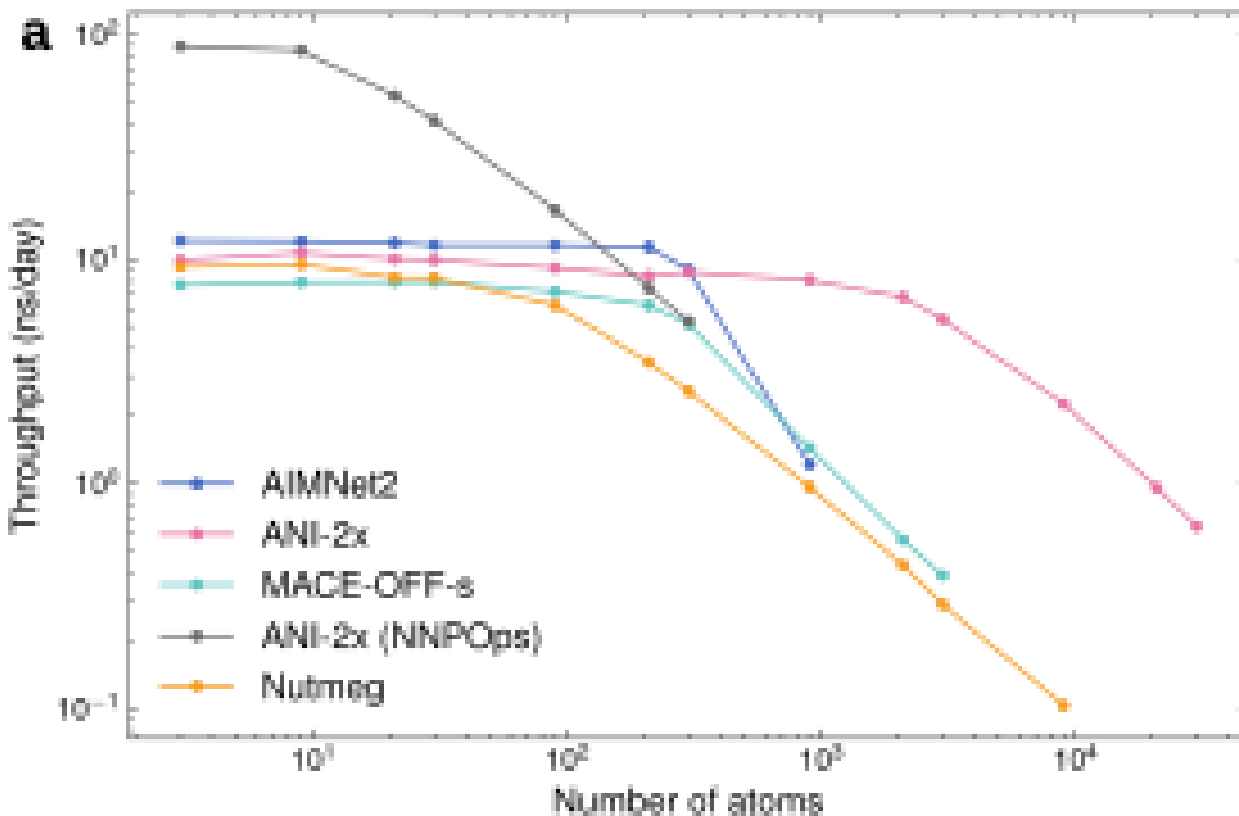
# Protein-Ligand Systems: ML region size matters

Catechol bound to  
T4 Lysozyme  
L99A/M102Q



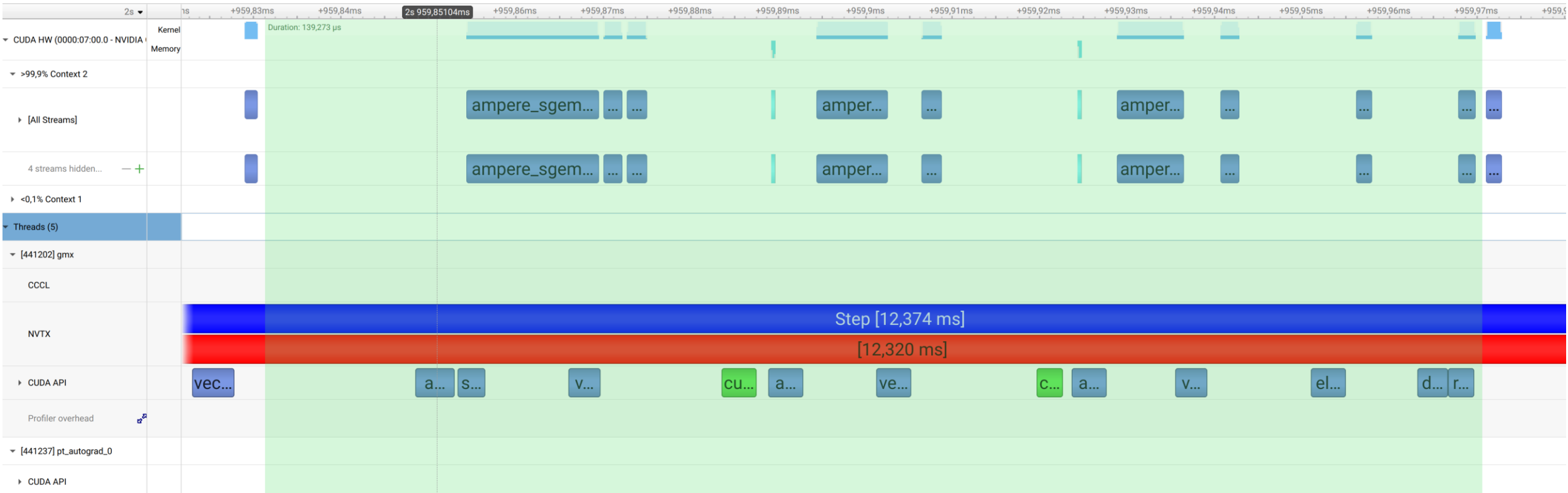
# Performance Scaling

NVIDIA RTX3070  
1 fs time step  
32-bit floating point



Smaller systems are dominated by libtorch latency

# Libtorch is often very inefficient – room to improve



# AI is not a threat to HPC - but our best friend!

- You don't need to train models yourself - just as you don't need to write all code
- Generative AI is amazing as a hypothesis/model generator
- It is less amazing as final arbiter of physical truth - combine with hard data
- AlphaFold/Chai-1 and friends are great at generating many models, but the ensembles are *not* anywhere near a correct Boltzmann distribution
- Beware of the hype - models are usually not as good in the real world
- But ... the speed of improvement is insane
- Being creative in using data often beats developing a really advanced AI model
- The vast majority of HPC users employ codes from others – use AI from others too!

# Acknowledgments

**GROMACS:** Szilárd Páll, Berk Hess, Andrey Alekseenko, Artem Zhmurov, Umair Sadiq, Lukas Müllender, Vedran Miletic, Cathrine Bergh, Yang Zhang, Petter Johansson and numerous other contributors

**Vendor co-design:** Alan Gray, Mahesh Doijade (Nvidia); Mark Abraham, Roland Shulz (Intel); Paul Bauer (AMD)

**Molecular simulation:** Magnus Lundborg, Tatiana Shugaeva, Sebastian Wingbermhle, Alessandra Villa, Farzaneh Jalalypour, Nandan Haloi, Anton Jansen, Samuel Eriksson Lidbrink, Rebecca Howard



Szilárd Páll  
KTH



Andrey Alekseenko  
KTH



Berk Hess  
KTH

Knut and Alice  
Wallenberg  
Foundation



European  
Research  
Council



**serc**  
Swedish e-Science Research Centre



**bioexcel**  
Center of Excellence for Computational Biomolecular Research



EuroHPC  
Joint Undertaking

## Take-home messages:

- SW is becoming like HW: No silver bullets remaining
- Constant fight: balance latency with GPU speedup
- MNNVL can be fantastic for traditional HPC apps
- GROMACS-2026 can do neural network potential simulations
- But the really big wins come from replacing the entire iterative approach with AI
- AI has a ton of important applications, but suffers from hype and lack of critical testing

