

# TOWARD REALIZING THE DREAM OF FREE-ENERGY BASED SMALL MOLECULE DESIGN



**John D. Chodera**

MSKCC Computational Biology Program

<http://www.choderalab.org>

Slides available at <http://www.choderalab.org>

## DISCLOSURES:

- Scientific Advisory Board, Schrödinger

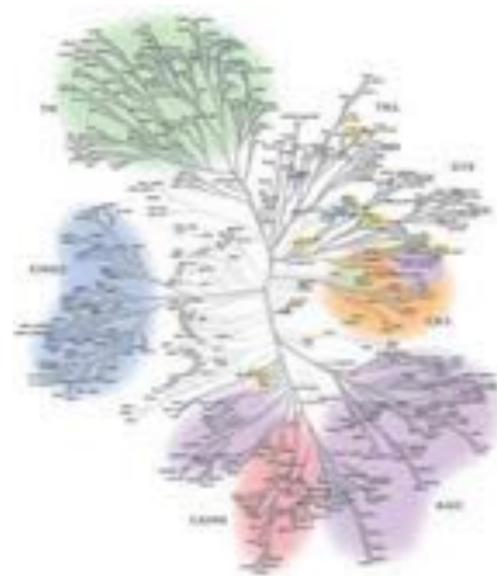
14 July 2017 - Telluride Free Energy Workshop

# SOMETIMES, DRUG DISCOVERY WORKS WELL

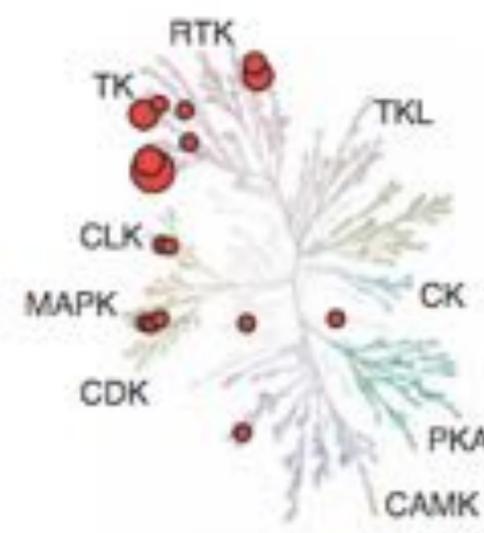
Bcr-Abl fusion constitutively activates ABL in CML patients, resulting in unchecked white blood cell proliferation



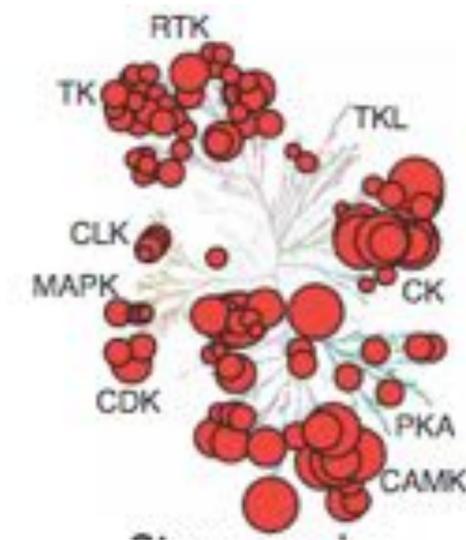
**imatinib** bound to **c-Abl** [PDB:1IEP]



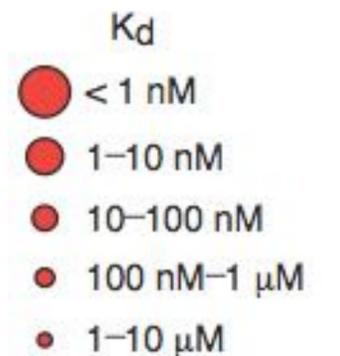
**human kinome**  
[518 kinases]



**imatinib**  
[blockbuster drug]



**staurosporine**  
[toxic natural product]



# DRUG DISCOVERY USUALLY ENDS IN FAILURE

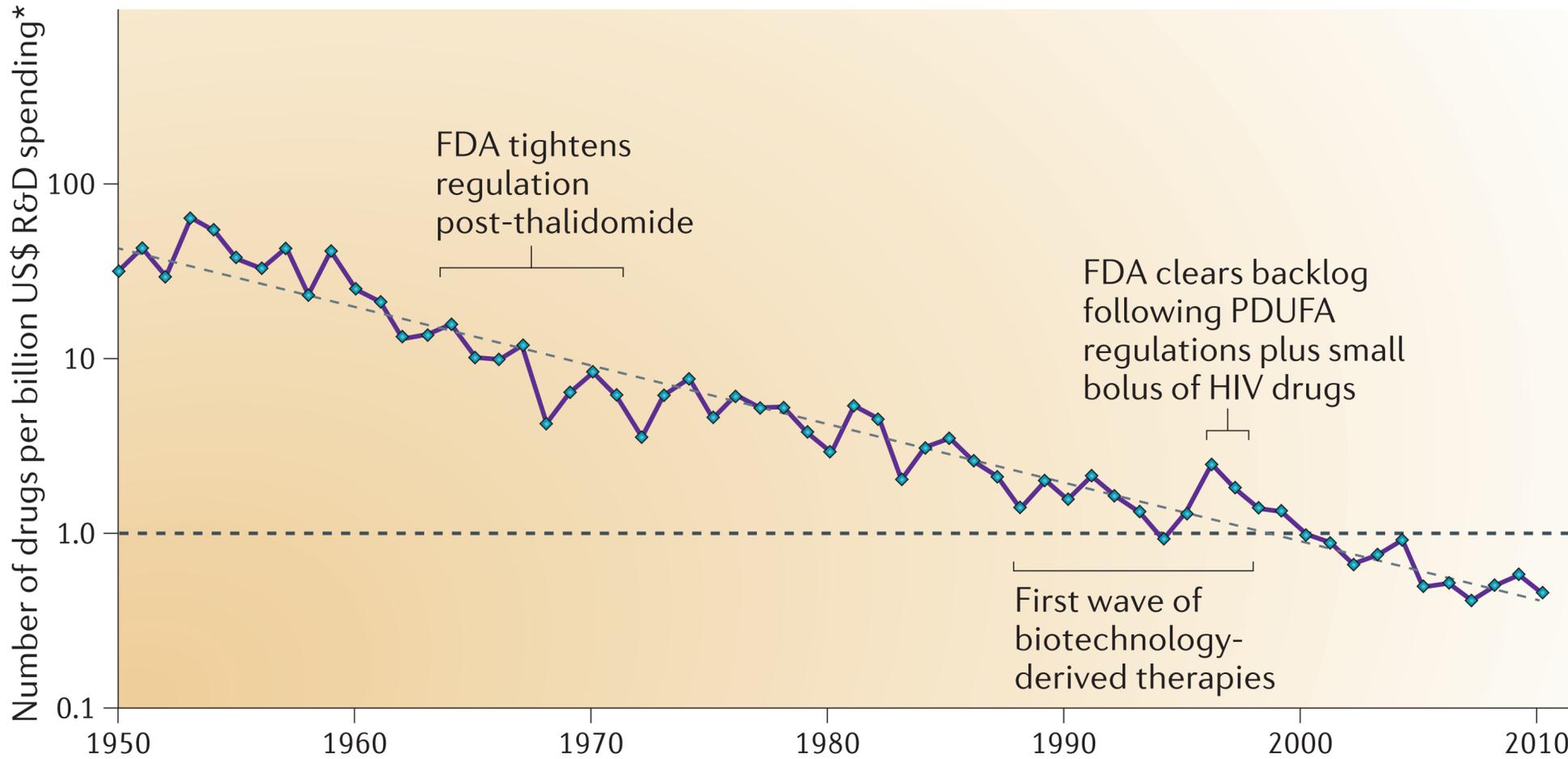
Total pharma R&D spending **doubled** to \$65B over 2000-2010

FDA approvals of new molecular entities **went down by half**

Number of truly innovative new molecules **remained constant at 5-6/year**

2010-2015 has seen large reductions in pharma R&D in the US

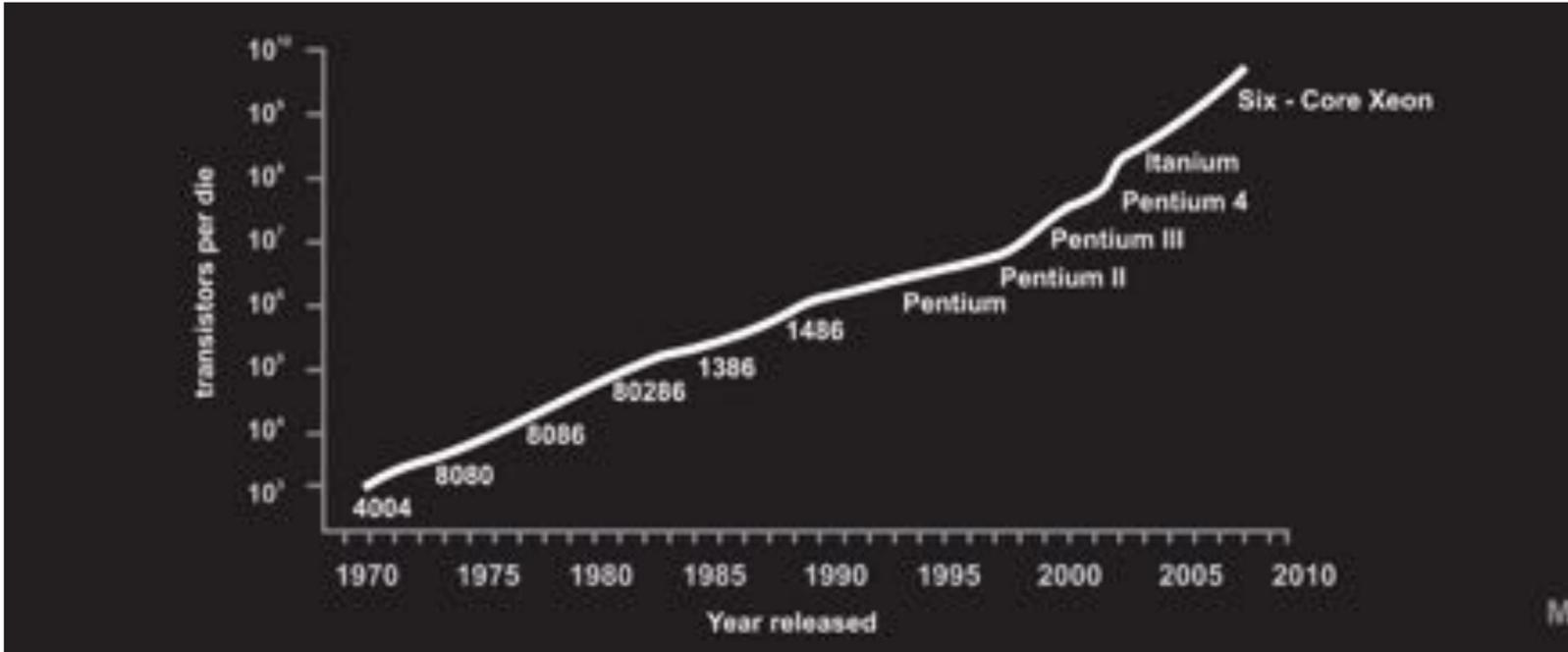
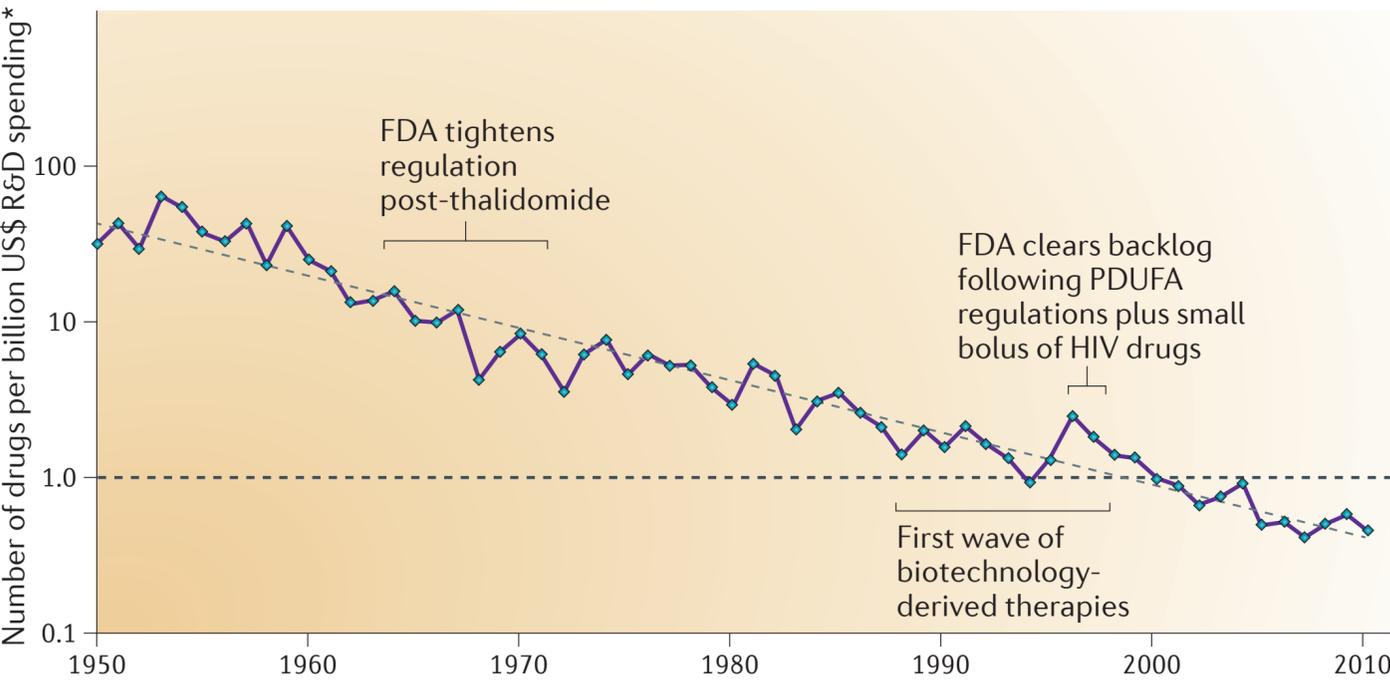
**a** Overall trend in R&D efficiency (inflation-adjusted)



# DRUG DISCOVERY USUALLY ENDS IN FAILURE



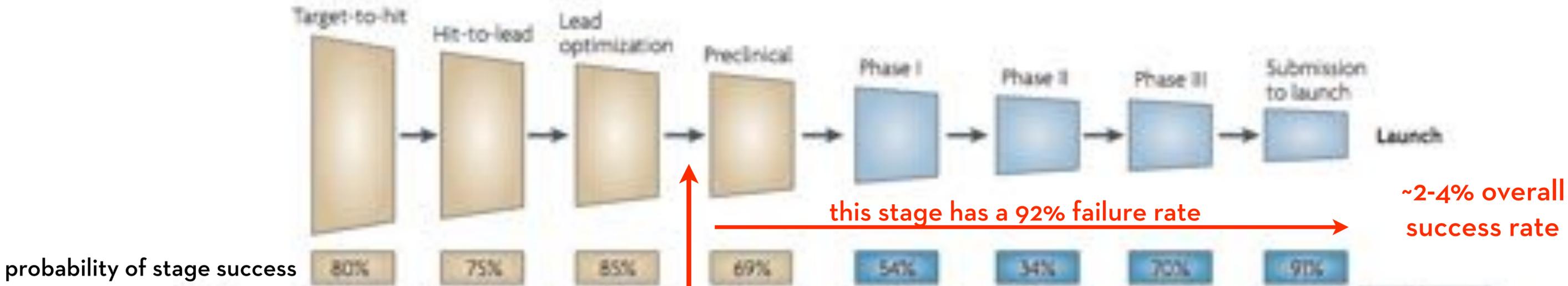
a Overall trend in R&D efficiency (inflation-adjusted)



**EROOM'S LAW**

**MOORE'S LAW**

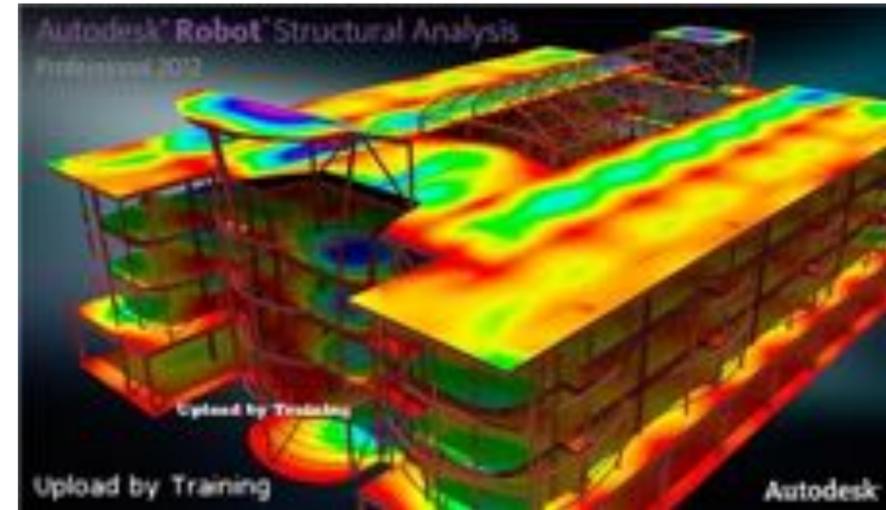
# DRUG DISCOVERY USUALLY ENDS IN FAILURE



4.5 years and \$219M

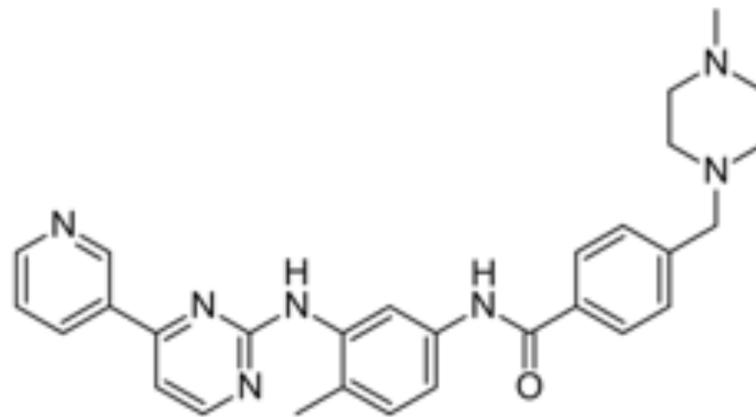
Paul et al. Nat. Rev. Drug Discover. 9:203, 2010.  
 Chodera et al. Curr. Opin. Struct. Biol., 21:150, 2011.

# WE REGULARLY **DESIGN** PLANES, BRIDGES, AND BUILDINGS ON COMPUTERS



$10^3 - 10^6$  parts

## **WHY NOT SMALL MOLECULE DRUGS?**

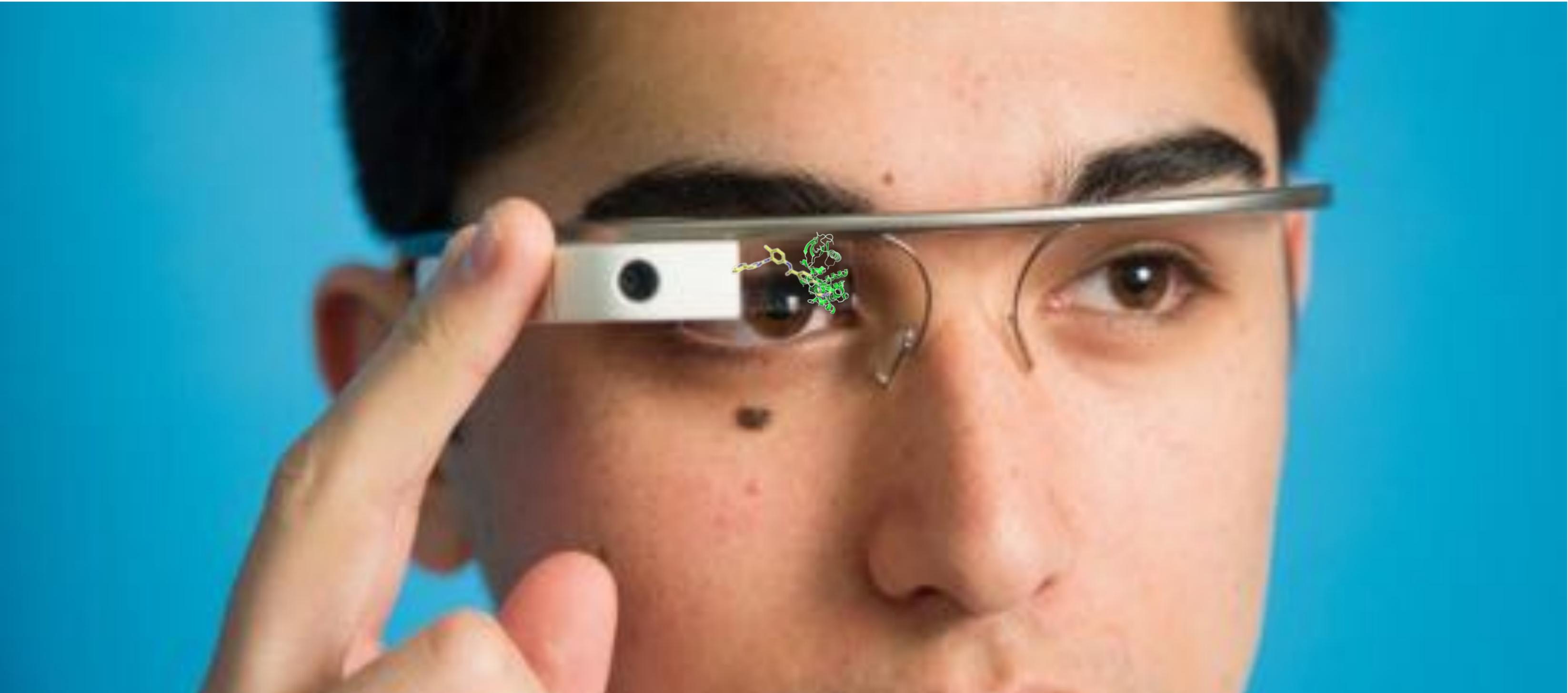


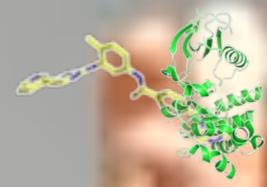
$< 10^2$  atoms

# HOW CAN WE BRING DRUG DESIGN INTO THE 21ST CENTURY?



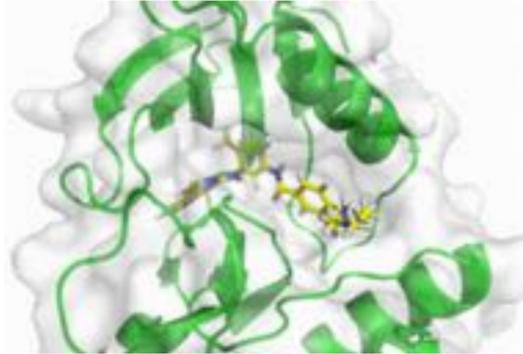
# HOW CAN WE BRING DRUG DESIGN INTO THE 21ST CENTURY?



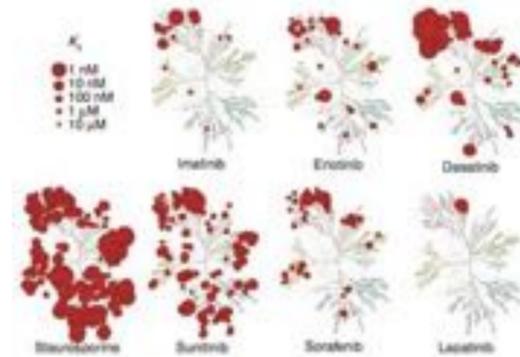


# CHODERA LAB

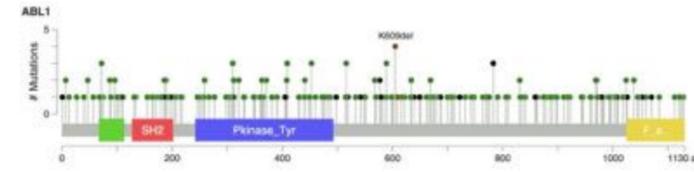
HOW CAN STATISTICAL MECHANICS PLAY A ROLE IN THE ERA OF GENOMICS AND BIOMEDICAL BIG DATA?



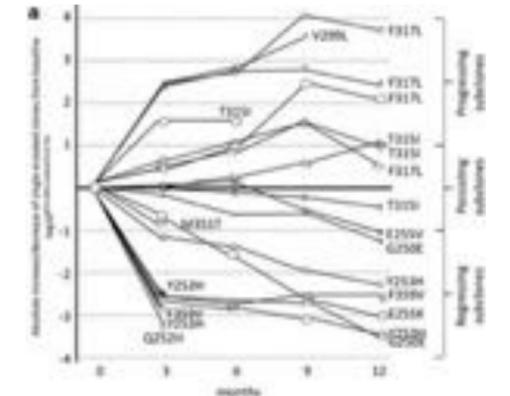
**SELECTIVE INHIBITOR DESIGN:  
TARGETS/ANTITARGETS**



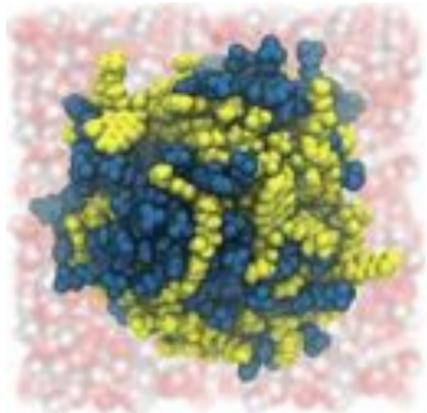
**KINASE INHIBITOR  
SELECTIVITY**



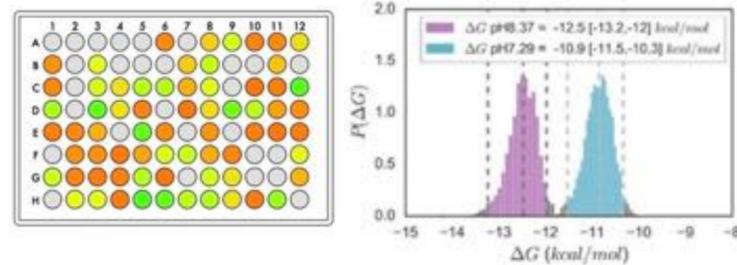
**PREDICTING DRUG  
SENSITIVITY/RESISTANCE**



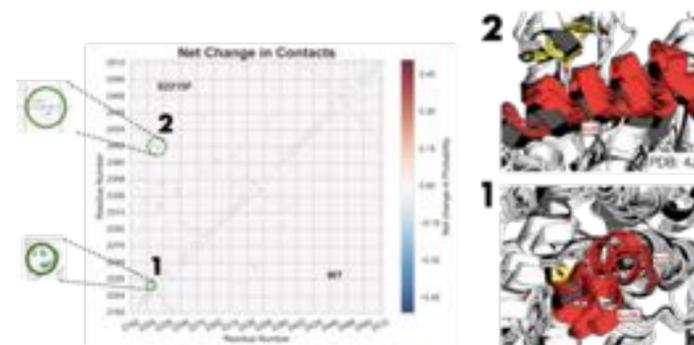
**ANTICIPATING  
DRUG RESISTANCE**



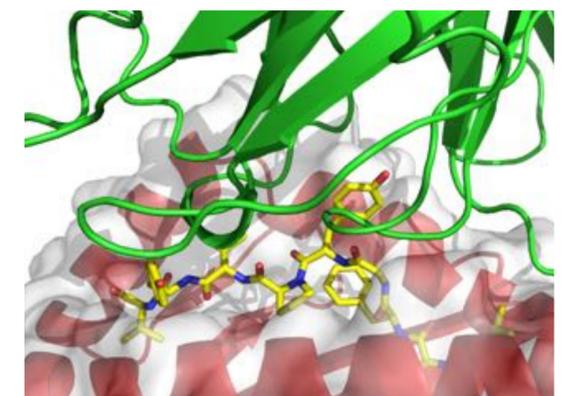
**NOVEL DRUG DELIVERY  
MODALITIES**



**AUTOMATED BIOPHYSICAL  
ASSAYS AND INFERENCE**

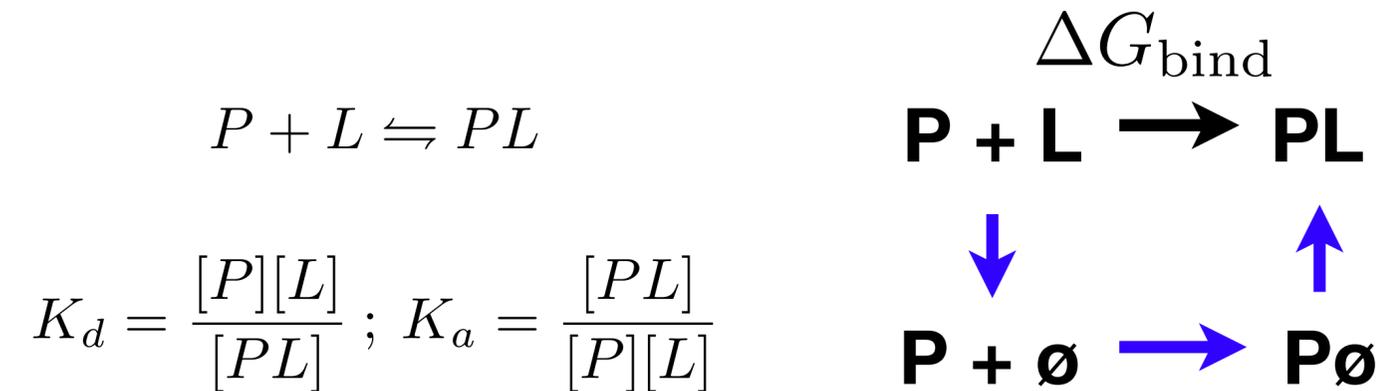


**MECHANISMS OF  
ONCOGENIC ACTIVATION**



**CANCER  
IMMUNOTHERAPY**

# COMPUTATION OF BINDING AFFINITIES TRANSFORMED INTO COMPUTATION OF RATIOS OF PARTITION FUNCTIONS



$$\Delta G_{PL}^o = -k_B T \ln \left[ \frac{C^o}{8\pi^2} \frac{Z_{PL}}{Z_{P\emptyset}} \frac{Z_{\emptyset}}{Z_L} \right]$$

(strong binding limit)

standard state  $C_0 = 1 M$

ligand in solvent  $\frac{Z_L}{Z_{\emptyset}} = \frac{\int_{\Gamma} dx e^{-\beta H_L(x; \lambda=1)}}{\int_{\Gamma} dx e^{-\beta H_L(x; \lambda=0)}}$

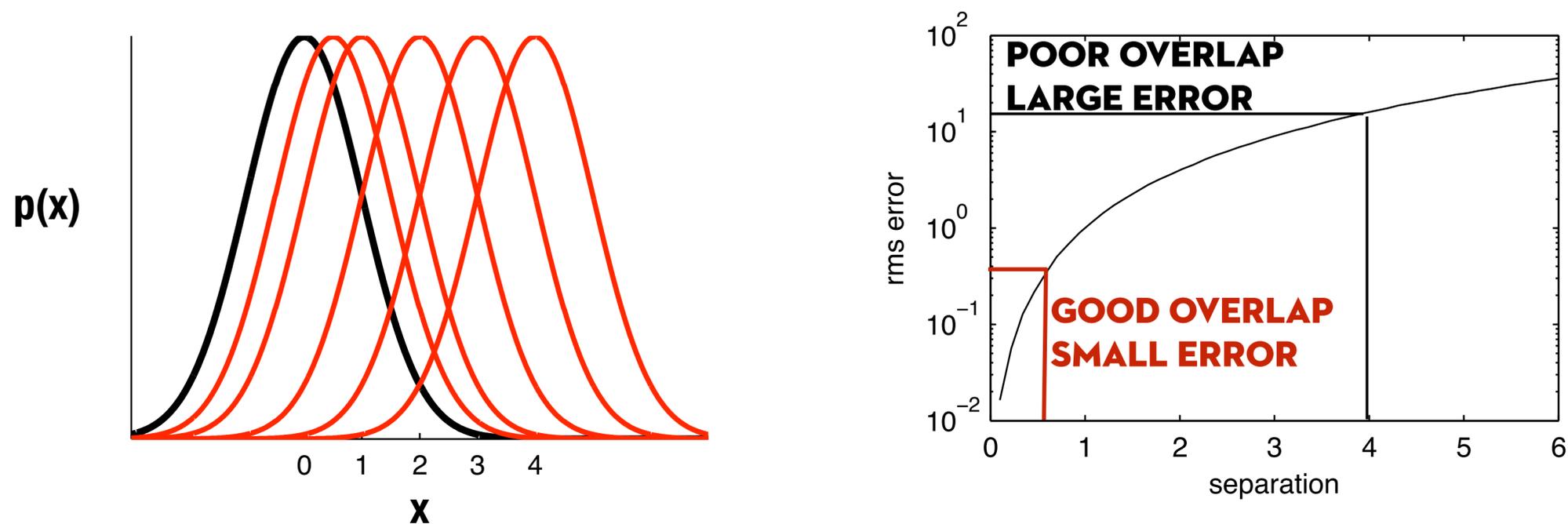
ligand in complex  $\frac{Z_{PL}}{Z_{P\emptyset}} = \frac{\int_{\Gamma} dx e^{-\beta H_{PL}(x; \lambda=1)}}{\int_{\Gamma} dx e^{-\beta H_{PL}(x; \lambda=0)}}$

**ISOMORPHIC TO COMPUTING RATIOS OF NORMALIZING CONSTANTS IN STATISTICAL INFERENCE OR RATIOS OF MODEL EVIDENCES IN INFERENCEAL MACHINE LEARNING**

**(A FIELD THAT IS BASICALLY PRINTING MONEY RIGHT NOW)**

# THE ADDITION OF **ALCHEMICAL INTERMEDIATES** GREATLY IMPROVES ESTIMATION EFFICIENCY

**COMPUTING FREE ENERGY DIFFERENCES BECOMES HARDER AS PHASE SPACE OVERLAP DECREASES; ERROR INCREASES RAPIDLY WITH DIMINISHING PHASE SPACE OVERLAP**

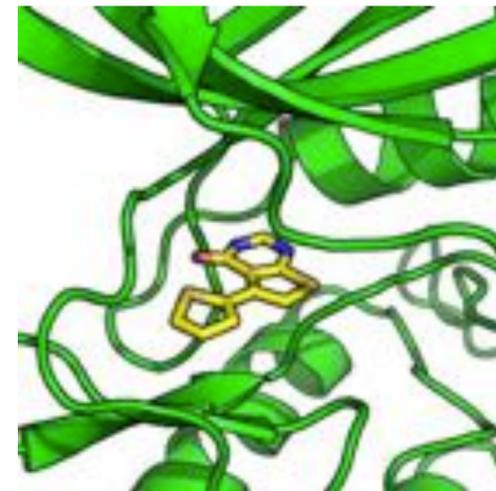
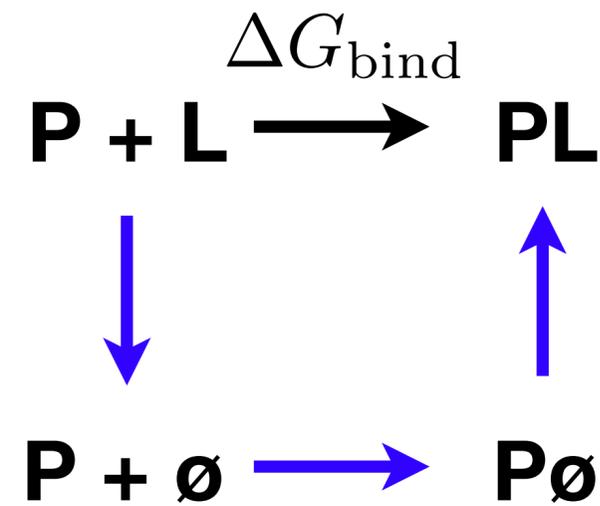


**INTRODUCING **ALCHEMICAL INTERMEDIATES** TO ENSURE RATIOS CAN BE EFFICIENTLY ESTIMATED**

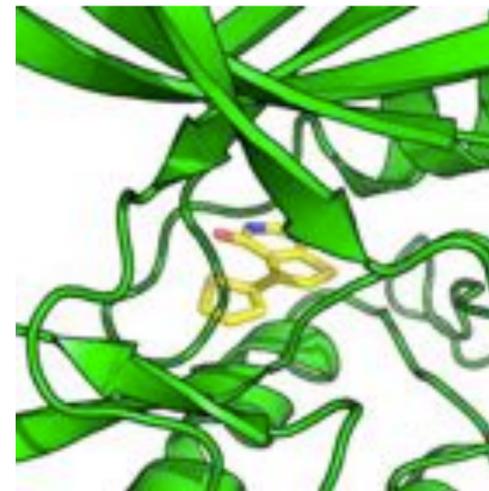
$$\Delta F_{1 \rightarrow N} = -\beta^{-1} \ln \frac{Z_N}{Z_1} = -\beta^{-1} \ln \frac{Z_2}{Z_1} \cdot \frac{Z_3}{Z_2} \cdots \frac{Z_N}{Z_{N-1}} = \sum_{n=1}^{N-1} \Delta F_{n \rightarrow n+1} \quad Z_n = \int d\mathbf{x} e^{-\beta U(\mathbf{x})}$$

# ALCHEMICAL FREE ENERGY CALCULATIONS PROVIDE A RIGOROUS WAY TO EFFICIENTLY COMPUTE BINDING AFFINITIES

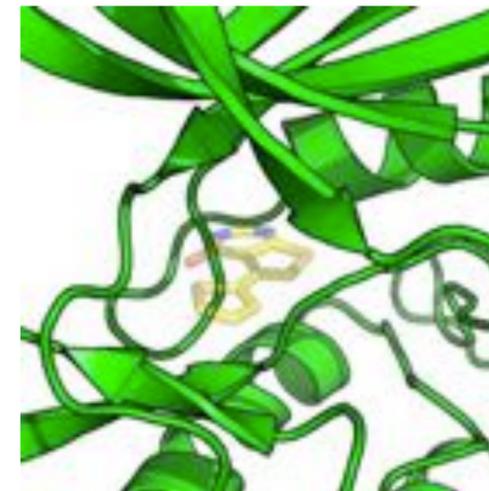
RUN 10-40 EQUILIBRIUM SIMULATIONS OF **ALCHEMICAL INTERMEDIATES**



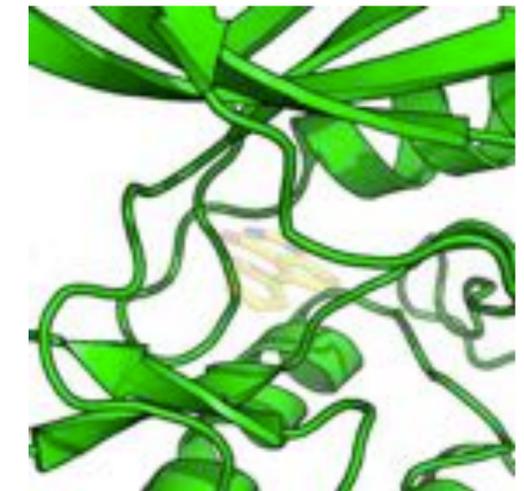
RESTRAINT IMPOSITION



DISCHARGING



STERIC DECOUPLING



NONINTERACTING

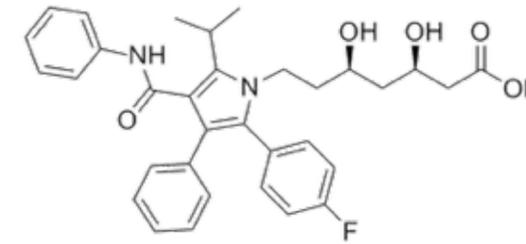
SAMPLE  $x_{kn} \sim \pi(x; \lambda_k)$  FROM  $\pi(x; \lambda) = Z_k^{-1} e^{-\beta H(x; \lambda)}$

REDUCES EFFORT BY **ORDERS OF MAGNITUDE**  
OVER SIMULATING DIRECT ASSOCIATION PROCESS

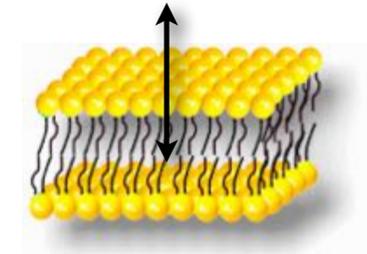
$$\Delta F_{1 \rightarrow N} = -\beta^{-1} \ln \frac{Z_N}{Z_1} = -\beta^{-1} \ln \frac{Z_2}{Z_1} \cdot \frac{Z_3}{Z_2} \cdots \frac{Z_N}{Z_{N-1}} = \sum_{n=1}^{N-1} \Delta F_{n \rightarrow n+1} \quad Z_n = \int d\mathbf{x} e^{-\beta U(\mathbf{x})}$$

# THE SAME ALCHEMICAL METHODS CAN ALSO COMPUTE MANY OTHER USEFUL PROPERTIES

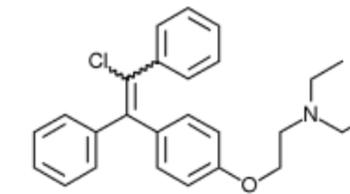
**PARTITION COEFFICIENTS (LOGP, LOGD) AND PERMEABILITIES**



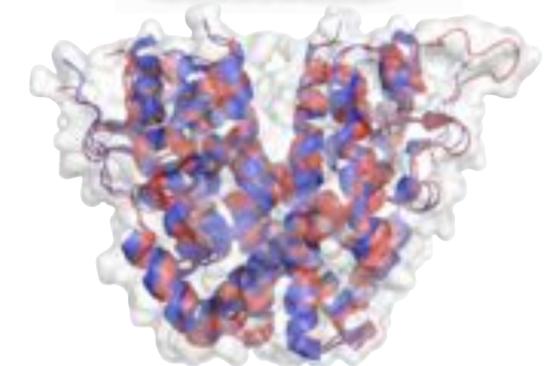
lipitor



**SELECTIVITY FOR SUBTYPES OR RELATED TARGETS/OFF-TARGETS**



clomifene

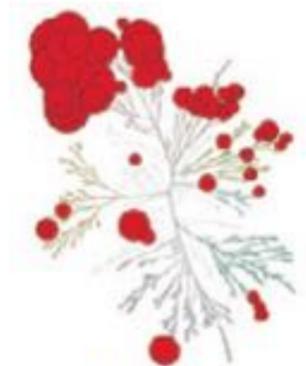


ER $\alpha$ / $\beta$

**LEAD OPTIMIZATION OF AFFINITY AND SELECTIVITY**

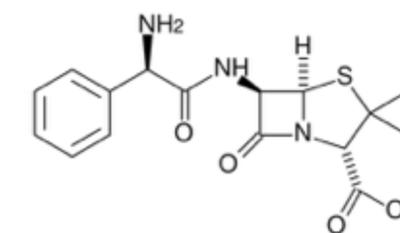


Imatinib

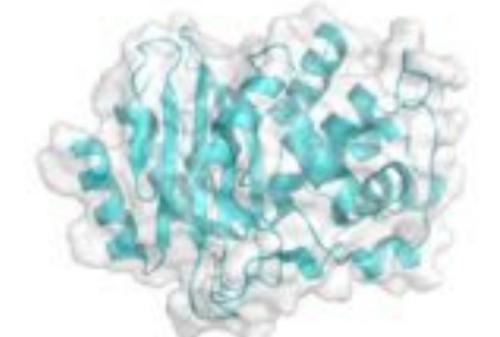


Dasatinib

**SUSCEPTIBILITY TO RESISTANCE MUTATIONS**



ampicillin



$\beta$ -lactamase

**ALSO SOLUBILITIES, POLYMORPHS, ETC.**

# **WHAT DO WE NEED TO DO TO ENABLE DESIGN BY FREE ENERGY METHODS?**

- 1. Speed up calculations without sacrificing fidelity**
- 2. Include relevant chemical effects**
- 3. Use a sufficiently accurate forcefield**

# WHAT **TIMESTEP**\* SHOULD I USE?

- \* **WHICH INTEGRATOR?**
- \* **HYDROGEN MASS REPARTITIONING?**
- \* **MULTIPLE TIMESTEP INTEGRATION?**
- \* **WHICH FORCE SPLITTING?**
- \* **SOLUTE/SOLVENT SPLITTING?**

# TO INTEGRATE LANGEVIN DYNAMICS ON A COMPUTER, WE HAVE TO DISCRETIZE EQUATIONS OF MOTION



## continuous Langevin dynamics

$$dr = v dt$$

$$dv = \frac{f(t)}{m} dt - \gamma v dt + \sqrt{\frac{2\gamma}{\beta m}} dW(t)$$



## discrete timestep Langevin integrator

$$v'_t = v_t^* + \frac{\Delta t}{2m} \left( F_t(r_t^*) - \gamma m v_t^* + \sqrt{\frac{2\gamma m}{\Delta t}} \xi_t \right)$$

$$r_t = r_t^* + \Delta t v'_t$$

$$v_t = \frac{1}{1 + \frac{\gamma \Delta t}{2}} \left[ v'_t + \frac{\Delta t}{2m} \left( F_t(r_t) + \sqrt{\frac{2\gamma m}{\Delta t}} \xi'_t \right) \right]$$

# NOT ALL INTEGRATORS ARE EQUAL

Continuous Langevin dynamics equations of motion

$$d \begin{bmatrix} \mathbf{x} \\ \mathbf{v} \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{v} \\ 0 \end{bmatrix}}_R dt + \underbrace{\begin{bmatrix} 0 \\ -M^{-1}\nabla U(\mathbf{x}) \end{bmatrix}}_V dt + \underbrace{\begin{bmatrix} 0 \\ -\gamma\mathbf{v}dt + \sigma M^{1/2}dW \end{bmatrix}}_O$$

Some of these parts can be integrated exactly:

$$(e^{t\mathcal{L}_R}\phi)(\mathbf{q}, \mathbf{p}) = \phi(\mathbf{q} + tM^{-1}\mathbf{p}, \mathbf{p})$$

$$(e^{t\mathcal{L}_V}\phi)(\mathbf{q}, \mathbf{p}) = \phi(\mathbf{q}, \mathbf{p} - t\nabla U(\mathbf{q}))$$

$$(e^{t\mathcal{L}_O}\phi)(\mathbf{q}, \mathbf{p}) = \int_{\mathcal{P}} \phi(\mathbf{q}, e^{-\gamma t}\mathbf{p} + \eta M^{1/2}\mathbf{x}) \frac{e^{-|\mathbf{x}|^2/2}}{(2\pi)^{N/2}} d\mathbf{x}$$



$$R : (\mathbf{x}, \mathbf{v}; h) \mapsto (\mathbf{x} + \mathbf{v}h, \mathbf{v})$$

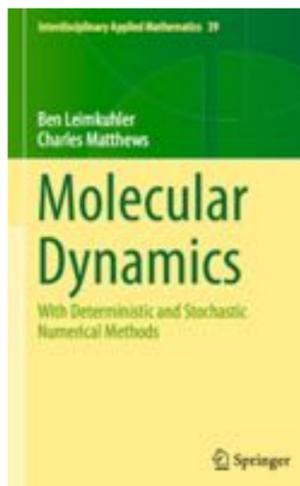
$$V : (\mathbf{x}, \mathbf{v}; h) \mapsto (\mathbf{x}, \mathbf{v} + (f(\mathbf{x})/m)h)$$

$$O : (\mathbf{x}, \mathbf{v}; h) \mapsto (\mathbf{x}, a_h\mathbf{v} + b_h\sqrt{k_B T/m}\xi)$$

We can approximate the Langevin propagator by various splittings:

$$e^{\Delta t[\mathcal{L}_O + \mathcal{L}_V + \mathcal{L}_R + \mathcal{L}_h]} \simeq e^{\frac{\Delta t}{2}\mathcal{L}_R} e^{\frac{\Delta t}{2}\mathcal{L}_h} e^{\frac{\Delta t}{2}\mathcal{L}_V} e^{\Delta t\mathcal{L}_O} e^{\frac{\Delta t}{2}\mathcal{L}_V} e^{\frac{\Delta t}{2}\mathcal{L}_h} e^{\frac{\Delta t}{2}\mathcal{L}_R} \quad \mathbf{RVOVR}$$

$$e^{\Delta t[\mathcal{L}_O + \mathcal{L}_V + \mathcal{L}_R + \mathcal{L}_h]} \simeq e^{\frac{\Delta t}{2}\mathcal{L}_R} e^{\frac{\Delta t}{2}\mathcal{L}_O} e^{\frac{\Delta t}{2}\mathcal{L}_h} e^{\Delta t\mathcal{L}_V} e^{\frac{\Delta t}{2}\mathcal{L}_h} e^{\frac{\Delta t}{2}\mathcal{L}_O} e^{\frac{\Delta t}{2}\mathcal{L}_R} \quad \mathbf{ROVOR}$$



# NOT ALL INTEGRATORS ARE EQUAL

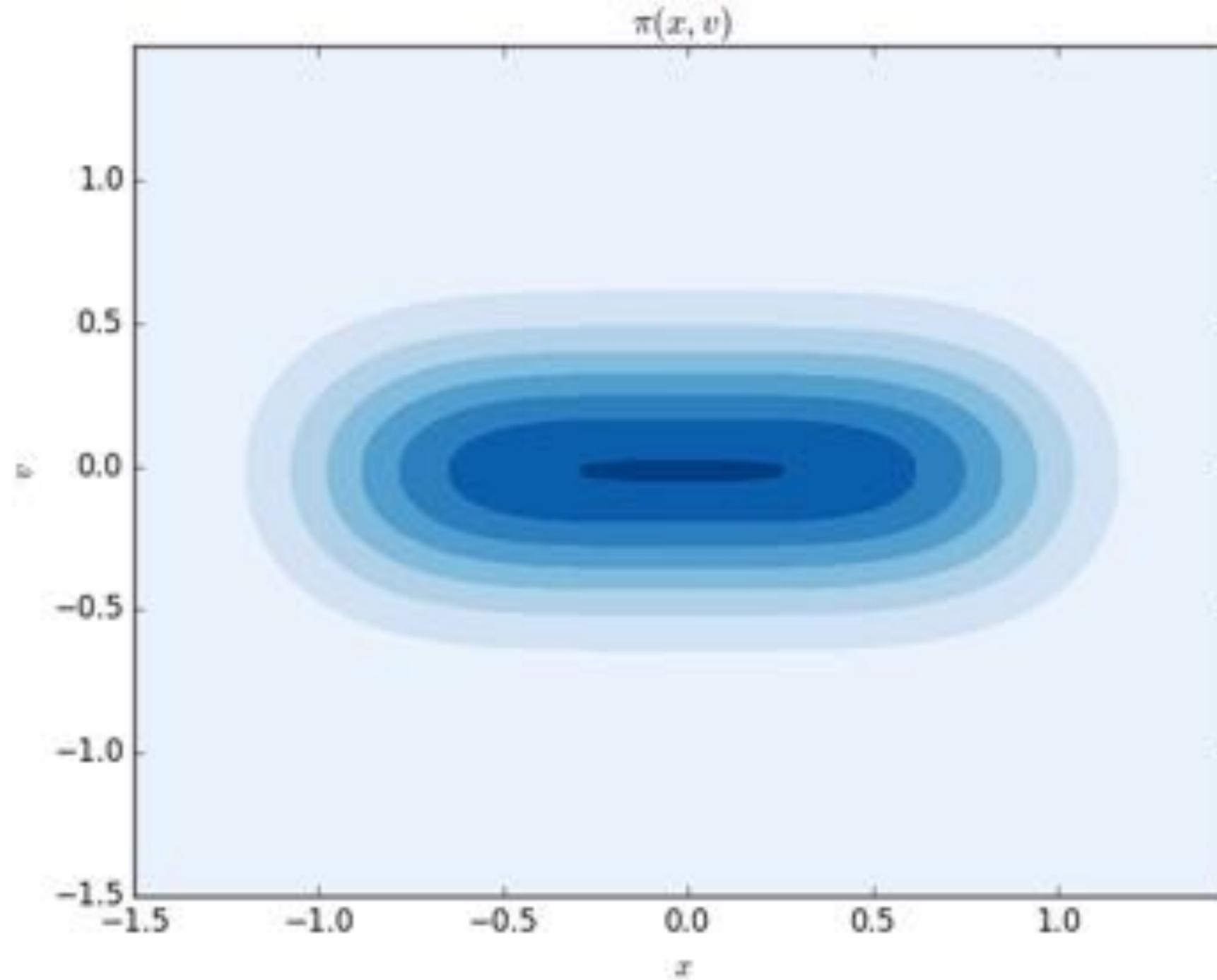
Some integrators are superior for purely kinetic properties:

Desideratum	OVRVO	ORVRO	RVQVR	VROIV	VOROV	ROVOR
zero-force MSD	exact	exact	exact	exact	exact	exact
zero-force MSV	exact	exact	exact	exact	exact	exact
zero-force VAC	exact	exact	exact	exact	exact	exact
uniform-force terminal drift	exact	exact	exact	exact	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error
linear-force MSD	exact at $n + \frac{1}{2}$	exact at $n + \frac{1}{2}$	exact at $n$	exact at $n$	$O(\Delta t^3)$ error	$O(\Delta t^2)$ error
linear-force MSV	exact at $n$	exact at $n$	exact at $n + \frac{1}{2}$	exact at $n + \frac{1}{2}$	$O(\Delta t^4)$ error at $n$	$O(\Delta t^4)$ error at $n + \frac{1}{2}$
linear-force virial	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error	$O(\Delta t^2)$ error
irreducible path action	yes simple	yes requires values at $n + \frac{1}{2}$	no may be infinite	no may be infinite	yes simple	no may be infinite
maintains Hamiltonian dependence for large $\gamma\Delta t$	yes	yes	yes	yes	no	no
for reducible dynamics, needs half as many random variables	yes	yes	yes	yes	no	no
in various limits, reduces to other popular integrators	yes	no	no	no	no	no

TABLE II. Comparison of properties for different splittings.

But what if we're only interested in configurational properties (like free energies)?

# ALL INTEGRATORS INTRODUCE TIMESTEP-DEPENDENT ERROR

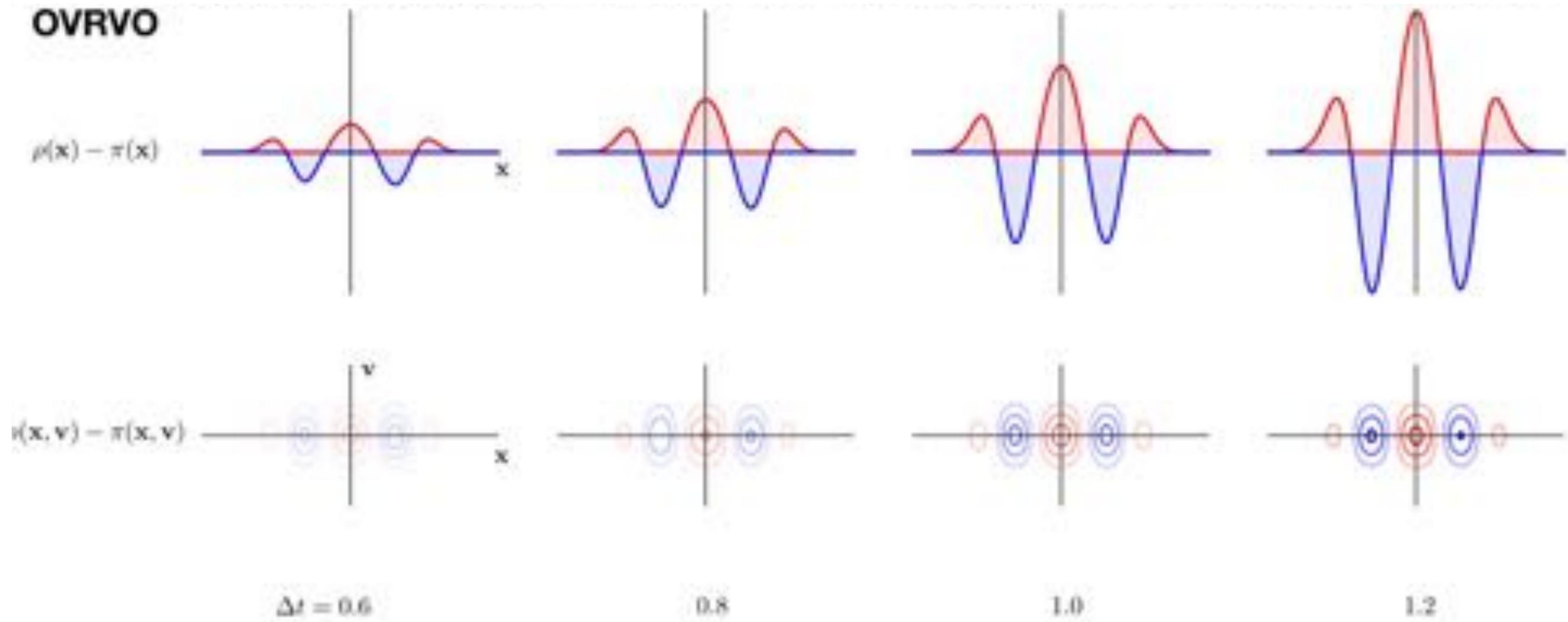


**SIMPLE  
QUARTIC  
POTENTIAL**



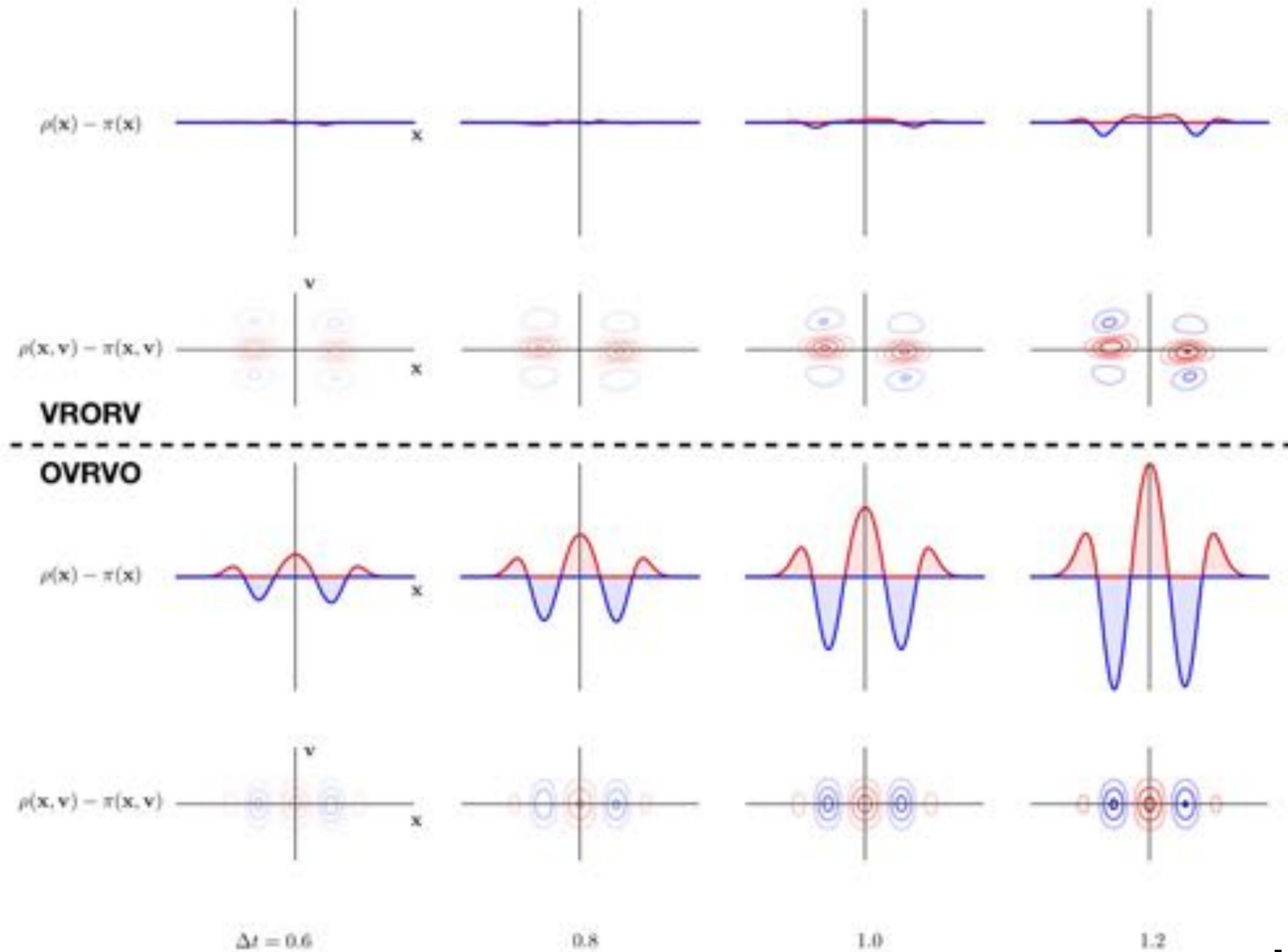
**JOSH FASS**

# ALL INTEGRATORS INTRODUCE TIMESTEP-DEPENDENT ERROR

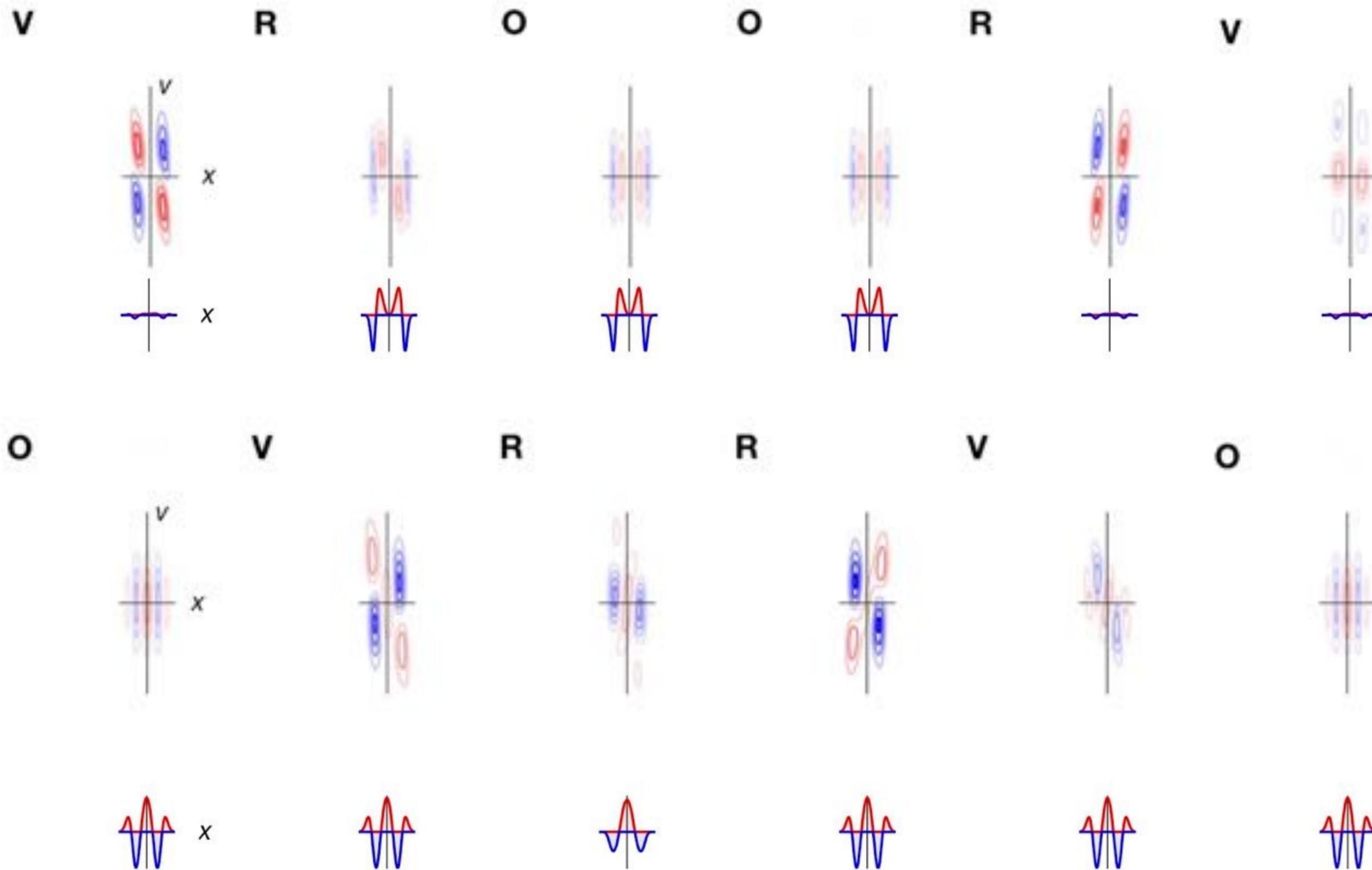


**JOSH FASS**

# NOT ALL INTEGRATORS ARE EQUAL



# ERRORS MOVE AROUND AT DIFFERENT STAGES OF THE INTEGRATOR CYCLE



**JOSH FASS**

# WHAT IF WE COULD MEASURE INTEGRATOR ERROR VIA A COMMON YARDSTICK?



Here's an example from the literature where the authors pick their favorite order parameters:

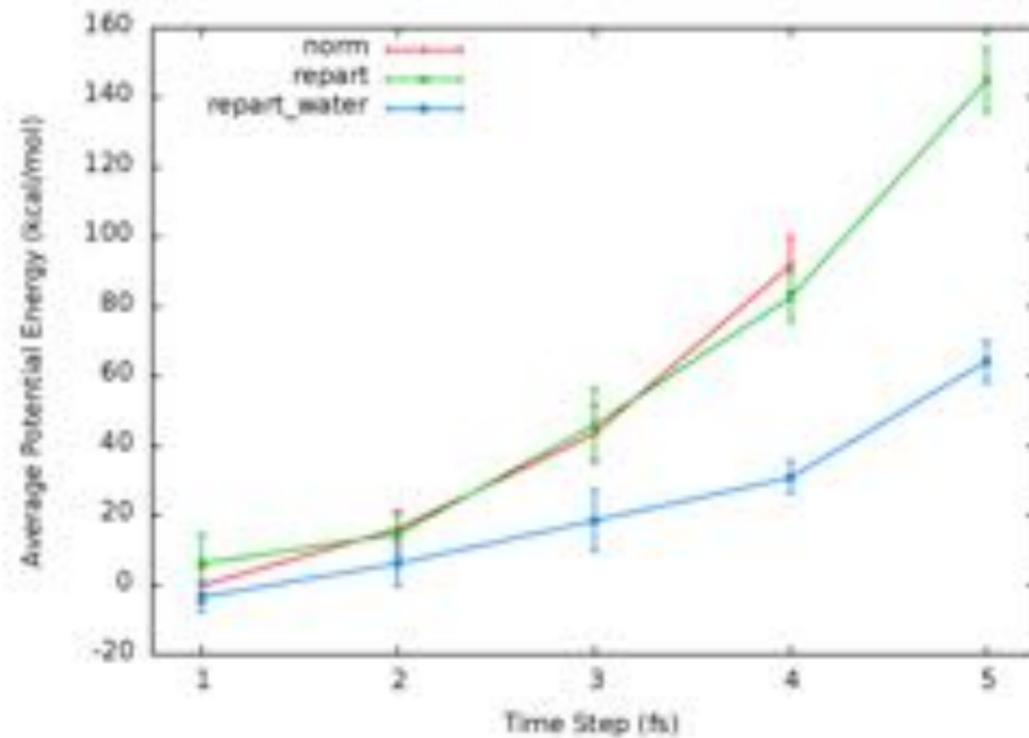


Figure 7. Average potential energies with different mass topology/time step combinations for the  $(Ala)_3$  peptide, relative to the norm/1 fs average, averaged over 10 trajectories of each trajectory type. The error bars show the standard deviation in the mean over the 10 trajectories.

**Protein-Structure.** Figure 8 shows a histogram of the total backbone root-mean-square deviation (RMSD) to the 4LYT

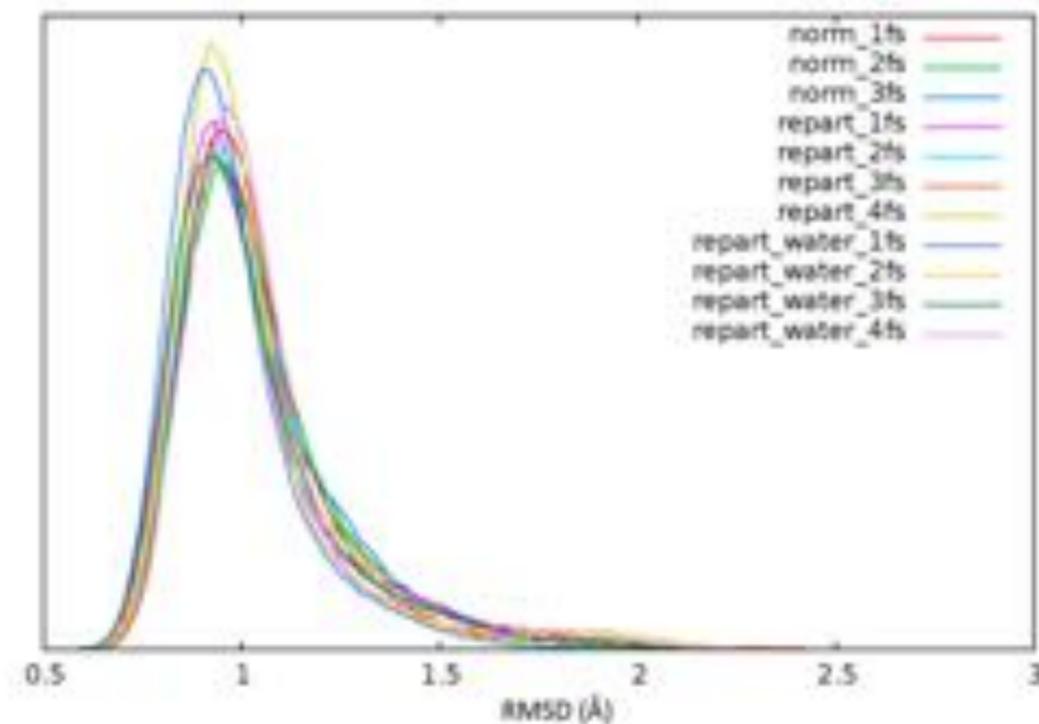


Figure 8. Backbone root-mean-square deviation (RMSD) to a crystal structure in HEWL for each trajectory type.

# WHAT IF WE COULD MEASURE INTEGRATOR ERROR VIA A COMMON YARDSTICK?

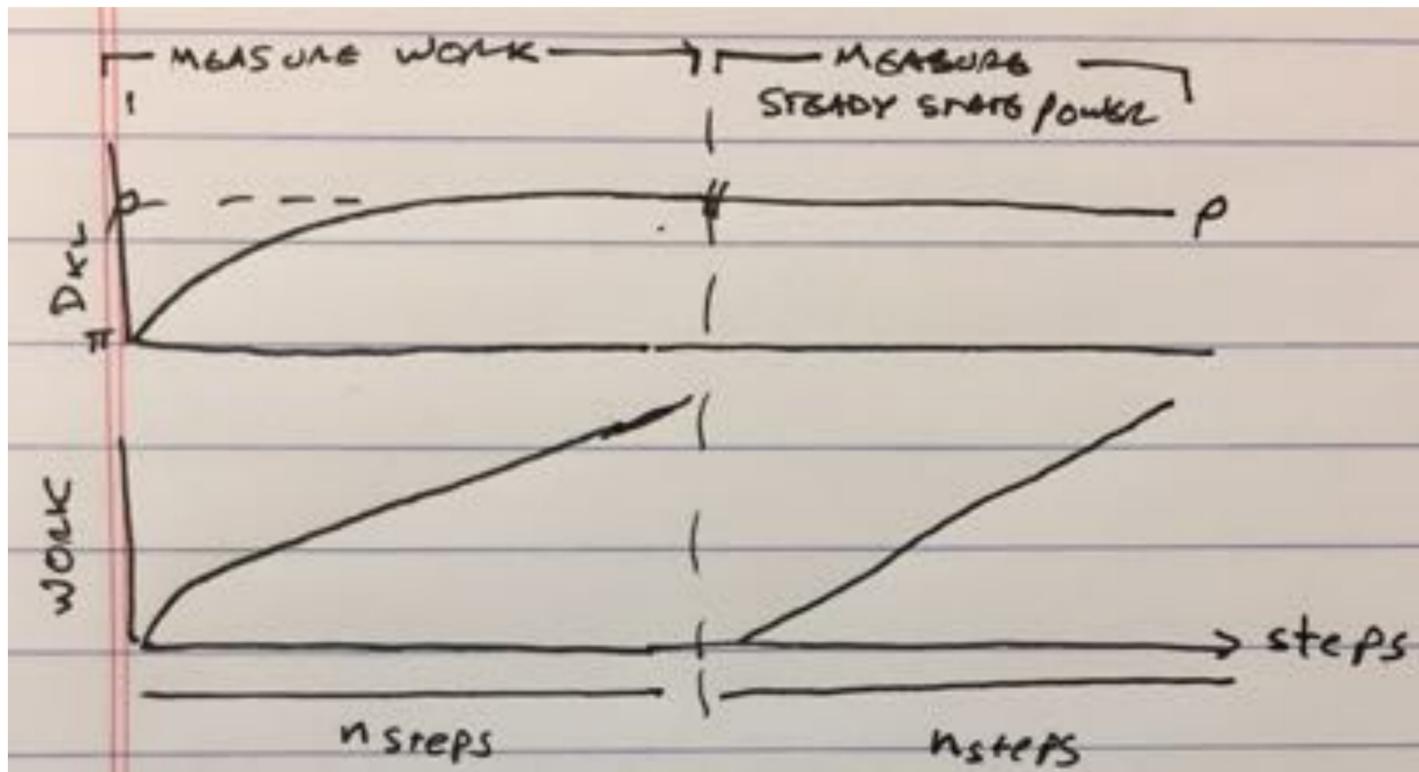


There is: the KL-divergence, which is isomorphic with a nonequilibrium free energy difference from equilibrium

$$\Delta F_{neq}(\Delta t) = (\langle E \rangle_{\Delta t} - TS_{\Delta t}) - (\langle E \rangle - TS) = \mathcal{D}_{\text{KL}}(\rho_{\Delta t} \parallel \pi) \equiv \int d\mathbf{x} \int d\mathbf{v} \rho_{\Delta t}(\mathbf{x}, \mathbf{v}) \ln \frac{\rho_{\Delta t}(\mathbf{x}, \mathbf{v})}{\pi(\mathbf{x}, \mathbf{v})}$$

We can estimate the nonequilibrium free energy deviation:

$$\Delta F_{neq} \equiv F_{neq} - F_{eq} = \frac{1}{2} [\langle W_{\text{shad}} \rangle - (t_f - t_i) \mathcal{P}_{ss}]$$



David A. Sivak  
Simon Fraser



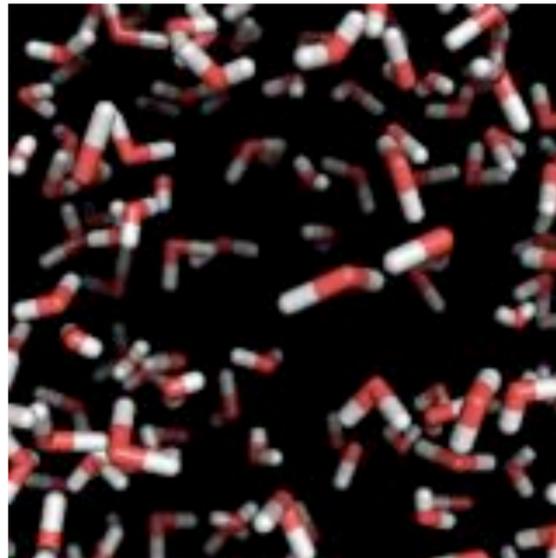
Gavin E. Crooks  
Rigetti Computing



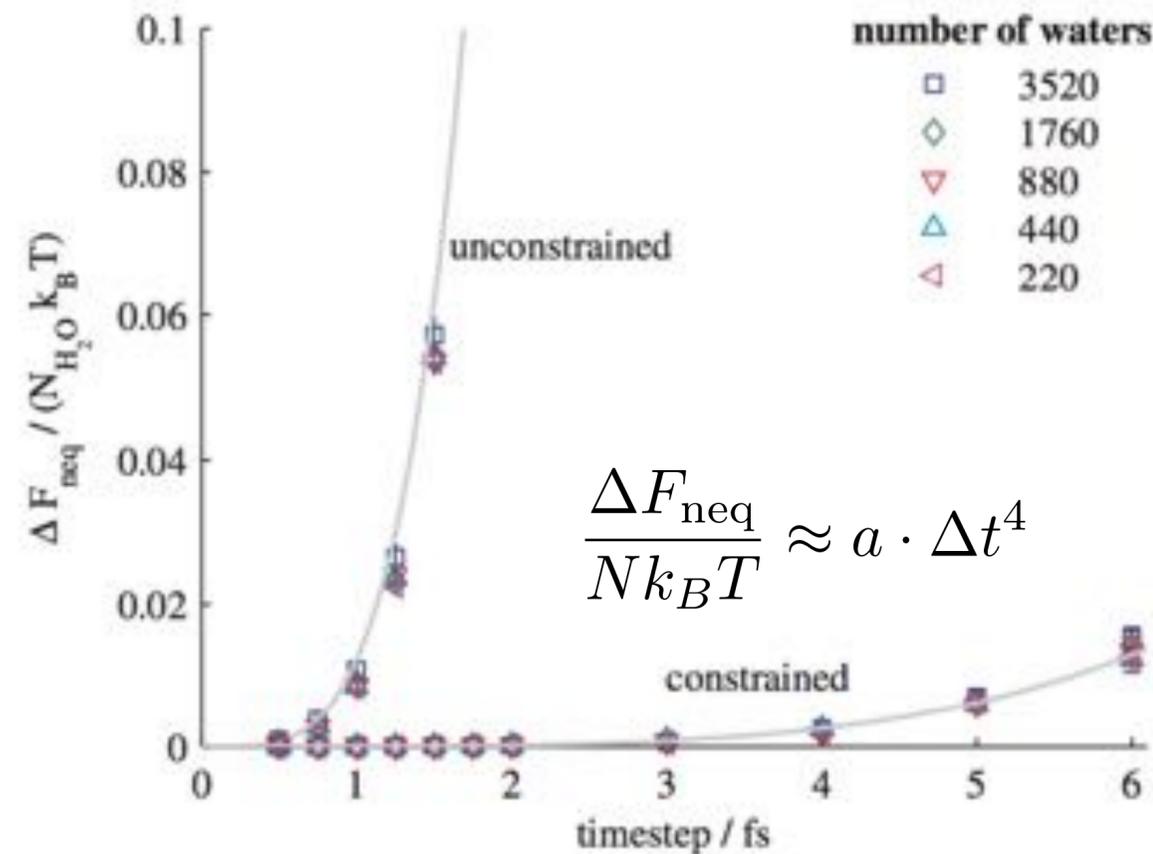
# WHAT IF WE COULD MEASURE INTEGRATOR ERROR VIA A COMMON YARDSTICK?



$$\Delta F_{neq}(\Delta t) = (\langle E \rangle_{\Delta t} - TS_{\Delta t}) - (\langle E \rangle - TS) = \mathcal{D}_{KL}(\rho_{\Delta t} \parallel \pi) \equiv \int d\mathbf{x} \int d\mathbf{v} \rho_{\Delta t}(\mathbf{x}, \mathbf{v}) \ln \frac{\rho_{\Delta t}(\mathbf{x}, \mathbf{v})}{\pi(\mathbf{x}, \mathbf{v})}$$



220 TIP3P waters at 298K and 1 atm

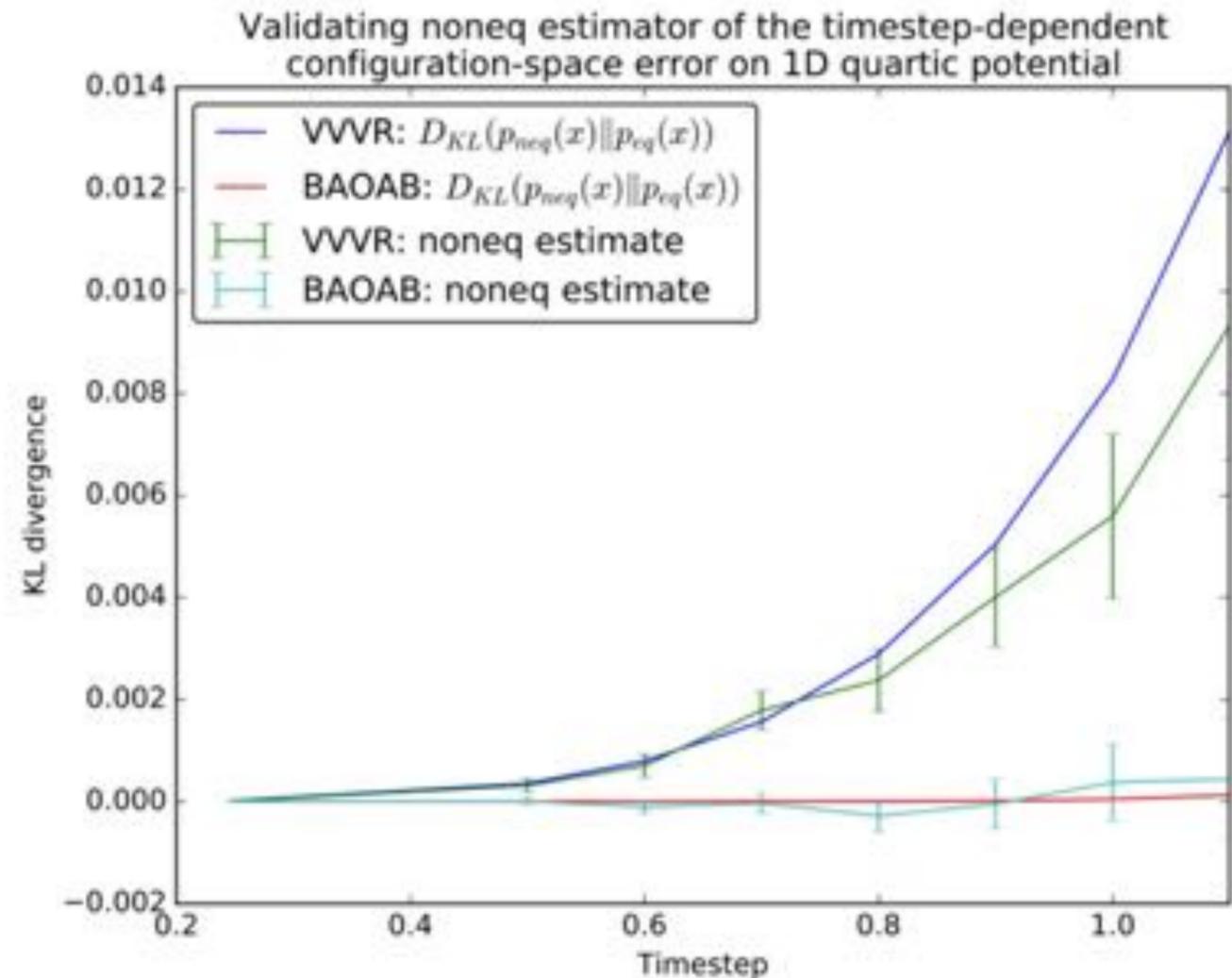
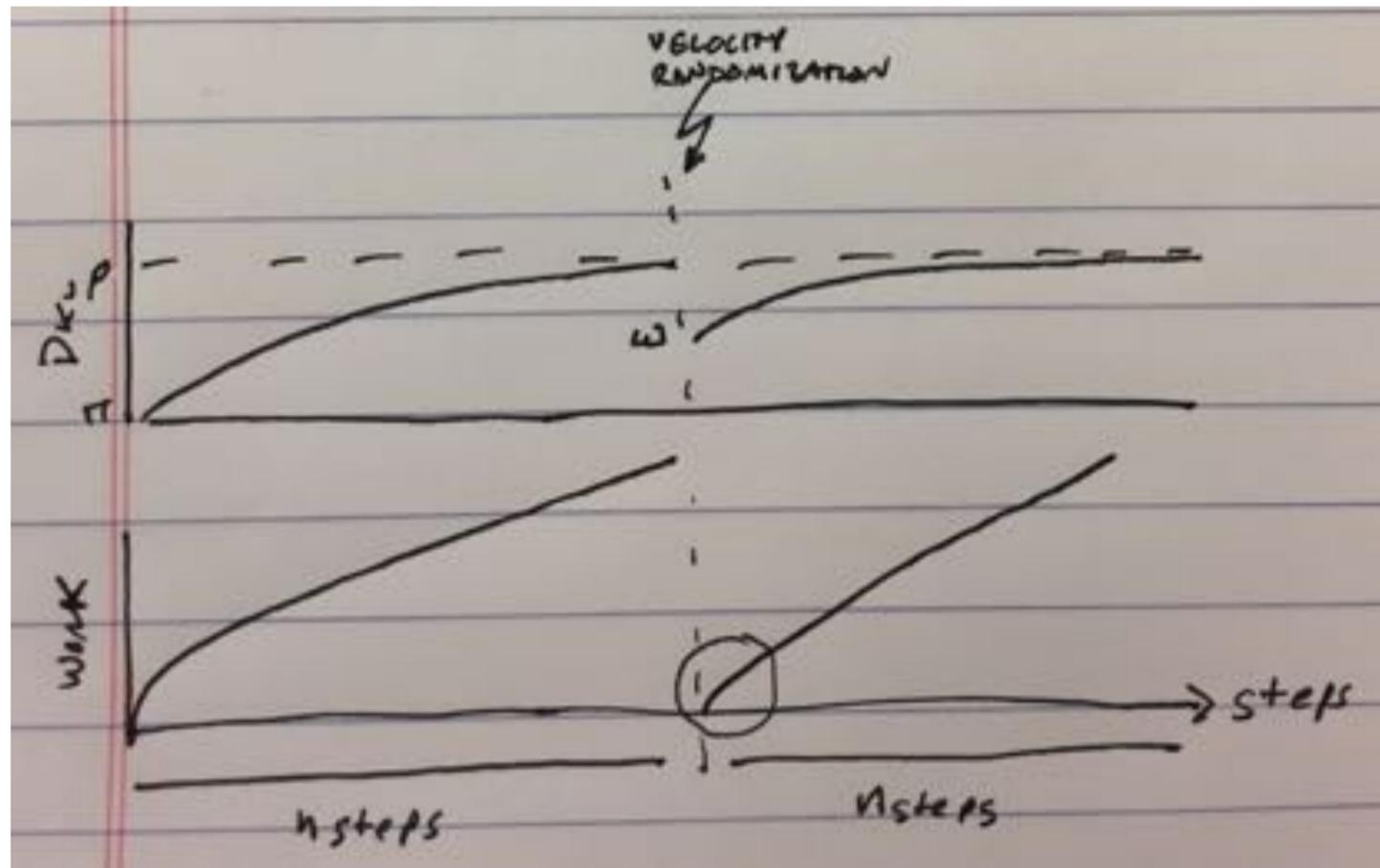


Phase space error grows as  $O(\Delta t^4)$ !

**But what about configuration space error?**



# WE CAN MEASURE ERROR IN JUST CONFIGURATION SPACE TOO



Gavin E. Crooks  
Rigetti Computing

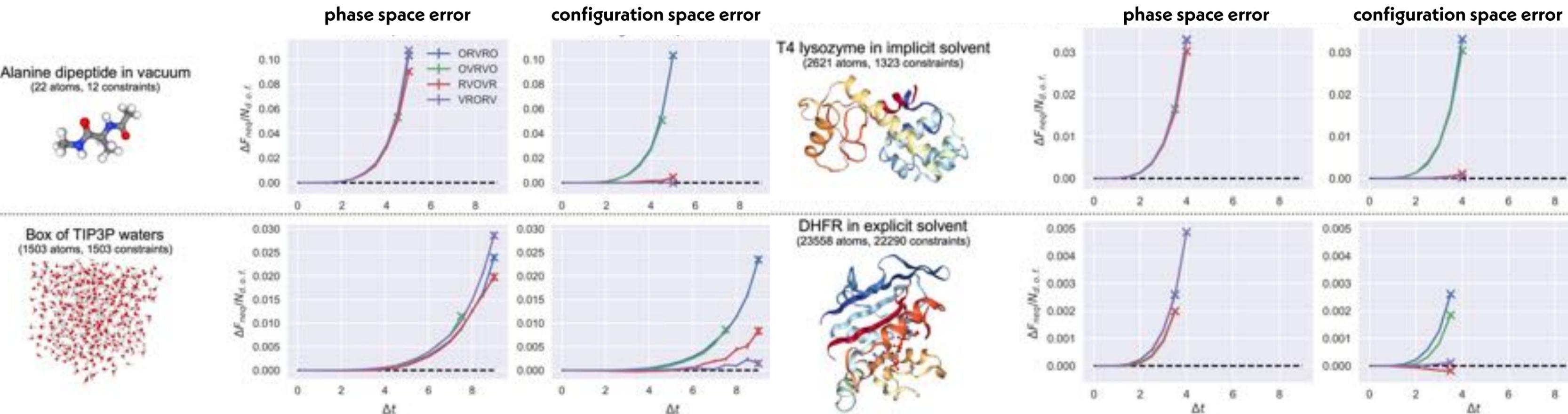
1. We can actually measure configuration space error!
2. The error of BAOAB (VRORV) seems really, really small...

JOSH FASS



# CERTAIN INTEGRATOR SPLITTINGS LEAD TO EXCEPTIONALLY SMALL CONFIGURATION-SPACE ERROR

$$\Delta F_{neq}(\Delta t) = (\langle E \rangle_{\Delta t} - TS_{\Delta t}) - (\langle E \rangle - TS) = \mathcal{D}_{\text{KL}}(\rho_{\Delta t} \| \pi) \equiv \int d\mathbf{x} \int d\mathbf{v} \rho_{\Delta t}(\mathbf{x}, \mathbf{v}) \ln \frac{\rho_{\Delta t}(\mathbf{x}, \mathbf{v})}{\pi(\mathbf{x}, \mathbf{v})}$$



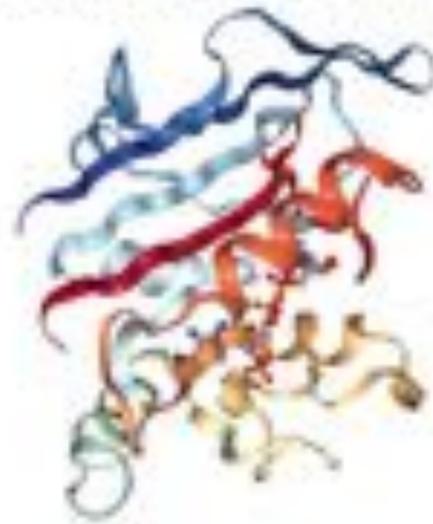
Take-home messages:

1. Error in VRORV almost too small to measure; MUCH better than other integrators
2. We can probably take larger timesteps than we're used to, but it really depends on the system!  
2 fs is not "good enough for everything": sometimes it's too big, often too small



# MULTIPLE TIMESTEP METHODS CAN HELP!

DHFR in explicit solvent  
(23558 atoms, 22290 constraints)

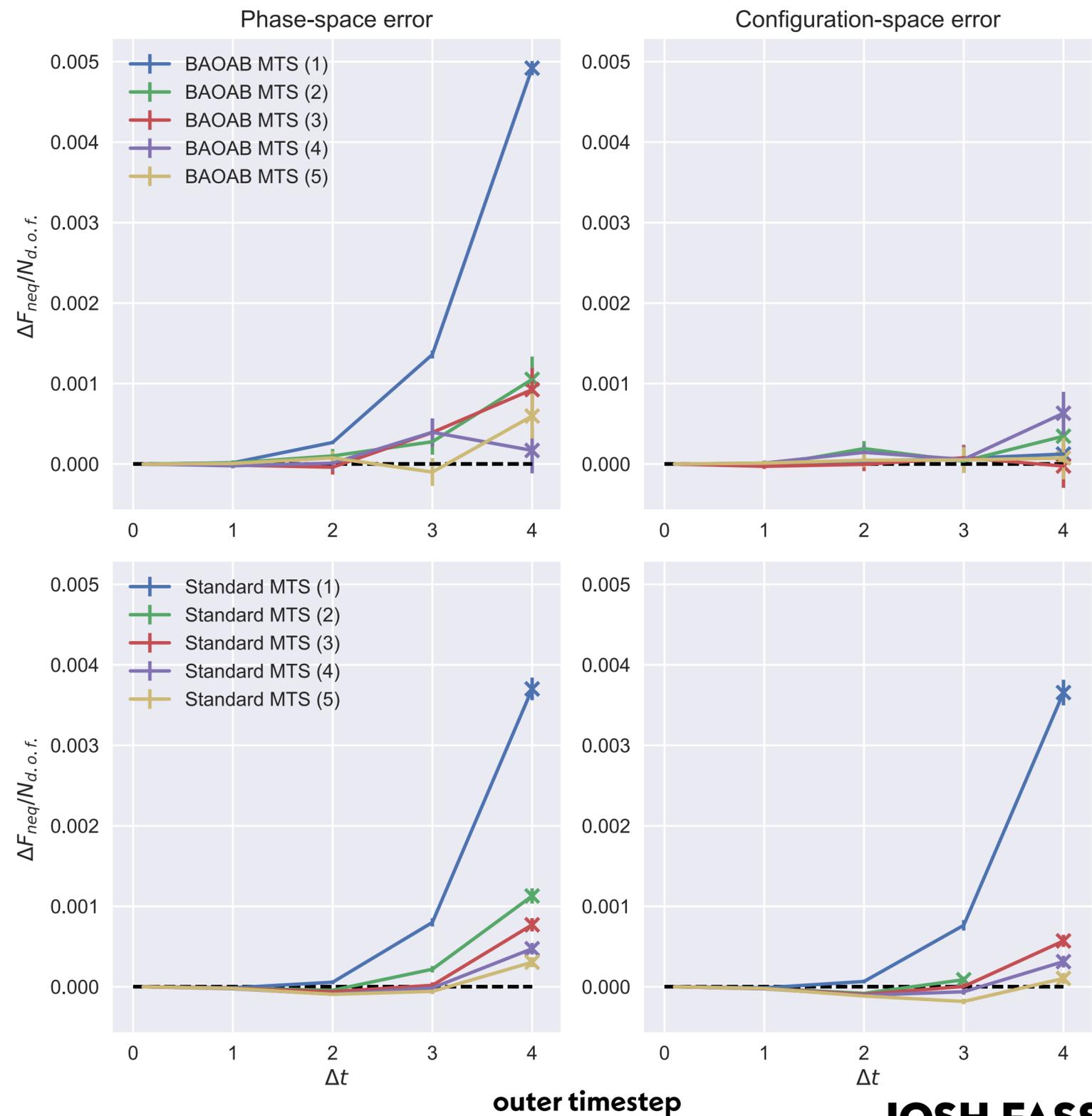


Multiple timestep scheme:

- inner timestep: valence forces
- outer timestep: nonbonded forces

GAFF / AMBER99SB-ILDN / TIP3P  
PME electrostatics

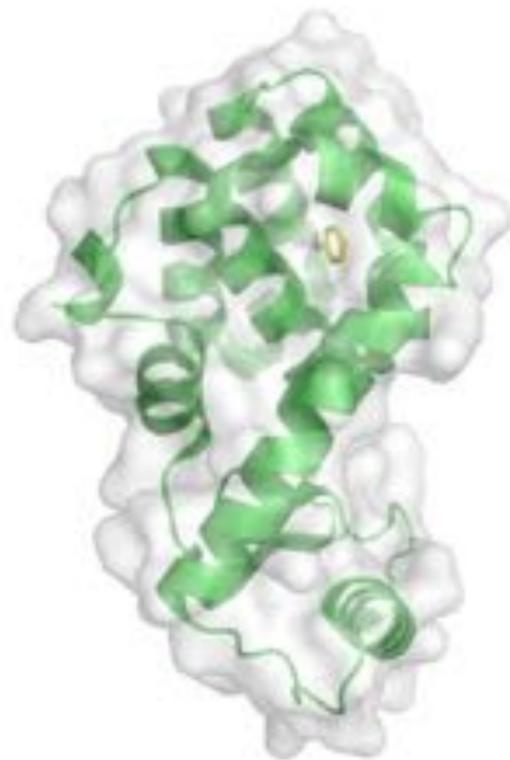
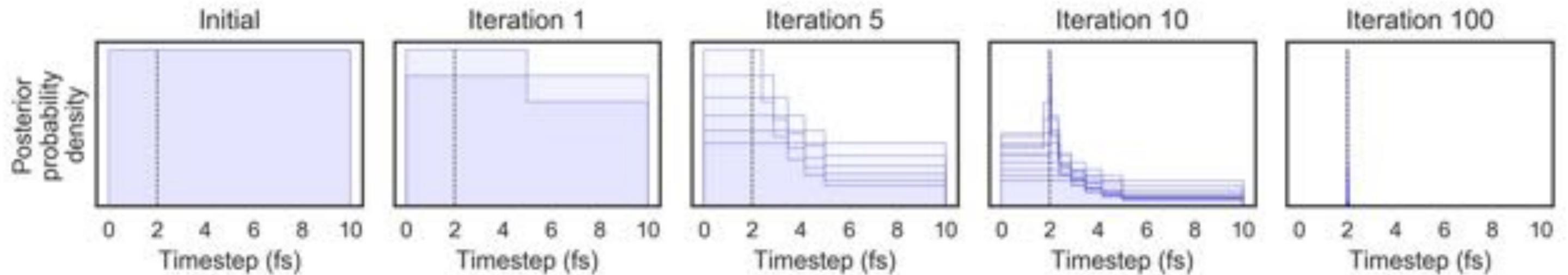
**[WARNING: PRELIMINARY DATA;  
NOT FULLY CONVERGED]**



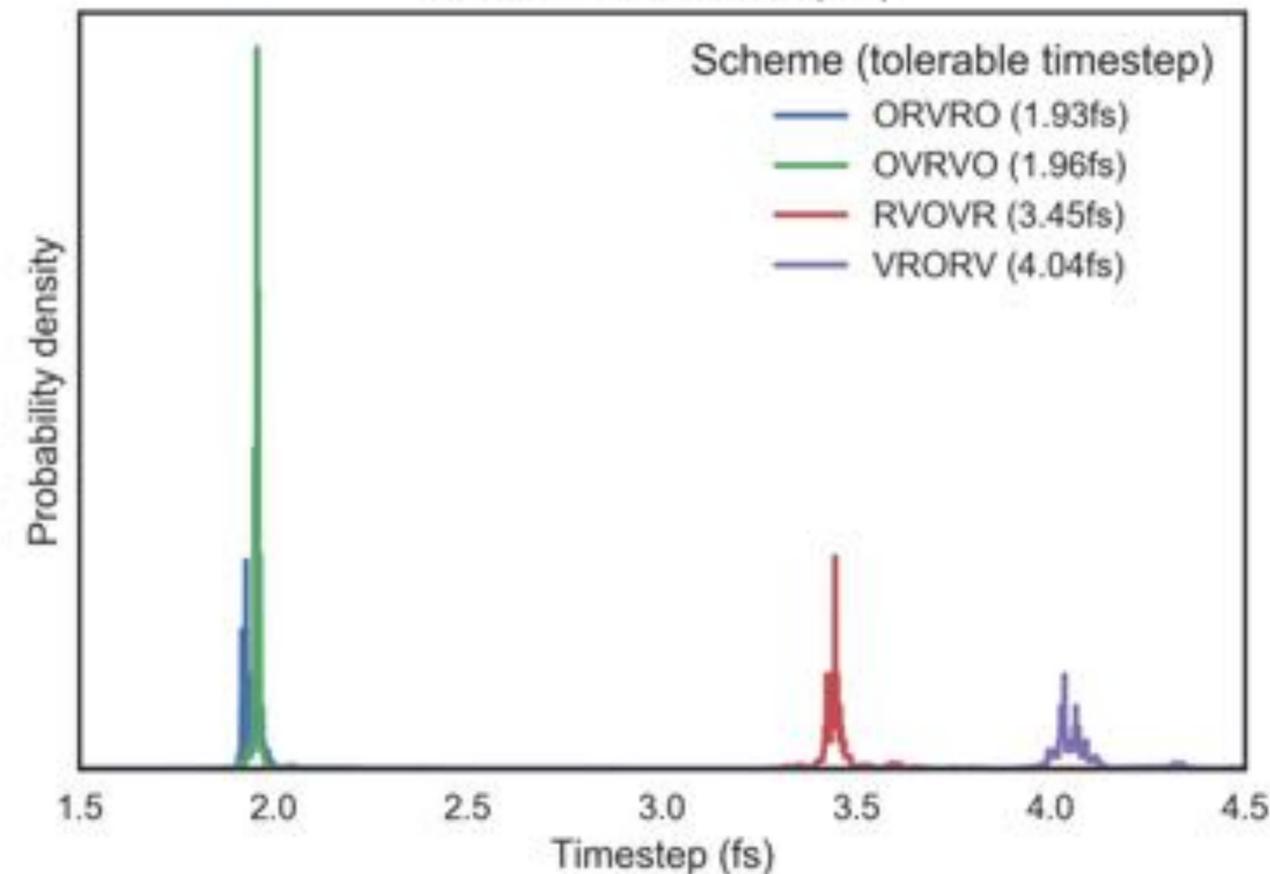
**JOSH FASS**



# PROBABILISTIC BISECTION ALLOWS FOR AUTOMATED SYSTEM-SPECIFIC DETERMINATION OF MAXIMUM TOLERATED TIMESTEP

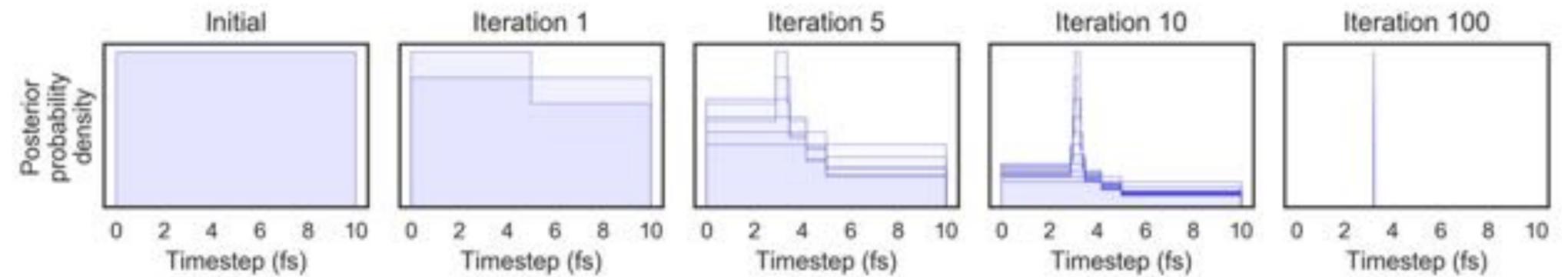
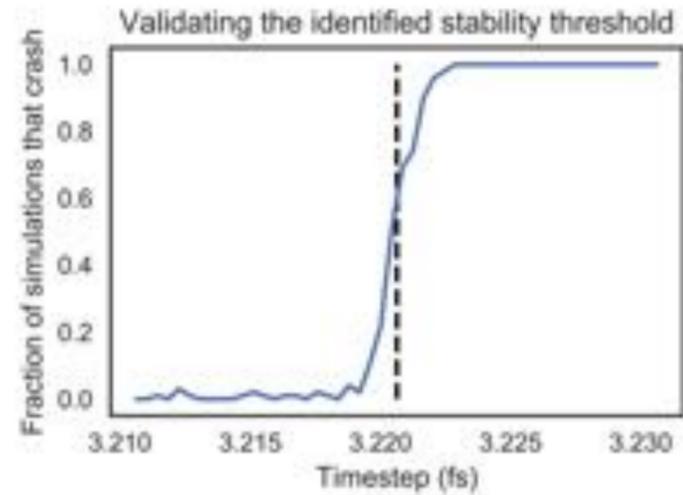


Lysozyme implicit max tolerable timestep  
Reference: OVRVO (2fs)

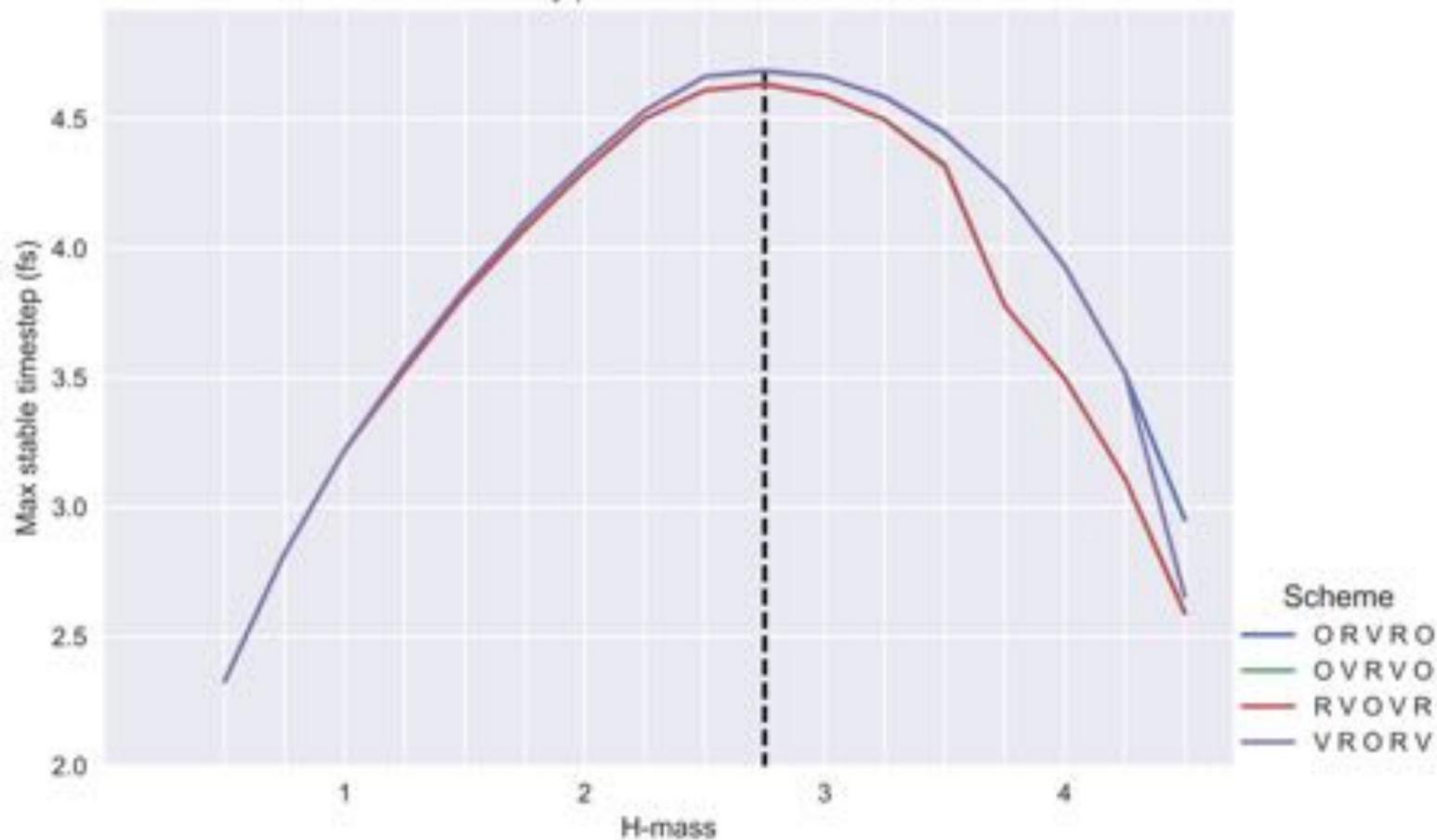


**JOSH FASS**

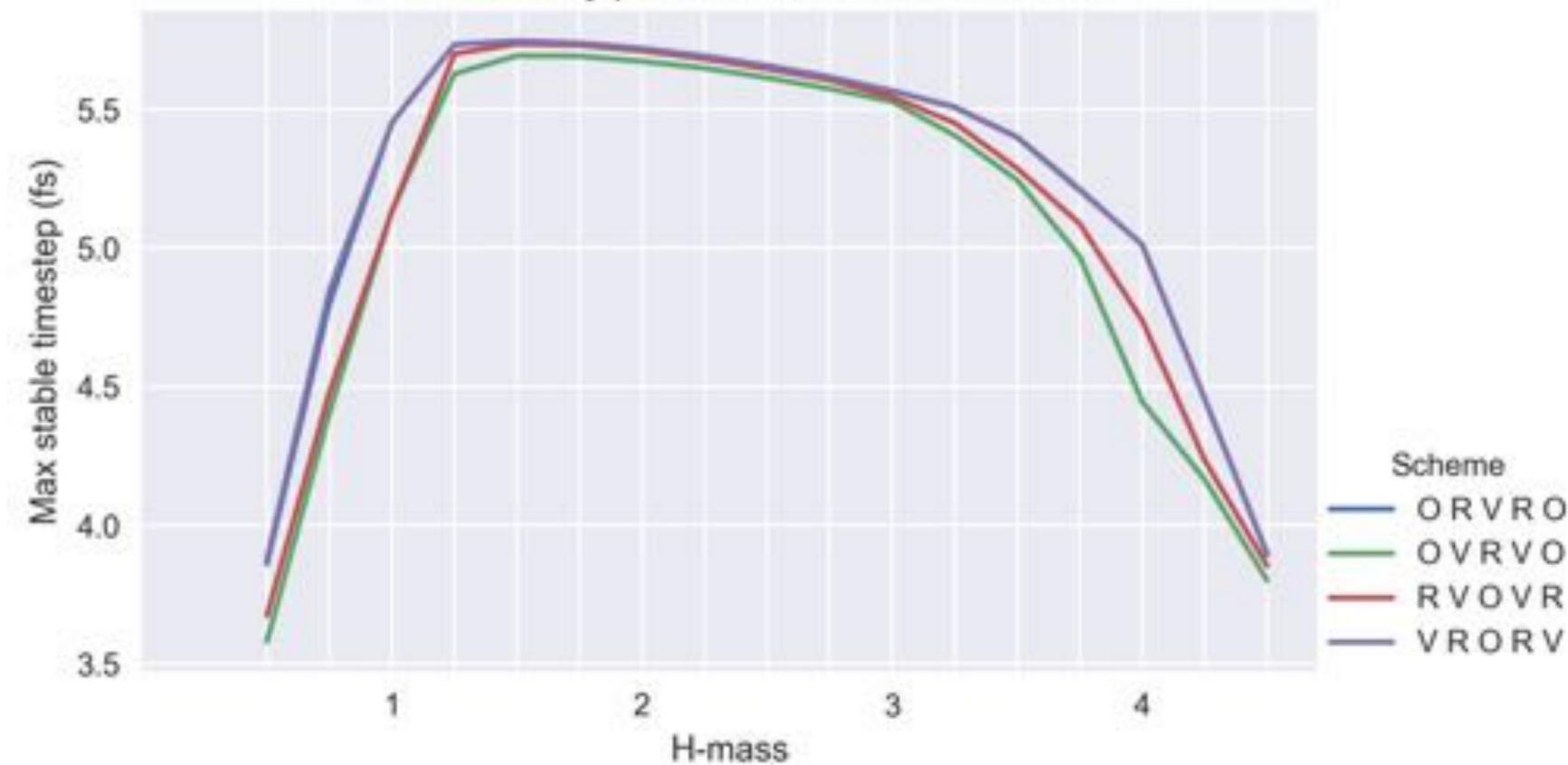
# PROBABILISTIC BISECTION CAN ALSO IDENTIFY MAXIMUM STABLE TIMESTEP



Stability thresholds on AlanineDipeptide (unconstrained) determined by probabilistic bisection search



Stability thresholds on LysozymeImplicit (constrained) determined by probabilistic bisection search



**HOW CAN WE IDENTIFY OPTIMAL  
THERMODYNAMIC PROTOCOLS?**

# MANY CHOICES EXIST FOR HOW TO DECOUPLE LIGAND FROM ENVIRONMENT

**$\lambda$ -DEPENDENT POTENTIAL MUST BE WELL-BEHAVED THROUGHOUT**

$$\Delta F = \int_0^1 d\lambda' \left\langle \frac{\partial V}{\partial \lambda} \right\rangle_{\lambda'}$$

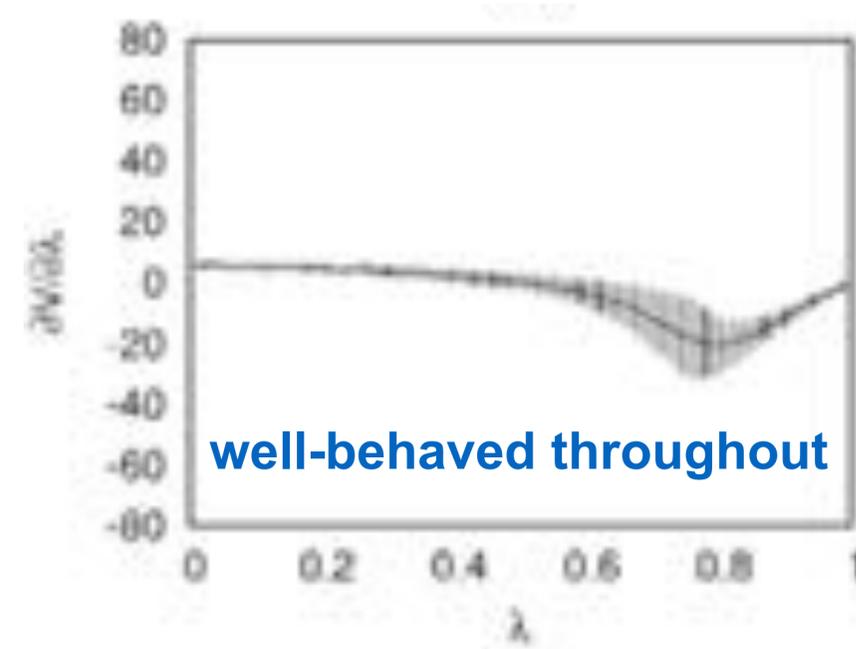
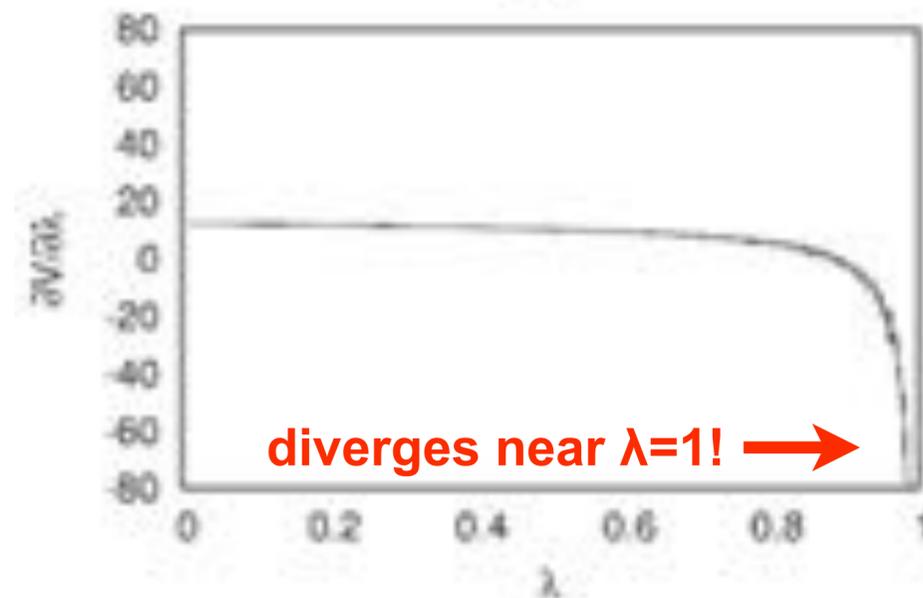
**THERMODYNAMIC INTEGRATION**

**LINEAR ALCHEMICAL SCALING**

$$U(r; \lambda) = 4\epsilon(1 - \lambda) \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right]$$

**“SOFT-CORE” FORM**

$$U(r; \lambda) = 4\epsilon(1 - \lambda) \left[ \frac{1}{[\alpha\lambda + (r/\sigma)^6]^2} - \frac{1}{[\alpha\lambda + (r/\sigma)^6]} \right]$$



# IS THERE AN OPTIMAL PROTOCOL?

THE **THERMODYNAMIC METRIC TENSOR** MEASURES HOW RAPIDLY THE EQUILIBRIUM DISTRIBUTION CHANGES AS CONTROL PARAMETERS ARE TWIDDLED.

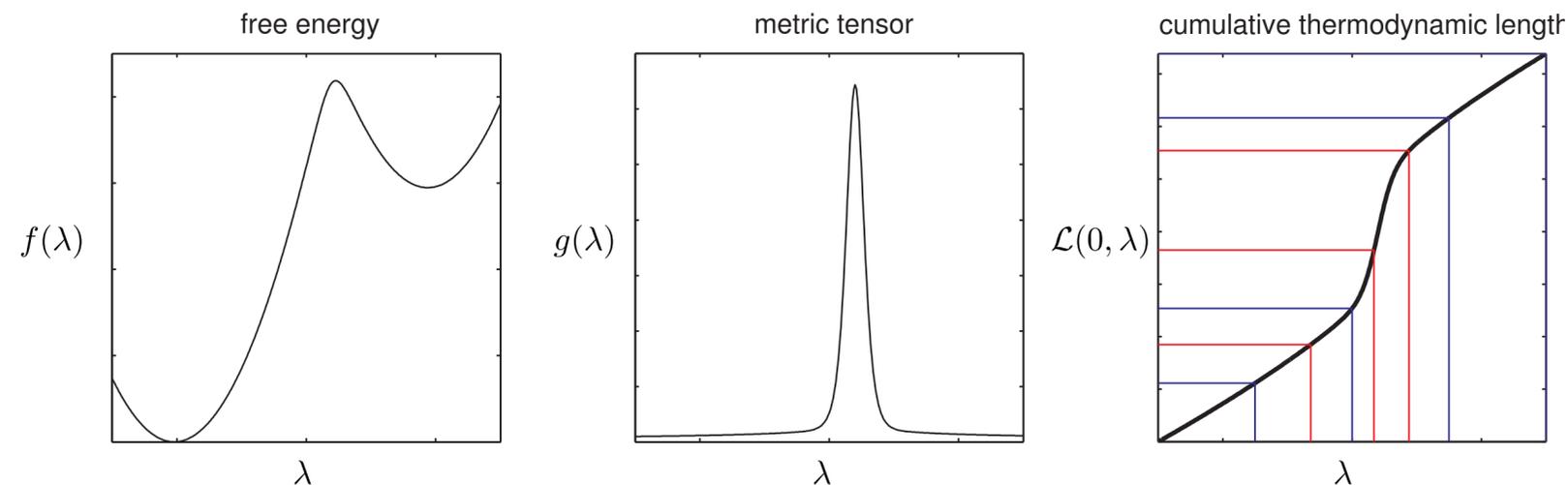
$$g_{ij}(\boldsymbol{\lambda}) \equiv \left\langle \frac{\partial \ln \pi(\boldsymbol{x}; \boldsymbol{\lambda})}{\partial \lambda_i} \frac{\partial \ln \pi(\boldsymbol{x}; \boldsymbol{\lambda})}{\partial \lambda_j} \right\rangle_{\boldsymbol{\lambda}}$$

THE **THERMODYNAMIC LENGTH** MEASURES HOW MUCH THE DISTRIBUTION HAS CHANGES FROM ONE VALUE OF CONTROL PARAMETERS TO ANOTHER.

$$\mathcal{L} \equiv \int_0^\tau dt \left( \dot{\boldsymbol{\lambda}}^T \boldsymbol{g} \dot{\boldsymbol{\lambda}} \right)^{1/2}$$

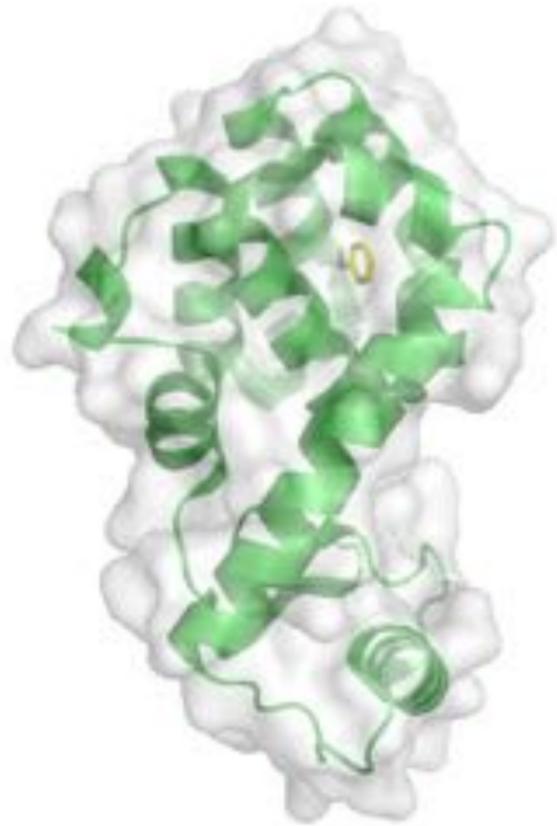
OPTIMAL PROTOCOLS ARE **GEODESICS IN THERMODYNAMIC METRIC SPACE**; THEY EQUALIZE THERMODYNAMIC LENGTH BETWEEN MEASUREMENTS.

$$\text{var}(\Delta \hat{f}) \geq N^{-1} \mathcal{L}(\lambda_a, \lambda_b)^2$$

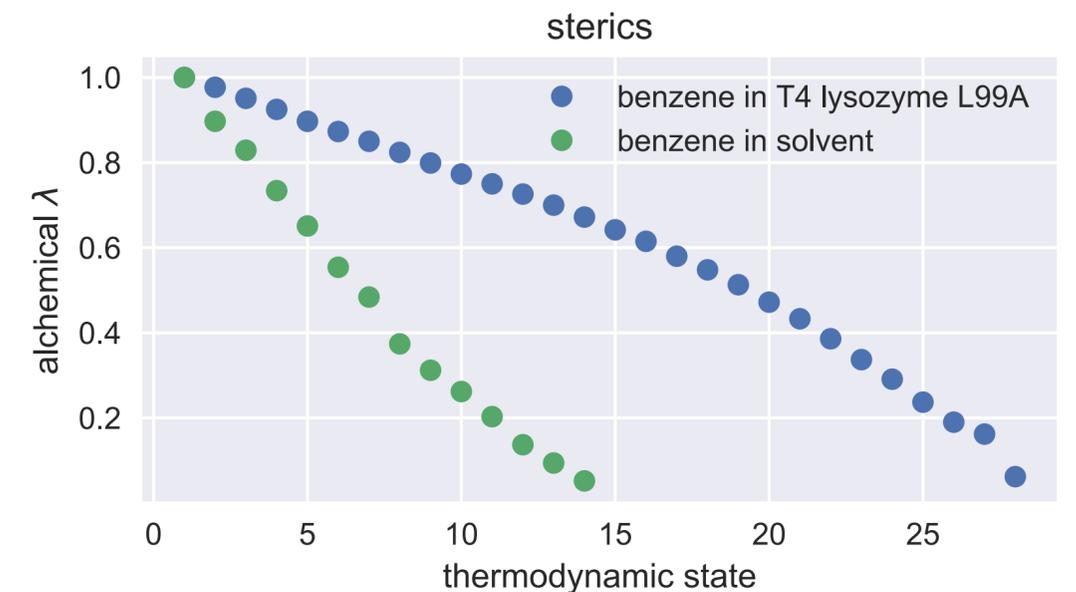
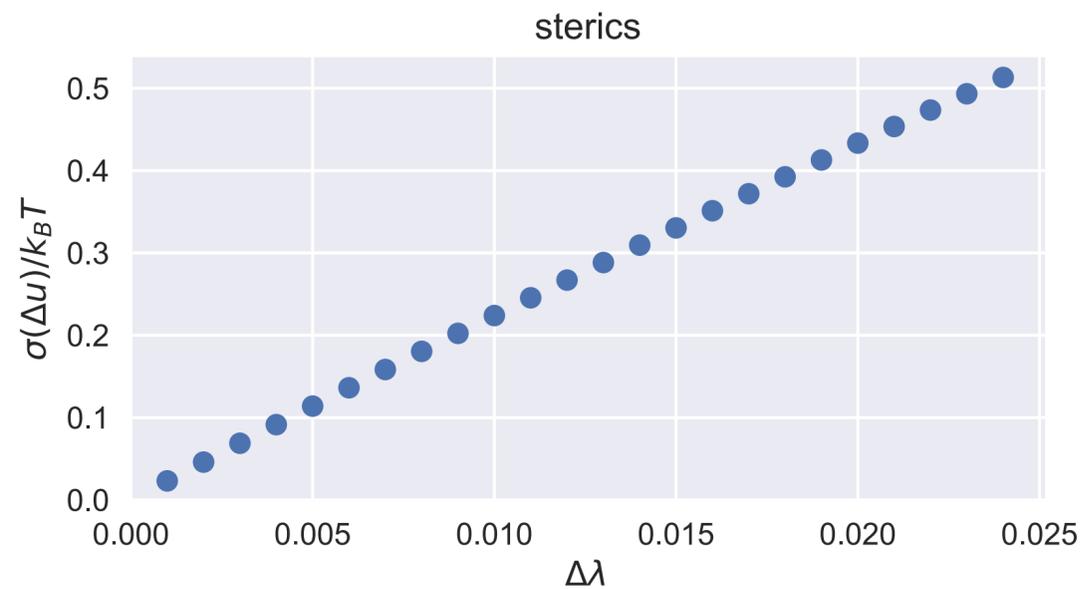
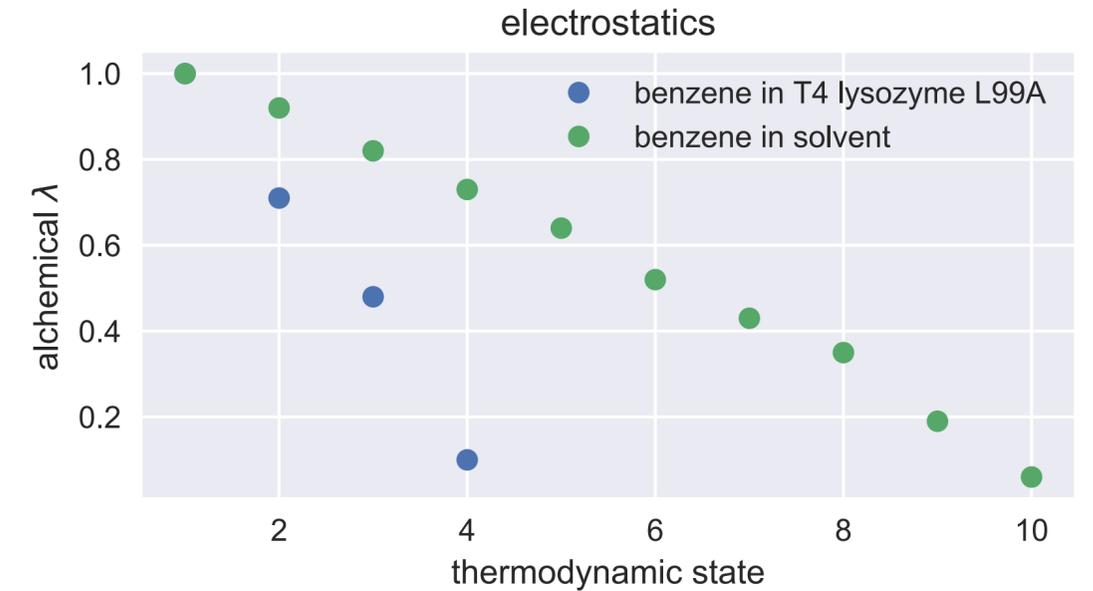
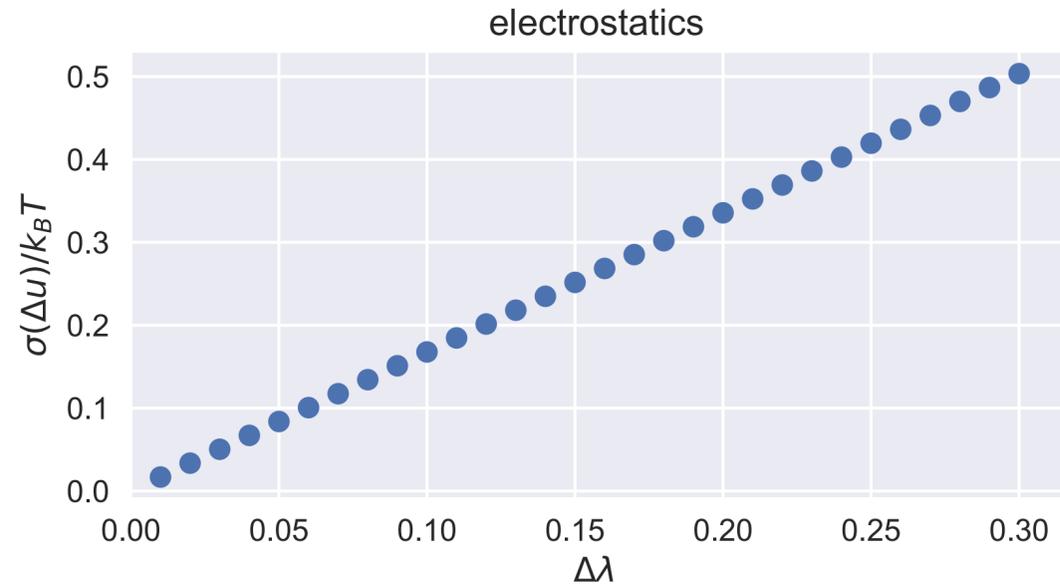


**OPTIMIZED PATHS CAN YIELD ORDERS OF MAGNITUDE REDUCTION IN VARIANCE!**

# A SIMPLE SCHEME FOR OPTIMIZING PROTOCOL GIVES A GOOD INITIAL GUESS



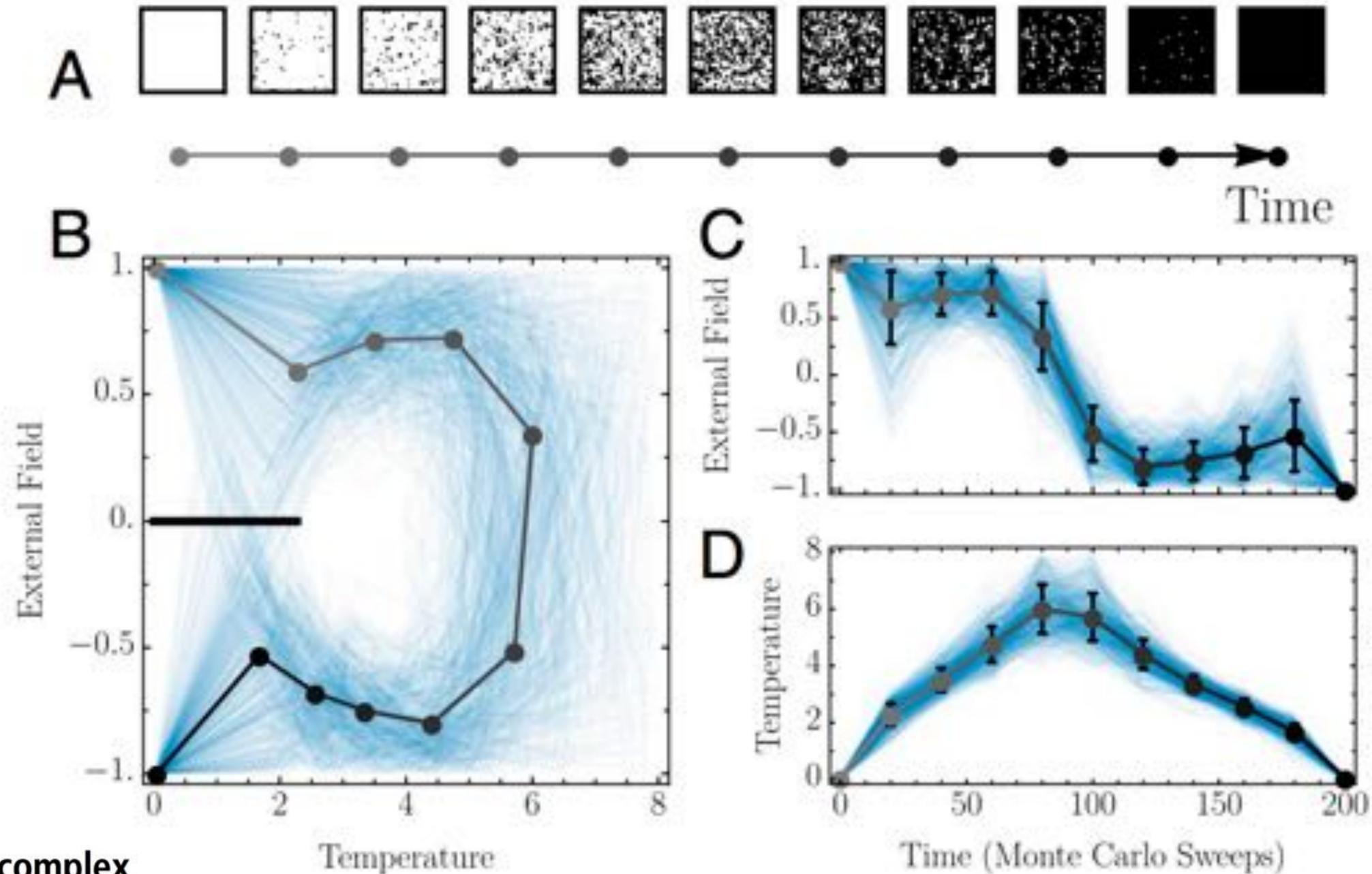
benzene : T4 lysozyme L99A  
GAFF / AMBER99SB-ILDN / TIP3P  
PME



**ANDREA RIZZI**  
CBM graduate student

(Suggested by David Minh and Huafeng Xu)

# ADAPTIVE PROTOCOL REFINEMENT WILL LIKELY BE FRUITFUL AREA OF RESEARCH



## Near-optimal protocols in complex nonequilibrium transformations

Todd R. Gingrich<sup>a,b,1</sup>, Grant M. Rotskoff<sup>c</sup>, Gavin E. Crooks<sup>d,e</sup>, and Phillip L. Geissler<sup>b,c,d,f</sup>

PNAS 113:10263, 2016.

# HOW CAN WE EXTRACT ALL INFORMATION FROM THE DATA?

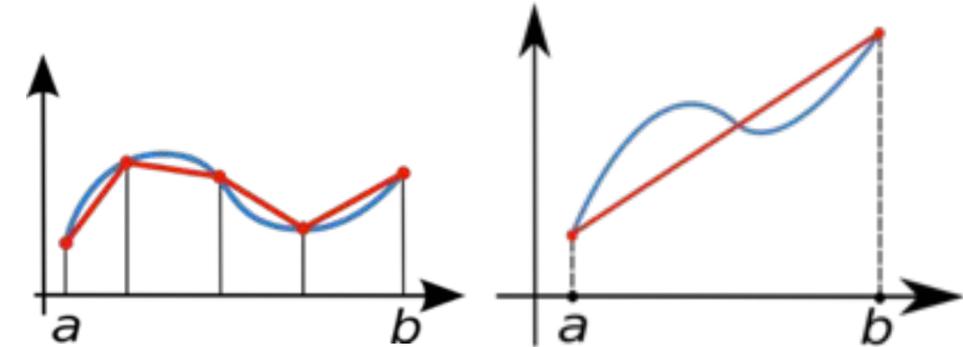
## NOT ALL ESTIMATORS ARE EQUAL

(IN BIAS AND STATISTICAL EFFICIENCY)

### THERMODYNAMIC INTEGRATION (TI)

$$\Delta F = \int_{\lambda_1}^{\lambda_2} d\lambda' \left\langle \frac{\partial H}{\partial \lambda} \right\rangle_{\lambda'} \approx \frac{\Delta \lambda}{2} \left[ \left\langle \frac{\partial H}{\partial \lambda} \right\rangle_{\lambda_1} + \left\langle \frac{\partial H}{\partial \lambda} \right\rangle_{\lambda_2} \right]$$

**QUADRATURE ERROR (BIAS) DIFFICULT TO QUANTIFY;  
DERIVATIVES OFTEN NOT AVAILABLE IN SIMULATION CODES**



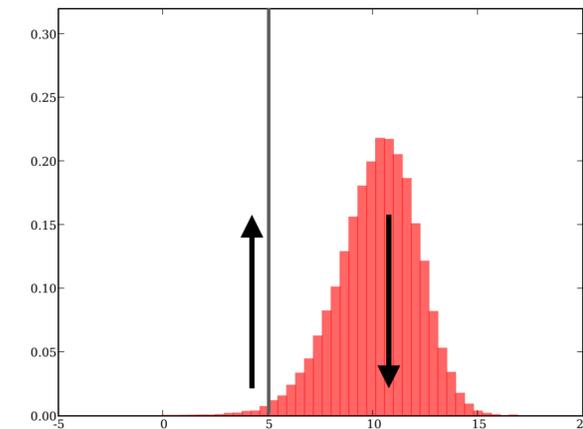
### EXPONENTIAL REWEIGHTING (EXP)

$$\Delta F = -\beta^{-1} \ln \left\langle e^{-\beta(U_2 - U_1)} \right\rangle_{\lambda_1} = +\beta^{-1} \ln \left\langle e^{-\beta(U_1 - U_2)} \right\rangle_{\lambda_2}$$

Zwanzig RW. *JCP* **22**:1420, 1954.

Shirts MR and Pande VS. *JCP* **122**:144107, 2005.

**SUFFERS FROM LARGE BIAS AND VARIANCE**



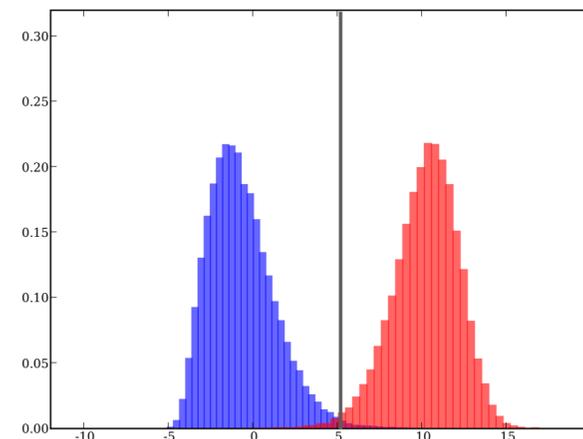
### BENNETT ACCEPTANCE RATIO (BAR)

$$\Delta F = -\beta^{-1} \ln \frac{\langle f(U_2 - U_1) \rangle_{\lambda_1}}{\langle f(U_1 - U_2) \exp[-\beta(U_2 - U_1)] \rangle_{\lambda_2}}$$

Bennett CH. *J Comput Phys* **22**:245, 1976.

Shirts MR, Bair E, Hooker G, and Pande VS. *PRL* **91**:140601, 2003.

**SUPERIOR, BUT ONLY APPLICABLE TO TWO STATES—WHY CAN'T WE USE ALL THE DATA?**



# THE **REDUCED POTENTIAL** GENERALIZES ALL THERMODYNAMIC STATES OF INTEREST

Define the **reduced potential** for a state  $k$  as a combination of terms

$$u_k(\mathbf{x}) = \beta_k [U_k(\mathbf{x}) + p_k V(\mathbf{x}) + \mu_k^T \mathbf{N}(\mathbf{x})]$$

with thermodynamic parameters for each state

$\beta_k$	inverse temperature
$U_k$	potential energy function
$p_k$	external pressure
$\mu_k$	chemical potential of exchangeable species

where

$\mathbf{x}$	microstate or configuration	
$V(\mathbf{x})$	volume of simulation box	
$\mathbf{N}(\mathbf{x})$	number of each chemical species in system	protonation states number of counter ions

The target configuration space density is given by

$$\pi_k(\mathbf{x}) = Z_k^{-1} \exp[-u_k(\mathbf{x})] \quad Z_k = \int d\mathbf{x} \exp[-u_k(\mathbf{x})]$$

Covers many common thermodynamic ensembles: NVT, NPT,  $\mu$ VT,  $\mu$ PT

# EXTENDED BRIDGE SAMPLING ESTIMATORS

## DEFINE A FAMILY OF USEFUL ESTIMATORS

Suppose we have  $K$  states defined by unnormalized distribution functions  $q(\mathbf{x})$ :

$$\pi_k(\mathbf{x}) = c_k^{-1} q_k(\mathbf{x}) \quad q_k(\mathbf{x}) > 0$$
$$c_k = \int d\mathbf{x} q_k(\mathbf{x})$$

Start with the identity

$$c_i \langle \alpha_{ij} q_j \rangle_i = \int_{\Gamma} d\mathbf{x} q_i(\mathbf{x}) \alpha_{ij}(\mathbf{x}) q_j(\mathbf{x}) = c_j \langle \alpha_{ij} q_i \rangle_j \quad i, j = 1, \dots, K$$

$\alpha_{ij}(\mathbf{x})$  arbitrary

Sum over all  $K$  states

$$\sum_{j=1}^K c_i \langle \alpha_{ij} q_j \rangle_i = \sum_{j=1}^K c_j \langle \alpha_{ij} q_i \rangle_j \quad i = 1, \dots, K$$

Substitute the empirical estimator for expectations

$$\sum_{j=1}^K \frac{\hat{c}_i}{N_i} \sum_{n=1}^{N_i} \alpha_{ij} q_j(\mathbf{x}_{in}) = \sum_{j=1}^K \frac{\hat{c}_j}{N_j} \sum_{n=1}^{N_j} \alpha_{ij} q_i(\mathbf{x}_{jn})$$

Defines a **family** of estimators parameterized by choice of  $\alpha_{ij}(\mathbf{x})$   
Ratios of normalization constants  $c_i/c_j$  determined by solving a set of  $K$  nonlinear equations.

# THE ASYMPTOTICALLY LOWEST-VARIANCE ESTIMATOR IS KNOWN!

Optimal choice of  $\alpha_{ij}(\mathbf{x})$  among large class of functions is:

$$\alpha_{ij}(\mathbf{x}) = \frac{N_j \hat{c}_j^{-1}}{\sum_{k=1}^K N_k \hat{c}_k^{-1} q_k(\mathbf{x})}$$

Covariance of any functions  $\phi, \psi$  of ratios  $\hat{c}_i/\hat{c}_j$  can be estimated from asymptotic covariance

$$\text{cov}(\hat{\phi}, \hat{\psi}) = \sum_{i,j=1}^K \frac{\partial \phi}{\partial \theta_i} \hat{\Theta}_{ij} \frac{\partial \psi}{\partial \theta_j}$$

where the **asymptotic covariance estimate** is computed in terms of  $\hat{\theta}_i \equiv -\ln \hat{c}_i$

$$\hat{\Theta} = [(\mathbf{W}^T \mathbf{W})^+ - \mathbf{N}]^+$$

$$\hat{\Theta}_{ij} \equiv \text{cov}(\hat{\theta}_i, \hat{\theta}_j)$$

$$W_{ni} = \frac{\hat{c}_i^{-1} q_i(x_n)}{\sum_{k=1}^K N_k \hat{c}_k^{-1} q_k(x_n)}$$

$$\mathbf{N} = \text{diag}\{N_1, N_2, \dots, N_K\}$$

code available at:  
<http://pymbar.org>

## MULTISTATE GENERALIZATION OF BENNETT ACCEPTANCE RATIO (MBAR)

**HOW CAN WE RUN CALCULATIONS QUICKLY  
AND RAPIDLY TEST NEW ALGORITHMS?**

# YANK: AN OPEN-SOURCE, COMMUNITY-ORIENTED PLATFORM FOR GPU-ACCELERATED FREE ENERGY CALCULATIONS



**NVIDIA GTX-1080 (\$650)**  
**9 TFLOP/S SINGLE PRECISION**

OpenMM speedup (GTX-1080) over 12-core Xeon X5650 CPU for DHFR

method	natoms	gromacs CPU	OpenMM GPU	speedup
GB/SA	2,489	2.54 ns/day	789 ns/day	<b>311 x</b>
RF	23,558	18.8 ns/day	572 ns/day	<b>30.4 x</b>
PME	23,558	6.96 ns/day	337 ns/day	<b>48.4 x</b>

<http://openmm.org> OpenMM 7.1.0 development snapshot benchmark  
gromacs benchmarks from <http://biowulf.nih.gov/apps/gromacs-gpu.html>



A free, open-source, extensible platform for free energy calculations and ligand design

**ANDREA RIZZI**  
CBM graduate student



**LEVI NADEN**  
Merck KGaA-sponsored postdoc

Docs » YANK

YANK 1.0 coming Aug 2017!

## YANK

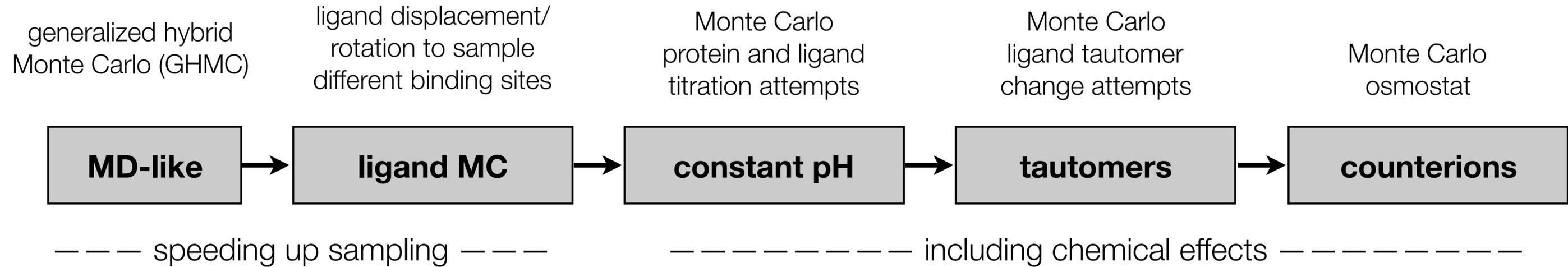
A GPU-accelerated Python framework for exploring algorithms for alchemical free energy calculations

### Features

- Modular Python framework for easily exploring new algorithms
- GPU-accelerated via the [OpenMM toolkit](#)
- [Alchemical free energy calculations](#) in both explicit and implicit solvent
- Hamiltonian exchange among alchemical intermediates with Gibbs sampling framework
- [General Markov chain Monte Carlo](#) framework for exploring enhanced sampling methods
- Built-in equilibration detection and convergence diagnostics
- Support for AMBER prmtop/inpcrd files
- Support for absolute binding free energy calculations
- Support for transfer free energies (such as hydration free energies)

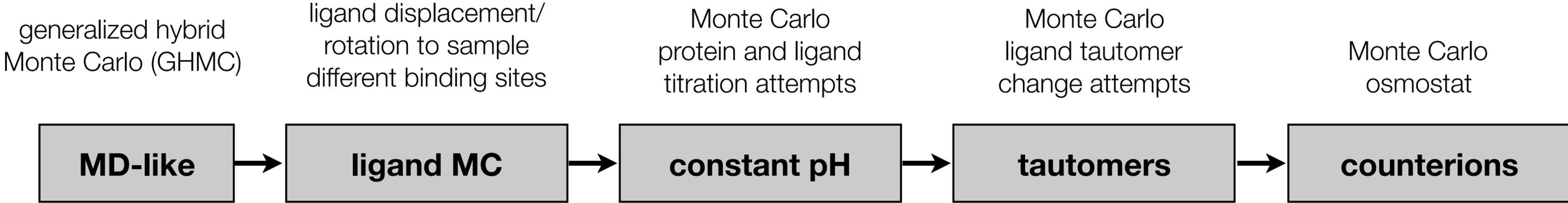
<http://www.getyank.org>

# MARKOV CHAIN MONTE CARLO (MCMC) FRAMEWORK PERMITS EFFICIENT MIXING\* OF MD AND MC



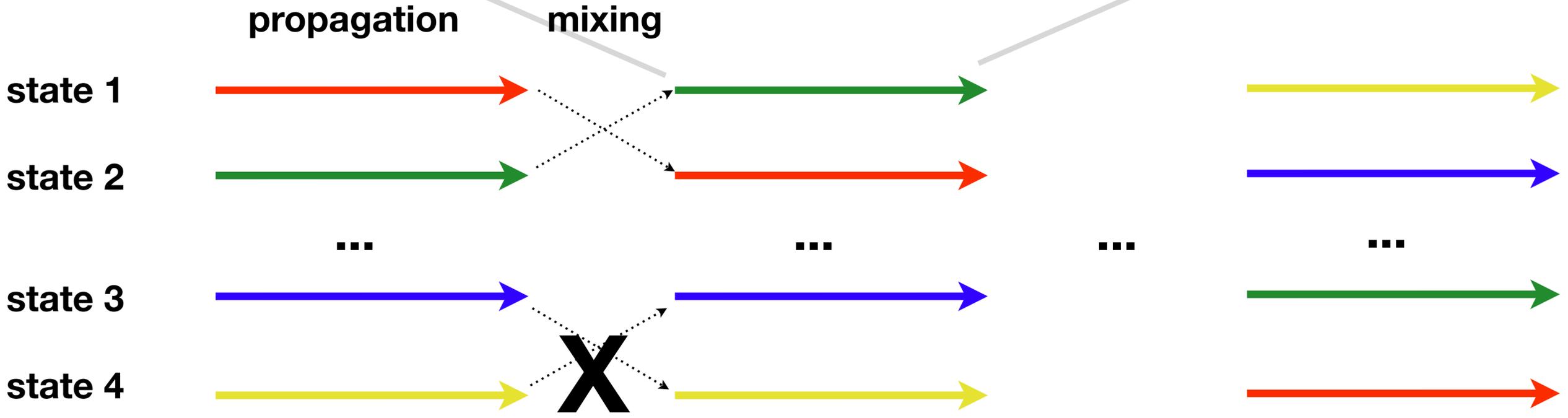
\* Must be careful to avoid mismatch between MD/MC densities.

# MCMC IS COMBINED WITH REPLICA EXCHANGE TO DECREASE CORRELATION TIMES



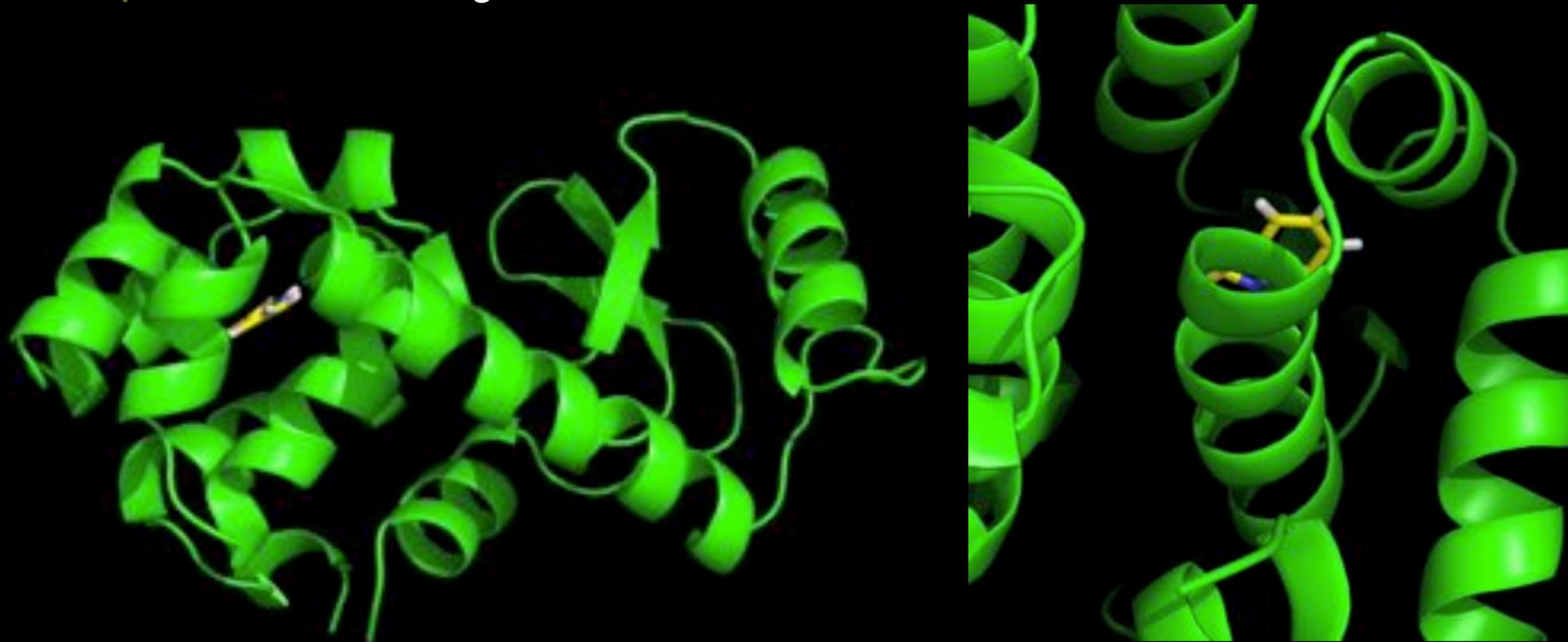
--- speeding up sampling ---      ----- including chemical effects -----

\* Must be careful to avoid mismatch between MD/MC densities.



# HAMILTONIAN EXCHANGE PROTOCOL ALLOWS FOR REPEATED BINDING/UNBINDING EVENTS AND REORIENTATION IN SITE

**solid** fully interacting  
**transparent** noninteracting



**indole** binding to T4 lysozyme L99A  
12 h on 2 NVIDIA Tesla M2090 GPUs  
Hamiltonian exchange with Gibbs sampling

# YAML SYNTAX ALLOWS FOR EASY AUTOMATION OF FREE ENERGY CALCULATIONS

```
options:
  verbose: true
  output_dir: test_kinase
  temperature: 310*kelvin
  pressure: 1*atmosphere
  constraints: HBonds
  number_of_iterations: 200
  minimize: yes

molecules:
  abl:
    filepath: !Combinatorial [2HY7.pdb, 3CS9.pdb]
    strip_protons: yes
  src:
    filepath: !Combinatorial [2010_A.pdb, 3EL7_A.pdb]
    strip_protons: yes

inhibitors:
  filepath: clinical-kinase-inhibitors.csv
  select: 0
  antechamber:
    charge_method: bcc
  epik:
    select: 0
    tautomerize: no
    ph: 7.4
    ph_tolerance: 5

# Run Schrodinger's tool Epik with default parameters and
# select the most likely protonation state for the molecule
# (in solution). More control over epik parameters is
# possible (see YML documentation).
```

# YAML SYNTAX ALLOWS FOR EASY AUTOMATION OF FREE ENERGY CALCULATIONS

```
solvents:  
  PME:  
    nonbonded_method: PME  
    nonbonded_cutoff: 1+nanometer  
    switch_distance: 0.9+nanometer  
    ewald_error_tolerance: 0.0003  
    clearance: 10+angstroms  
    positive_ion: Na+  
    negative_ion: Cl-  
  GBSA:  
    nonbonded_method: NoCutoff # Implicit solvent using GBSA/OBC2 model  
    implicit_solvent: OBC2  
    implicit_solvent_salt_concentration: 1.0+mole/liter  
    solute_dielectric: 1.5  
    solvent_dielectric: 80.0  
  VACUUM:  
    nonbonded_method: NoCutoff
```

# YAML SYNTAX ALLOWS FOR EASY AUTOMATION OF FREE ENERGY CALCULATIONS

```
systems:  
  kinase-inhibitor:  
    receptor: !Combinatorial [abl, src]  
    ligand: !Combinatorial [imatinib, bosutinib]  
    solvent: !Combinatorial [PME, GBSA]  
    pack: yes  
    leap:  
      parameters: [oldff/leaprc.ff99SBildn, leaprc.gaff]
```

```
  imatinib-hydration:  
    solute: imatinib  
    solvent1: PME  
    solvent2: vacuum  
    leap:  
      parameters: [leaprc.ff14SB, leaprc.gaff]
```

```
amber-system:  
  phase1_path: [complex.prmtop, complex.inpcrd]  
  phase2_path: [solvent.prmtop, solvent.inpcrd]  
  ligand_dsl: resname MOL  
  solvent: RF
```

```
# System files for the complex phase.  
# System files for the solvent phase.  
# MDTraj DSL string to select the ligand.  
# Specify how to model the solvent.
```

```
gromacs-system:  
  phase1_path: [complex.top, complex.gro]  
  phase2_path: [solvent.top, solvent.gro]  
  gromacs_include_dir: include/
```

```
# Optional path to the directory containing the files  
# included in .top files.
```

```
  ligand_dsl: resname MOL  
  solvent: RF
```

# YAML SYNTAX ALLOWS FOR EASY AUTOMATION OF FREE ENERGY CALCULATIONS

```
protocols:
  absolute-binding:
    complex:
      alchemical_path:
        lambda_electrostatics: [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]
        lambda_sterics: [1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.0]
    solvent:
      alchemical_path:
        lambda_electrostatics: [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]
        lambda_sterics: [1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1, 0.0]

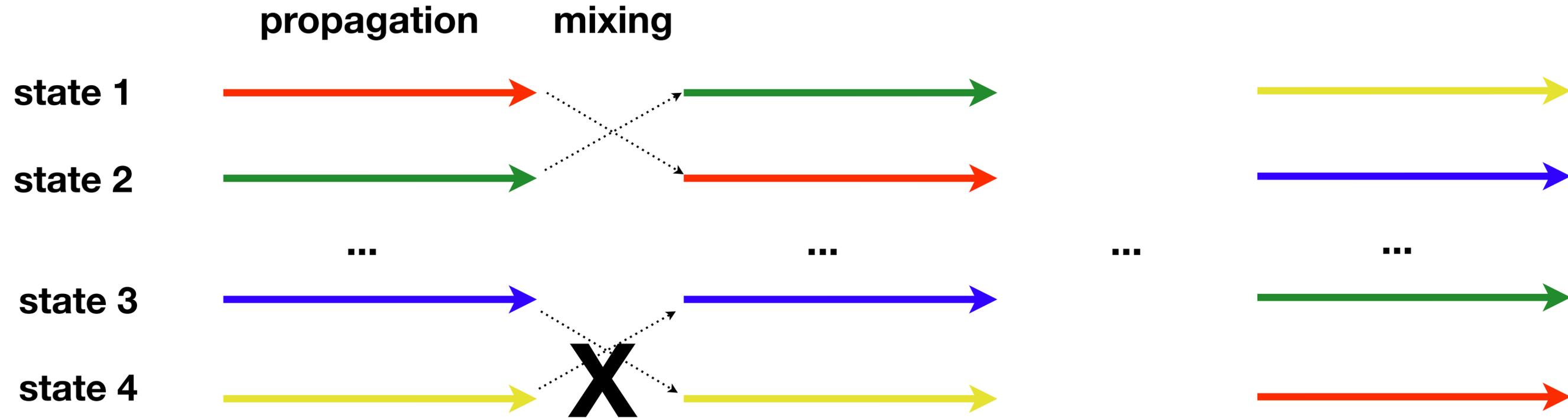
experiments:
  system: kinase-inhibitor
  protocol: absolute-binding
  restraint:
    type: !Combinatorial [Harmonic, FlatBottom] # Optimally, apply a restraint to the ligand to keep
    # it close to the receptor. Possible types are null,
    # Harmonic, FlatBottom, and Barostat.
  options:
    temperature: !Combinatorial [298.0*kelvin, 307.0*kelvin] # All options that can be specified in the "options"
    # section can be included here.
```

# INSTALLING YANK IS EASY

```
MINICONDA="Miniconda3-latest-Linux-x86_64.sh"  
wget https://repo.continuum.io/miniconda/$MINICONDA  
bash $MINICONDA -b -p $HOME/miniconda  
export PATH="$HOME/miniconda/bin:$PATH"
```

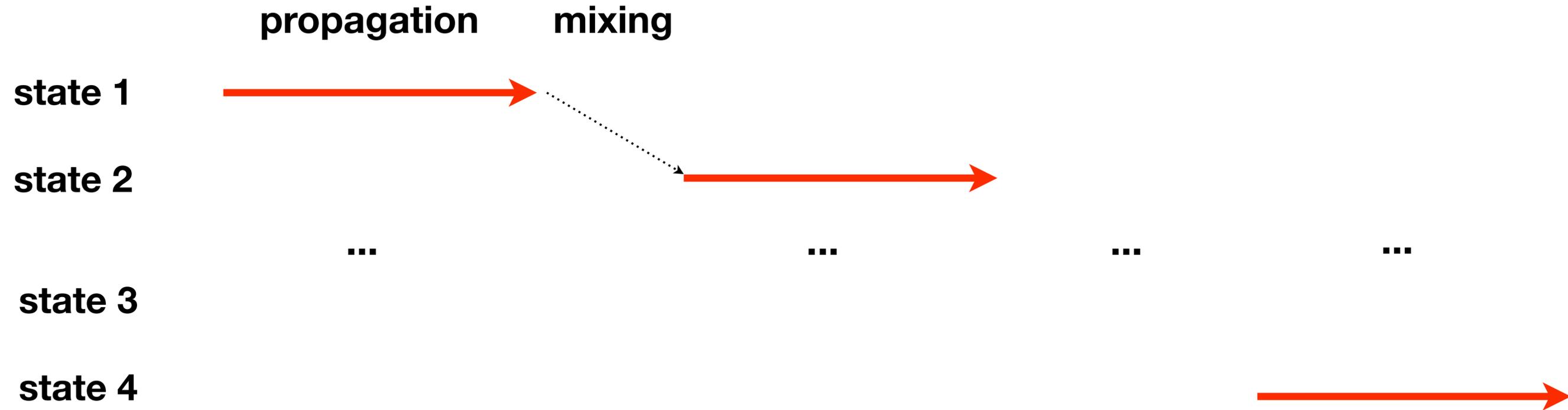
```
conda install -c omnia yank
```

# DO WE NEED ALL THOSE REPLICAS?



## HAMILTONIAN EXCHANGE

# DO WE NEED ALL THOSE REPLICAS?



## EXPANDED ENSEMBLE

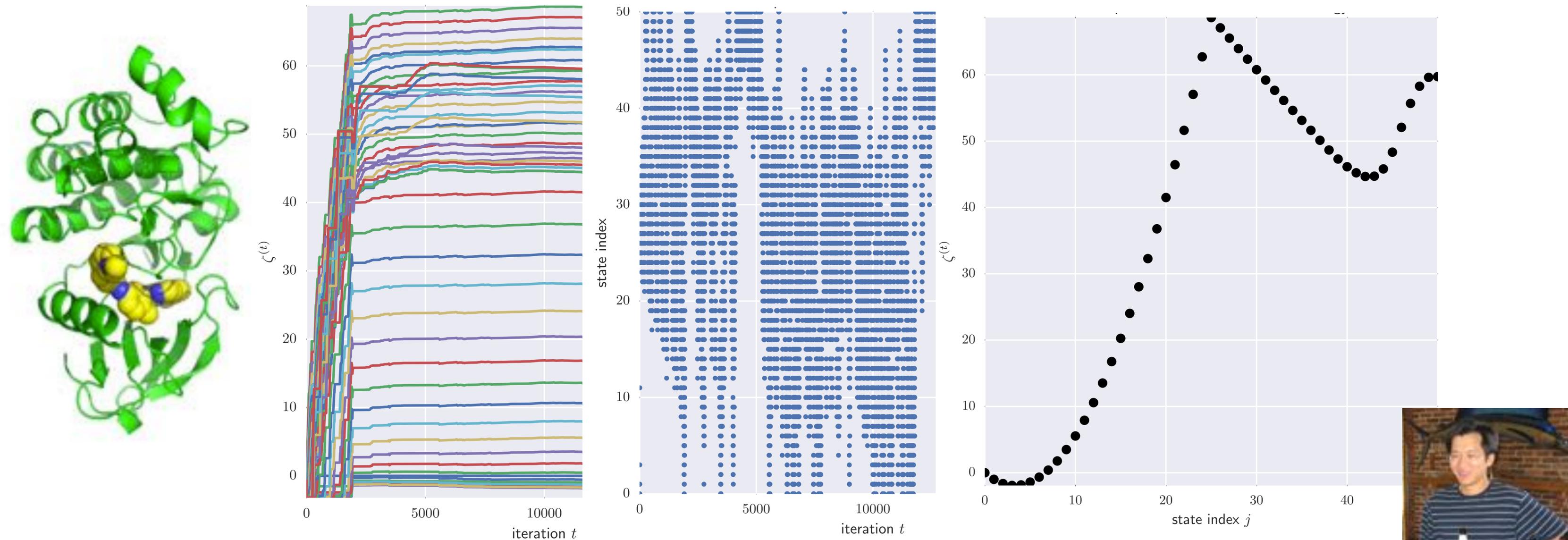
$$\pi(x, k) \propto \exp[-u_k(x) + g_k]$$

One caveat: We need to guess the weights  $g$   
(which are unfortunately the free energies we are trying to compute!)

# SELF-ADJUSTED MIXTURE SAMPLING (SAMS)

Provably asymptotically optimal algorithm for estimating free energies from a single simulation!

Tan Z. J. Comp. Graph. Stat. <http://dx.doi.org/10.1080/10618600.2015.1113975>



Similar to simulated scaling (Wei Yang), but provably optimal.

**ZHIQIANG TAN**  
Rutgers

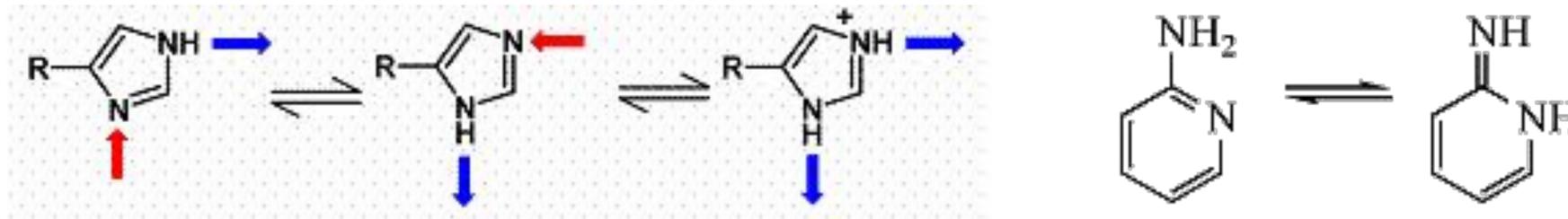


# PREDICTIONS FAIL FOR THREE REASONS

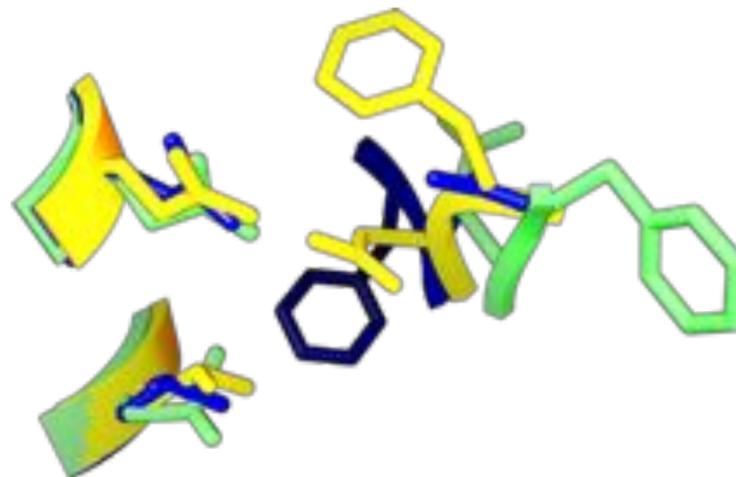
1. The **forcefield** does a poor job of modeling the physics of our system

$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r (r - r_{eq})^2 + \sum_{\text{angles}} K_\theta (\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$

2. We're missing some **essential chemical** in our simulations (e.g. protonation states, tautomers, covalent association)



3. We haven't **sampled** all of the relevant conformations



# PREDICTIONS FAIL FOR THREE REASONS

1. The **forcefield** does a poor job of modeling the physics of our system

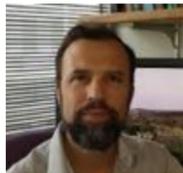
$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r (r - r_{eq})^2 + \sum_{\text{angles}} K_\theta (\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$

# THE OPEN FORCEFIELD GROUP

<https://github.com/open-forcefield-group>



**DAVID MOBLEY**  
**UCI**



**MICHAEL GILSON**  
**UCSD**



**MICHAEL SHIRTS**  
**UNIVERSITY OF COLORADO, BOULDER**



**CHRISTOPHER BAYLY**  
**OPENEYE SCIENTIFIC**



**JOHN CHODERA**  
**SKI/MSKCC**

# DIRECT CHEMICAL PERCEPTION ESCAPES THE SHACKLES OF ATOM TYPE CENTRIC FORCEFIELDS

```
<NonbondedForce coulomb14scale="0.833333" lj14scale="0.5" sigma_unit="angstroms" epsilon_unit="kilocalories_per_mole">  
  <Atom smirks="#1:1" rmin_half="1.4870" epsilon="0.0157"/>  
  <Atom smirks="#1:1-#6" rmin_half="1.4870" epsilon="0.0157"/>  
  ...  
</NonbondedForce>
```

```
<HarmonicBondForce length_unit="angstroms" k_unit="kilocalories_per_mole/angstrom**2">  
  <Bond smirks="#6X4:1-#6X4:2" length="1.526" k="620.0"/>  
  <Bond smirks="#6X4:1-#1:2" length="1.090" k="680.0"/>  
  ...  
</HarmonicBondForce>
```

```
<HarmonicAngleForce angle_unit="degrees" k_unit="kilocalories_per_mole/radian**2">  
  <Angle smirks="[a,A:1]-#6X4:2-[a,A:3]" angle="109.50" k="100.0"/>  
  <Angle smirks="#1:1-#6X4:2-#1:3" angle="109.50" k="70.0"/>  
</HarmonicAngleForce>
```

```
<BondChargeCorrections method="AM1" increment_unit="elementary_charge">  
  <BondChargeCorrection smirks="#6X4:1-#6X3a:2" increment="+0.0073"/>  
  <BondChargeCorrection smirks="#6X4:1-#6X3a:2-#7" increment="-0.0943"/>  
  <BondChargeCorrection smirks="#6X4:1-#8:2" increment="+0.0718"/>  
</BondChargeCorrections>
```

**SMARTS IS LIKE REGULAR EXPRESSIONS FOR SMILES**  
**CHECK OUT OUR 332-LINE SMALL MOLECULE FORCEFIELD!**

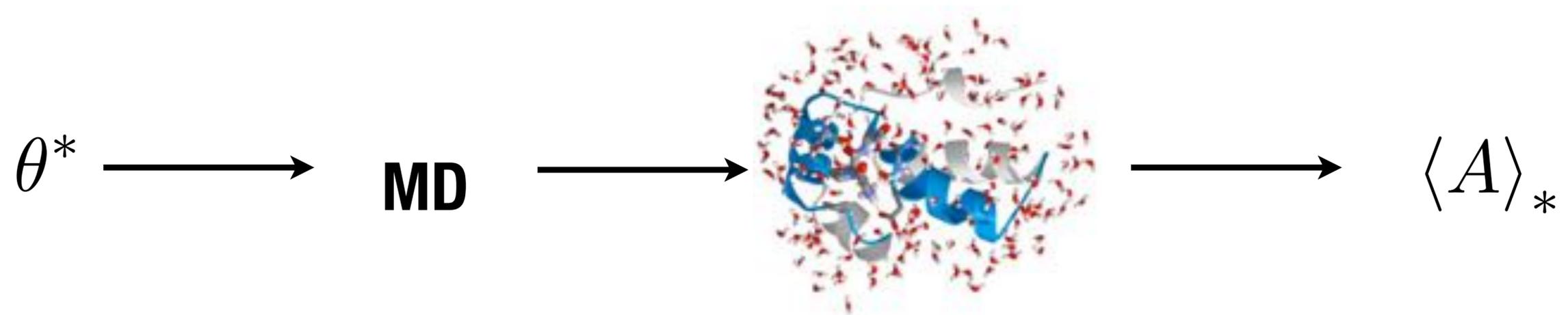
**implementation:** <http://github.com/open-forcefield-group/openforcefield>

**forcefield:** <http://github.com/open-forcefield-group/smirnoff99Frosst>

# EXAMPLE: TIP3P WATER

```
<?xml version='1.0' encoding='ASCII'?>
<SMIRNOFF version="0.1">
  <!-- SMIRks Native Open Force Field (SMIRNOFF) file -->
  <Date>2017-04-29</Date>
  <Author>J. D. Chodera, MSKCC; A. Rizzi, Weill Cornell; C. C. Bannan, UC Irvine</Author>
  <!-- SMIRNOFF file implementing TIP3P water model. -->
  <NonbondedForce coulomb14scale="0.833333" lj14scale="0.5" sigma_unit="nanometers" epsilon_unit="kilojoules_per_mole">
    <!-- TIP3P water oxygen with charge override -->
    <Atom smirks="#1-[#8X2H2+0:1]-[#1]" id="n1" sigma="0.31507524065751241" epsilon="0.635968" charge="-0.834"/>
    <!-- TIP3P water hydrogen with charge override -->
    <Atom smirks="#1:1-[#8X2H2+0]-[#1]" id="n2" sigma="1" epsilon="0" charge="0.417"/>
  </NonbondedForce>
  <Constraints distance_unit="angstroms">
    <!-- constrain water O-H bond to equilibrium bond length (overrides earlier constraint) -->
    <Constraint smirks="#1:1-[#8X2H2+0:2]-[#1]" id="c1" distance="0.9572"/>
    <!-- constrain water H...H, calculating equilibrium length from H-O-H equilibrium angle and H-O equilibrium bond lengths -->
    <Constraint smirks="#1:1-[#8X2H2+0]-[#1:2]" id="c2" distance="1.5139806545247014"/>
  </Constraints>
</SMIRNOFF>
```

# THE OLD WAY



One set of parameters in, one computed result out

# THE BAYESIAN WAY

Bayes rule provides a **probability measure** over unknown parameters given data and an automated way to **update** parameters given new experimental data

$$p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$$

$\mathcal{D}$  data

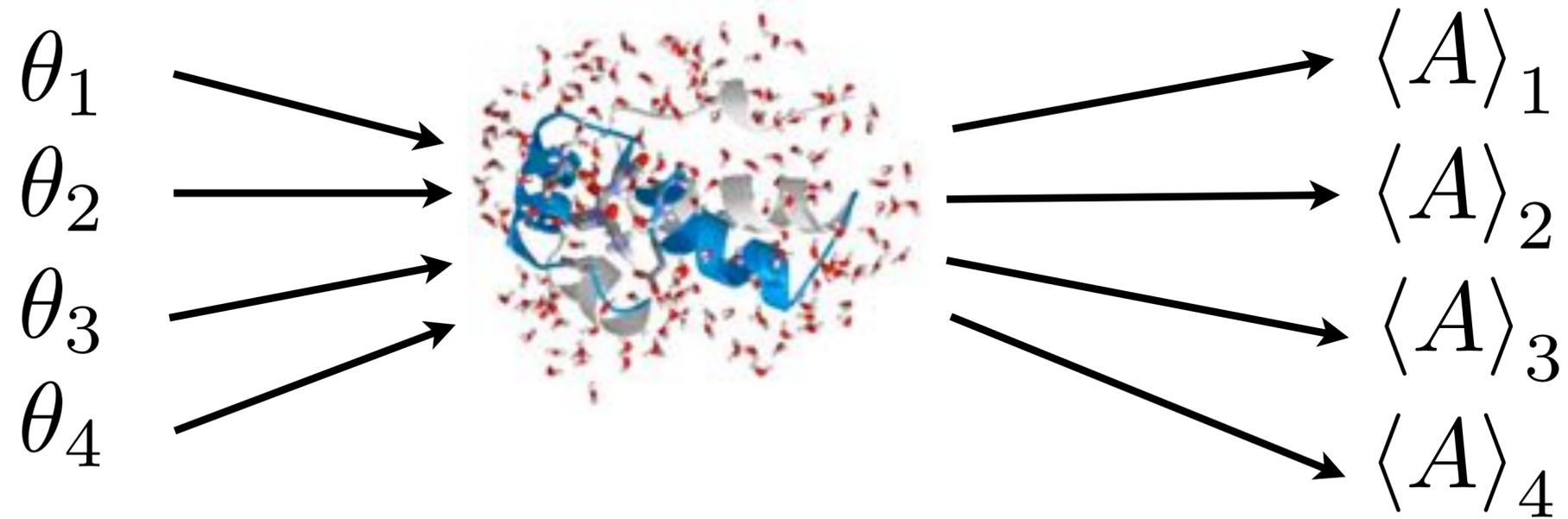
$\theta$  forcefield

$p(\theta|\mathcal{D})$  posterior

$p(\mathcal{D}|\theta)$  data model

$p(\theta)$  prior on forcefield parameters

# THE BAYESIAN WAY



Multiple parameter sets in, multiple estimates out

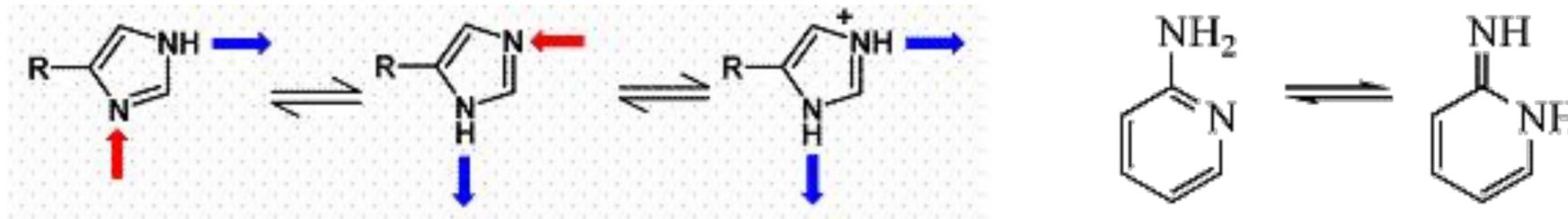
We can estimate both **statistical** and **systematic** components of computed results

# PREDICTIONS FAIL FOR THREE REASONS

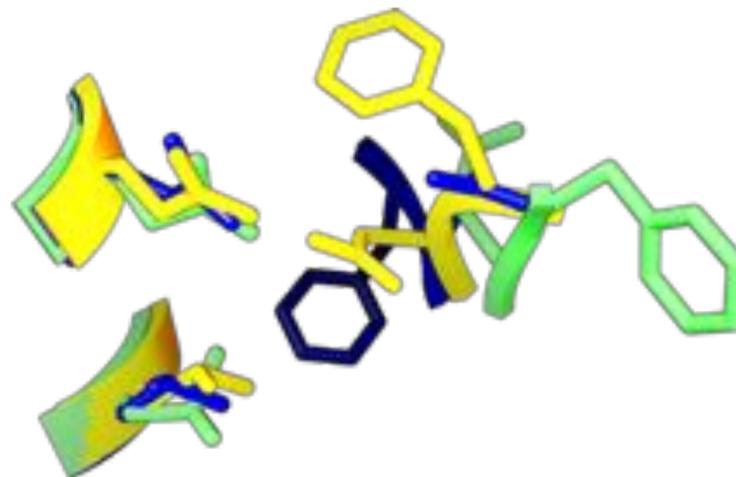
1. The **forcefield** does a poor job of modeling the physics of our system

$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r (r - r_{eq})^2 + \sum_{\text{angles}} K_\theta (\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$

2. We're missing some **essential chemical** in our simulations (e.g. protonation states, tautomers, covalent association)

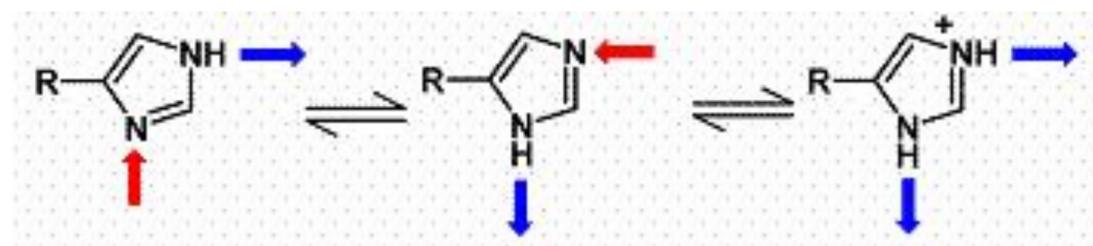


3. We haven't **sampled** all of the relevant conformations



# PREDICTIONS FAIL FOR THREE REASONS

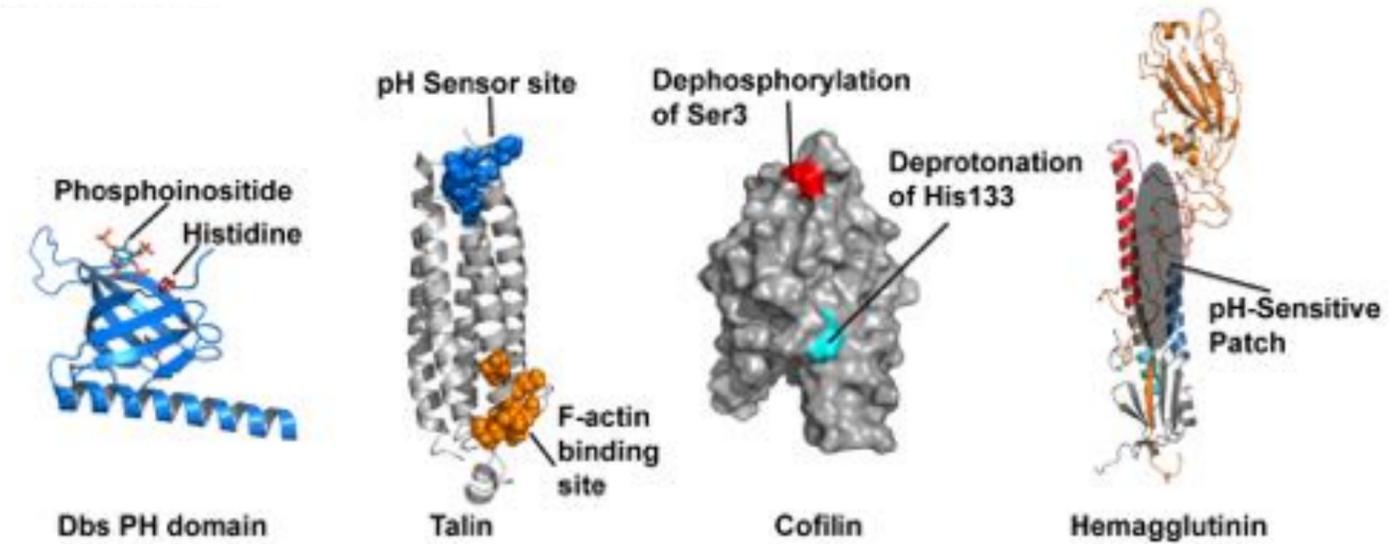
2. We're missing some **essential chemical** in our simulations (e.g. protonation states, tautomers, covalent association)



# PROTONATION STATE EFFECTS ARE UBIQUITOUS IN BIOLOGY

proteins

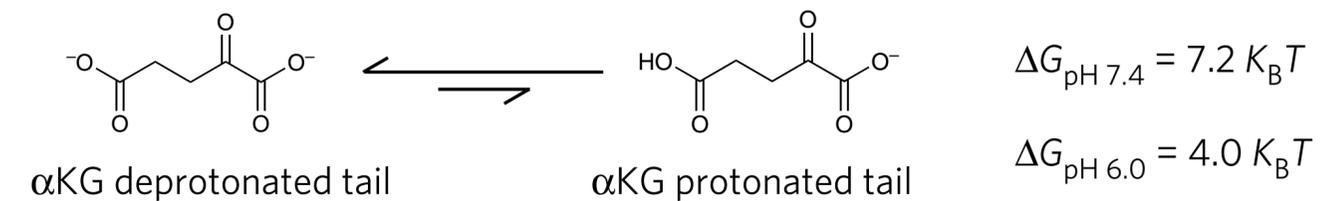
ionizable residues constitute 29% of protein residues



**Diane Barber** // Considering protonation as a post-translational modification regulating protein structure and function. *Annu. Rev. Biophys.* 42:289, 2013.

metabolites

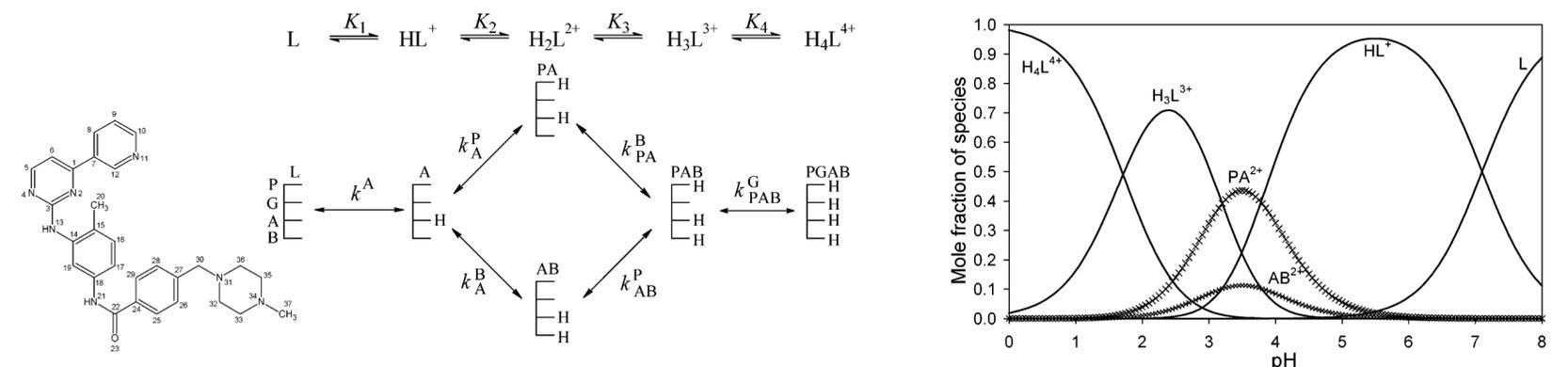
often possess ionizable groups for solubility



**Craig Thompson** // L-2-hydroxyglutarate production arises from noncanonical enzyme function at acidic pH. *Nat. Chem. Biol.* 13:494, 2017.

drugs

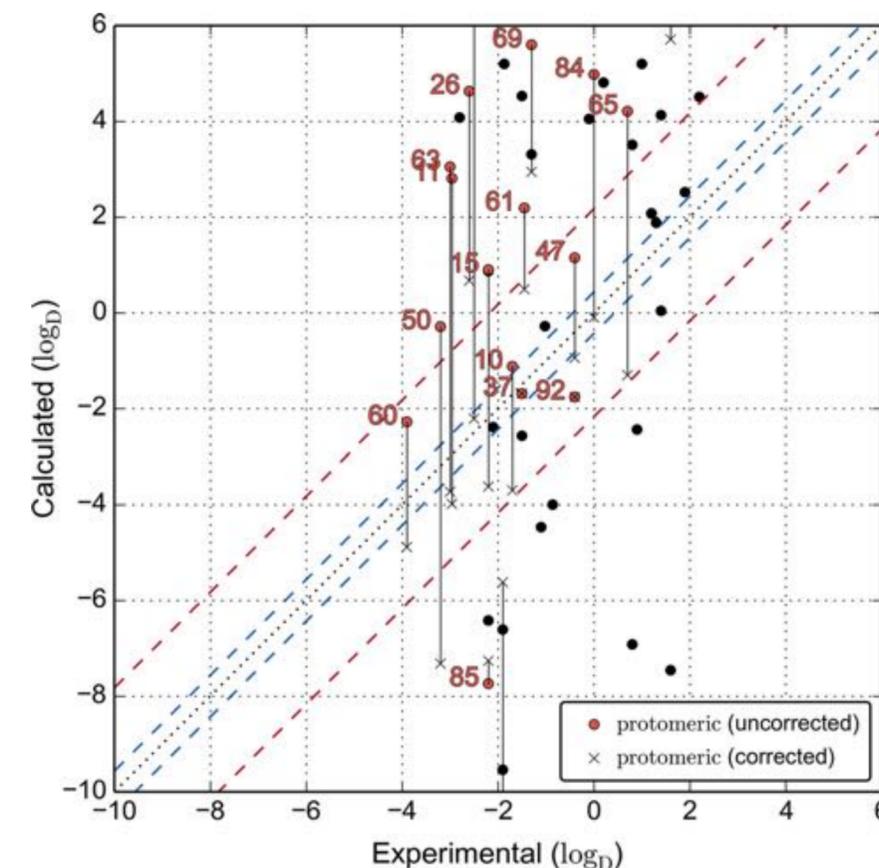
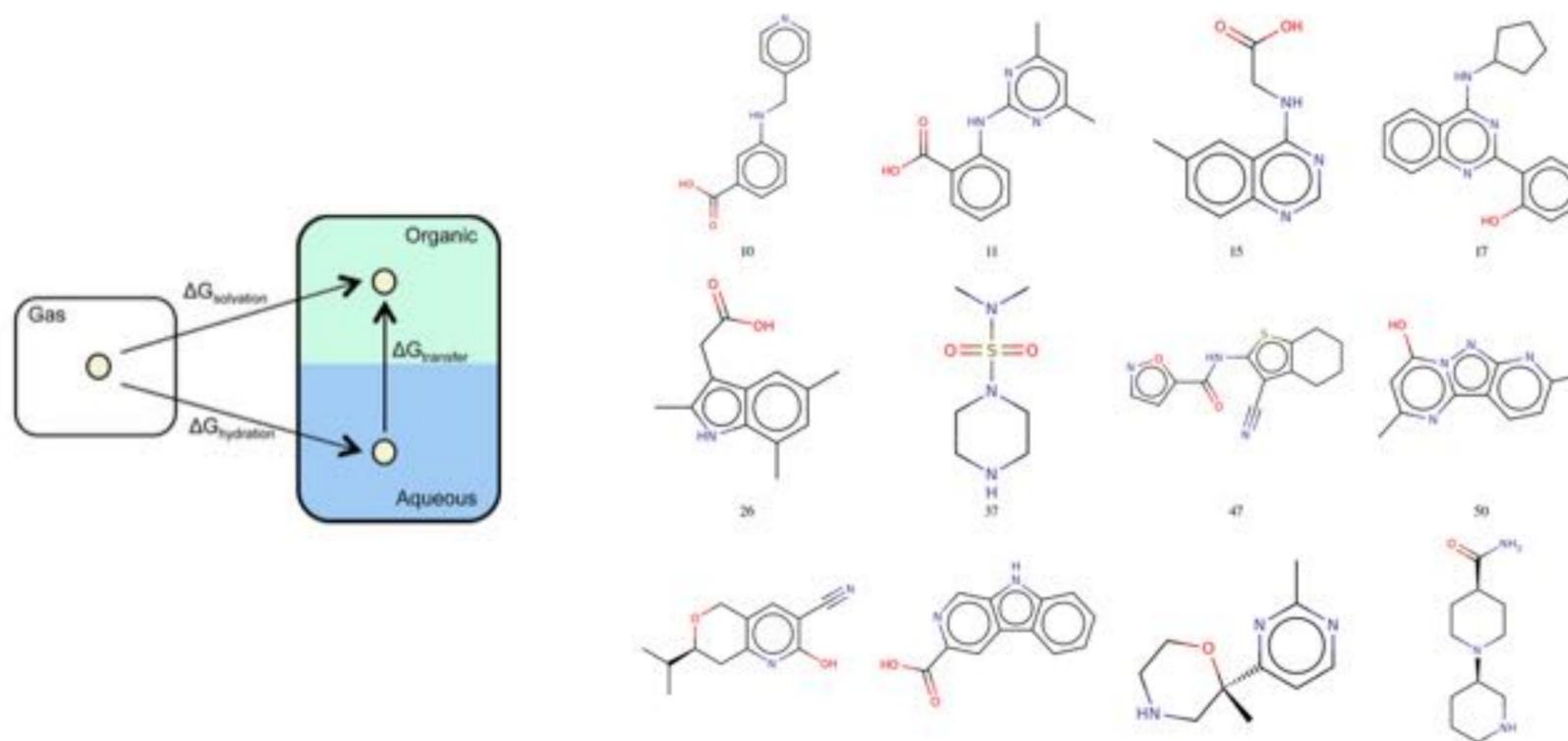
frequently contain many titratable sites  
>25% of approved drugs have more than one tautomer



**Béla Noszál** // Acid-base profiling of imatinib (Gleevec) and its fragments. *J. Med. Chem.* 48:249, 2005.

# NEGLECT OF PROTONATION STATE EFFECTS LEADS TO LARGE ERRORS IN TRANSFER FREE ENERGIES

Blind challenge eliminates slow protein sampling issue to test accuracy of small molecule transfer energetics modeling

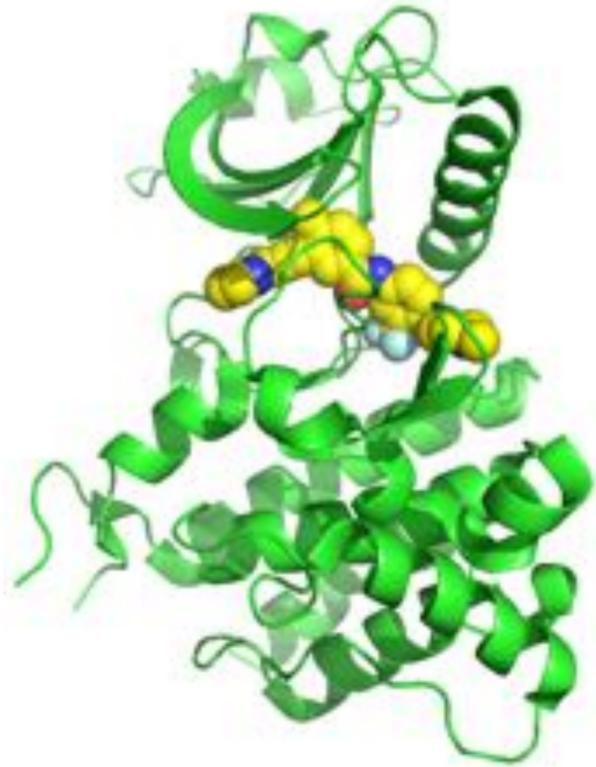


Protonation state effects contribute up to 10 kcal/mol for some compounds  
Including protonation state effects reduced RMSE by 3.5 kcal/mol

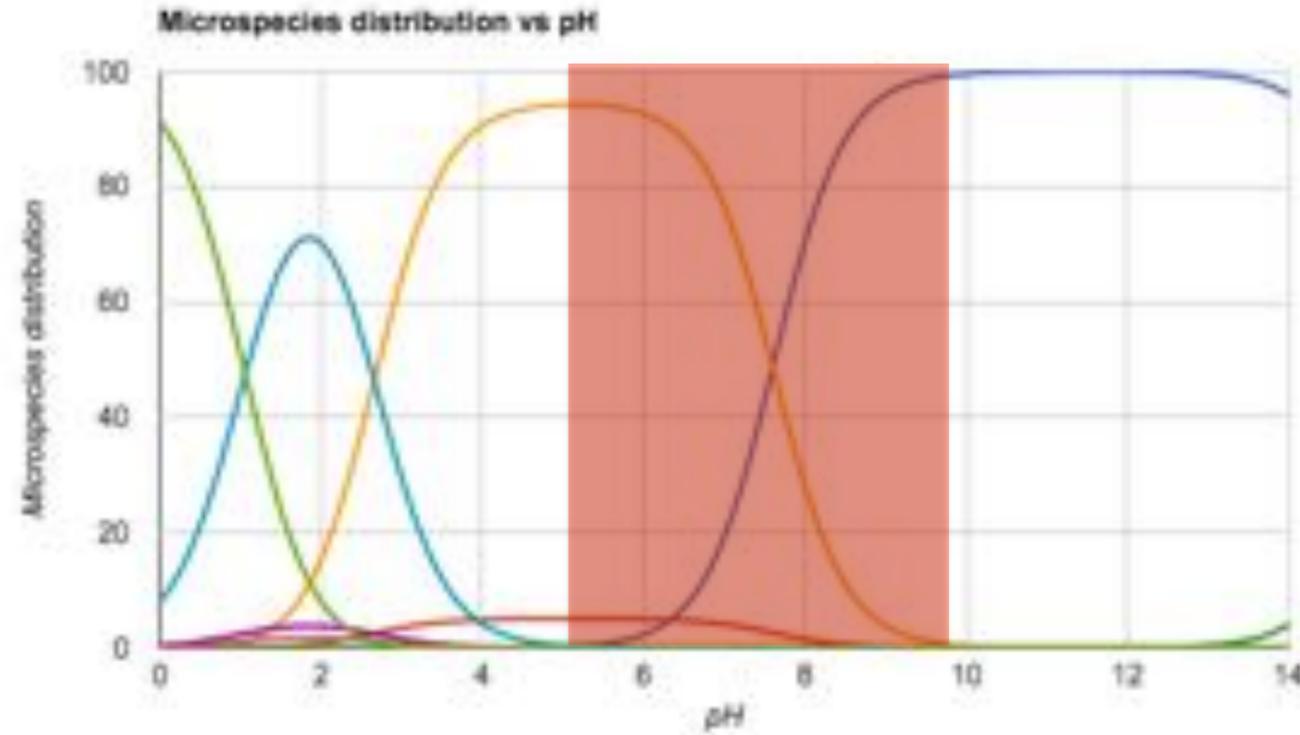
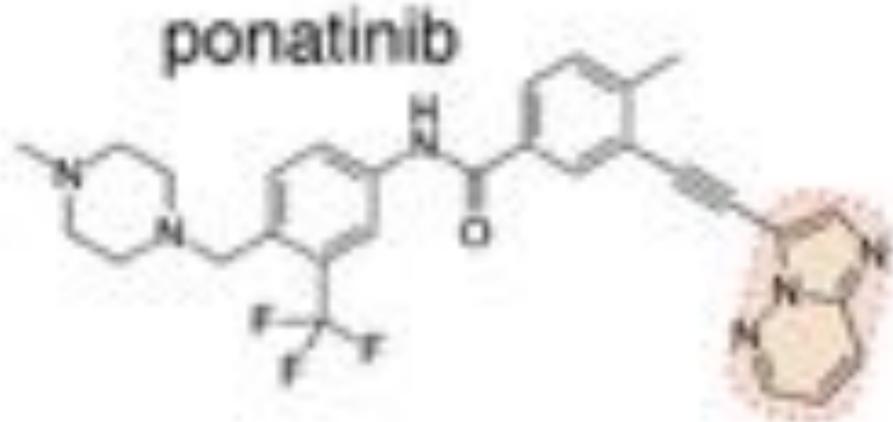
**David Mobley** // Calculating partition coefficients of small molecules in octanol/water and cyclohexane/water. *J. Chem. Theor. Comput.* 12:4015, 2016.

**Bernie Brooks** // Blind prediction of distribution in the SAMPL5 challenge with QM based promoter and pKa corrections. *J. Comput. Aided Mol. Des.* 30:1087, 2016.

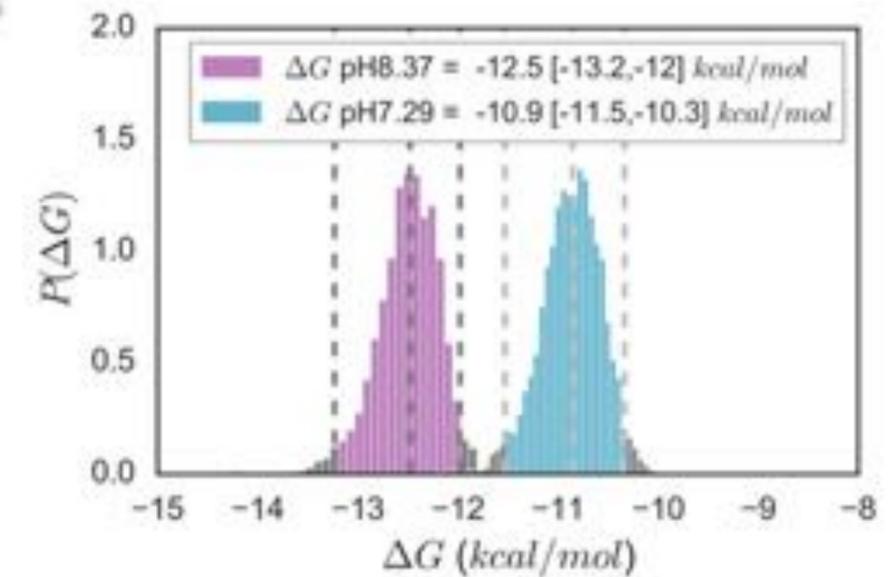
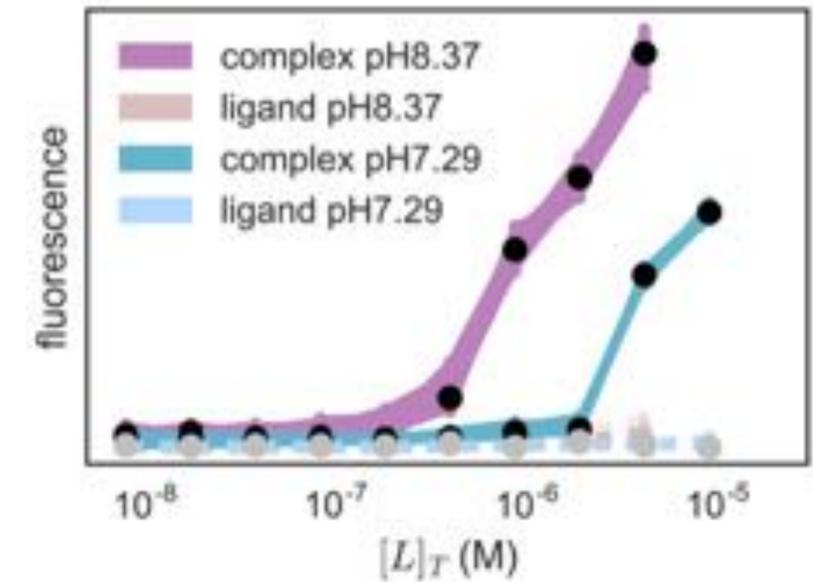
# KINASE:INHIBITOR BINDING CAN BE PH-SENSITIVE



**ponatinib:DDR1**  
(pdbid:3ZOS)



Change pH 7.3 → 8.4



$\Delta\Delta G = 1.6$  kcal/mol

**SONYA HANSON**



# CONTINUUM ELECTROSTATICS PREDICTIONS SUGGEST PROTONATION STATE EFFECTS ARE PERVASIVE

pdbid	Inhibitor	kinase	$\Delta$ protein	$\Delta$ inhibitor	$\Delta$ protomer	pdbid	Inhibitor	kinase	$\Delta$ protein	$\Delta$ inhibitor	$\Delta$ protomer
3UE4	Bosutinib	ABL	0	0.5	YES	4HUO	Erlotinib	EGFR (Inactive)	0	0	YES
2GQG	Dasatinib	ABL	-0.1	0.6	YES	2ITD	Gefitinib	EGFR G719S	0.1	0.48	YES
4XEY	Dasatinib	ABL	0.12	0.82	YES	3UG2	Gefitinib	EGFR G719S/T790M	0.2	0.13	NO
2HYY	Imatinib	ABL	-0.2	-0.01	NO	4G5P	Afatinib	EGFR T790M	-0.5	-0.01	NO
3PYY	Imatinib	ABL	-0.28	0.01	NO	4I22	Gefitinib	EGFR T790M/L858R	0.2	-0.18	NO
3CS9	Nilotinib	ABL	0.1	0.06	NO	4V01	Ponatinib	FGFR1	-2.6	0.06	YES
3DX2	Ponatinib	ABL	-0.6	0.02	NO	4V04	Ponatinib	FGFR1	-1.2	0.05	YES
3IK3	Ponatinib	ABL T315I	-0.63	0.06	NO	4QRC	Ponatinib	FGFR4	-1	0.02	YES
3ADX	Alectinib	ALK	0	0.13	NO	4TYJ	Ponatinib	FGFR4	-0.3	0.03	NO
4MKC	Ceritinib	ALK	0.7	0	NO	4UXQ	Ponatinib	FGFR4	-0.36	0.02	NO
2XP2	Crizotinib	ALK	-0.04	-0.77	YES	3LXX	Tofacitinib	JAK3	-0.02	-0.07	NO
4ANQ	Crizotinib	ALK G1269A	-0.1	-0.76	YES	4U0I	Ponatinib	KIT	-0.1	0.03	NO
2YFX	Crizotinib	ALK L1196M	-0.1	-0.77	YES	4AN2	Cobimetinib	MEK1	0	0.01	NO
4AN5	Crizotinib	ALK L1196M/G1269A	-0.1	-0.77	YES	4LMN	Cobimetinib	MEK1	0	0	NO
4XV2	Dabrafenib	BRAF	-0.92	-0.31	NO	2WGJ	Crizotinib	MET	-0.06	-1.05	YES
5CSW	Dabrafenib	BRAF	0.4	0.65	YES	4AG8	Axitinib	VEGFR2	0.12	0	NO
5HIE	Dabrafenib	BRAF	1	-0.25	NO	4AGC	Axitinib	VEGFR2	0	0	NO
2EUF	Palbociclib	CDK6	-0.08	-0.28	NO	3WZD	Lenvatinib	VEGFR2	0.18	0	NO
3ZOS	Ponatinib	DDR1	-1.5	-0.23	NO	3CJG	Pazopanib	VEGFR2	0.34	-0.02	NO
4G5I	Afatinib	EGFR	-0.18	-0.98	YES	2QU5	Regorafenib	VEGFR2	0.5	-0.07	YES
1M17	Erlotinib	EGFR	0.2	0	NO	3WZE	Sorafenib	VEGFR2	0.1	-0.01	NO
4WKQ	Gefitinib	EGFR	0	0.65	YES	4ASD	Sorafenib	VEGFR2	0.2	-0.01	NO
1XKK	Lapatinib	EGFR	-0.26	-0.54	YES	4AGD	Sunitinib	VEGFR2	0.32	-0.99	YES
4ZAU	Osimertinib	EGFR	-0.3	0.02	NO						
2ITY	Gefitinib	EGFR	0.08	-0.04	NO						
2ITZ	Gefitinib	EGFR	0	0.09	NO						

proton gain

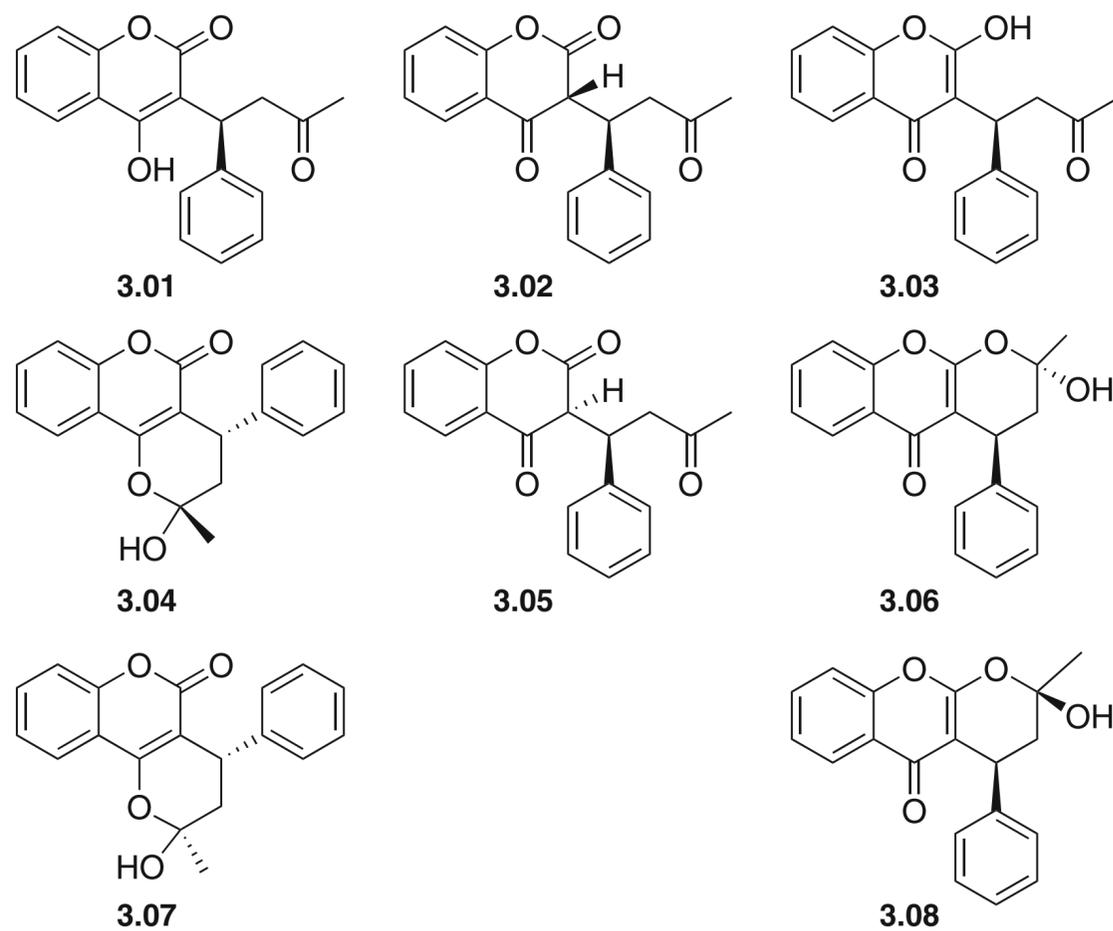
proton loss

tautomer shift

MARILYN GUNNER  
SALAH SALAH



# LET'S NOT FORGET TAUTOMERS



tautomers of warfarin

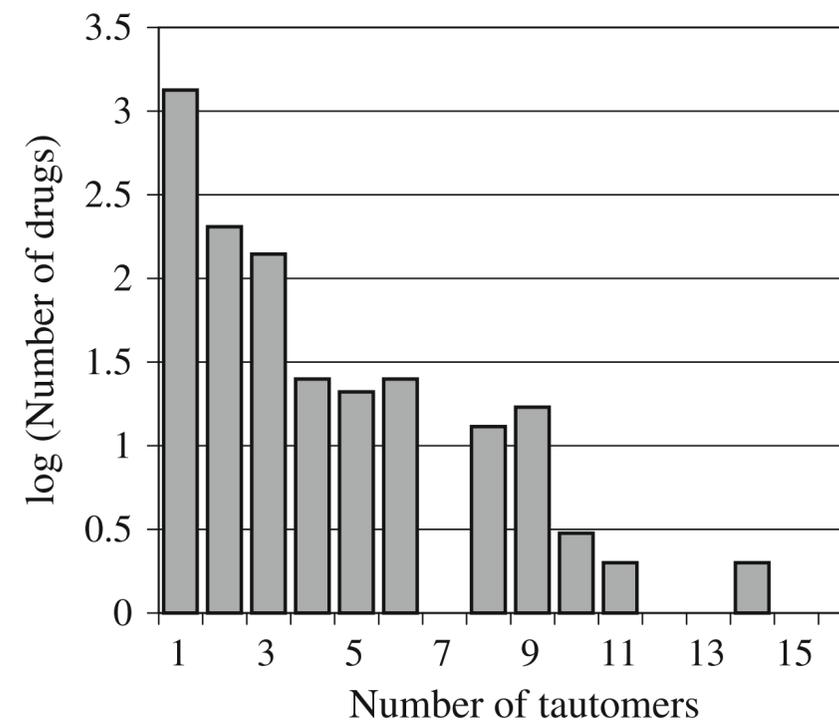


Fig. 13 The frequency distribution of tautomers of marketed drugs

**MORE THAN HALF OF ALL DRUGS  
HAVE 2 OR MORE TAUTOMERS**

# WE CAN SIMULATE AN EXPANDED ENSEMBLE WHERE PROTONATION STATES ARE DYNAMIC

Define the **reduced potential** for a state  $k$  as a combination of terms

$$u_k(\mathbf{x}) = \beta_k [U_k(\mathbf{x}) + p_k V(\mathbf{x}) + \mu_k^T \mathbf{N}(\mathbf{x})]$$

$\mu_k$  chemical potential of exchangeable species

$\mathbf{N}(\mathbf{x})$  number of each chemical species in system     **protons**

The target configuration space density is given by

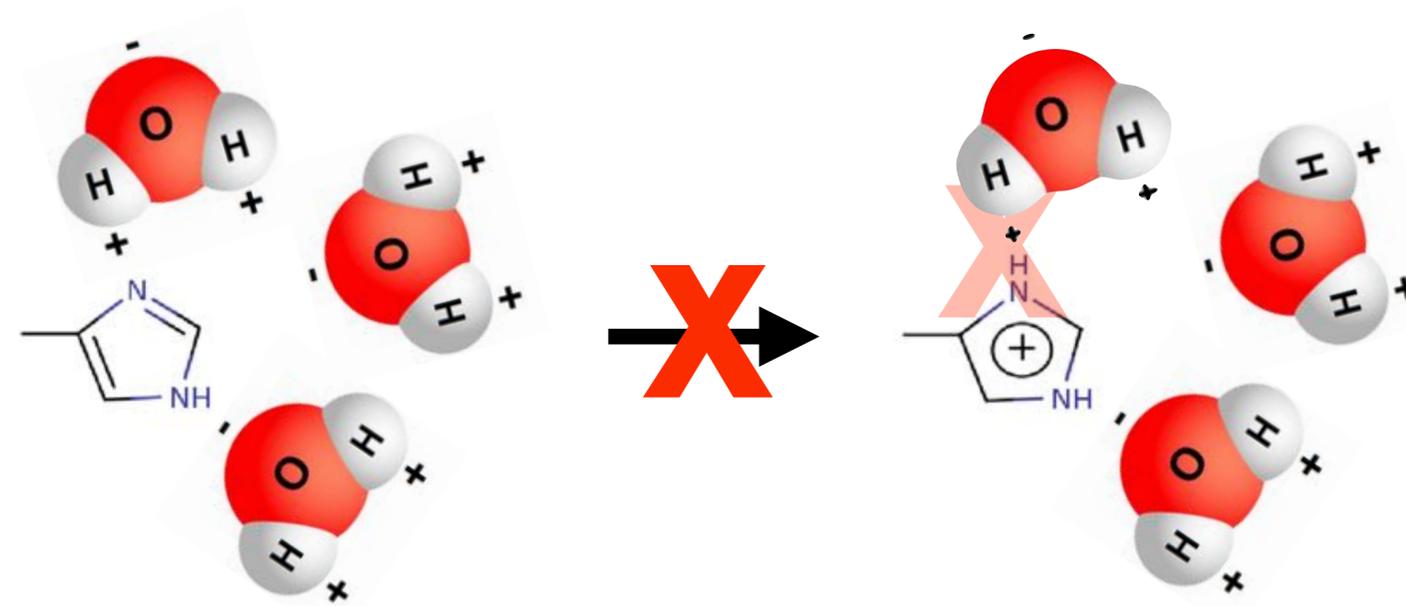
$$\pi_k(\mathbf{x}) = Z_k^{-1} \exp[-u_k(\mathbf{x})] \quad Z_k = \int d\mathbf{x} \exp[-u_k(\mathbf{x})]$$

We can sample from this ensemble using **Gibbs sampling**

$\mathbf{x}_{n+1} \sim \pi(\mathbf{x} | \mathbf{N}_n)$      **update positions with MD/MC at fixed protonation state**

$N_{n+1} \sim \pi(\mathbf{N} | \mathbf{x}_{n+1})$      **update protonation states with MC at fixed position**

# PROBLEM: ACCEPTANCE RATES FOR PROTONATION STATE CHANGES ARE ESSENTIALLY ZERO IN EXPLICIT SOLVENT



**INSTANTANEOUS CHANGES IN PROTONATION STATE ARE HIGHLY UNFAVORABLE DUE TO SOLVATION SHELL**

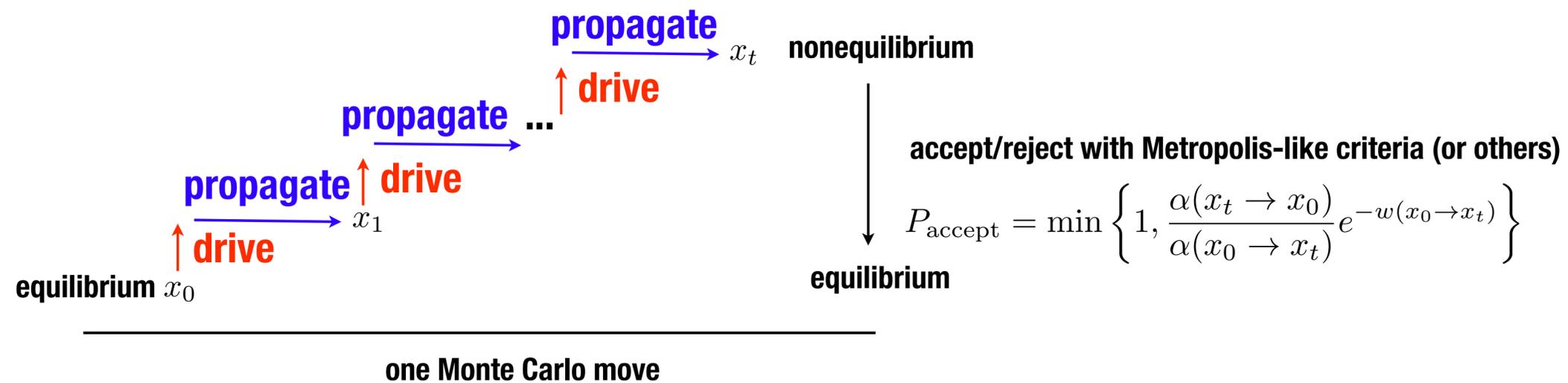
# NONEQUILIBRIUM CANDIDATE MONTE CARLO CAN ASTRONOMICALLY BOOST MC ACCEPTANCE RATES

Algorithm:

**Drive** some degrees of freedom or thermodynamic parameters in small steps, accumulating work

**Propagate** others using Metropolis MC or molecular dynamics

**Accept or reject** final configuration with Metropolis-like criterion



Follows earlier ideas by Manuel Athenes (work-bias Monte Carlo), Harry Stern (constant pH simulation), Chris Jarzynski (switching replica temperatures), and Jerome Nilmeier (approximate).

**Nilmeier, Crooks, Minh, Chodera. PNAS 108:E1009, 2011.**

# THE ACCEPTANCE CRITERIA CAN BE GENERALIZED TO ALLOW ARBITRARY PROPAGATION KERNELS

$x_t \in \mathcal{X}$  **configurations** in some state space  
 $\alpha_t(x, y)$  **perturbation kernel**  $\sum_y \alpha_t(x, y) = 1$   
 $K_t(x, y)$  **propagation kernel**  $\sum_x \pi_t(x) K_t(x, y) = \pi_t(y)$   
 $\Lambda \equiv \{\alpha_0, K_0, \dots, \alpha_T, K_T\}$  **switching protocol**  
 $X \equiv \{x_0, x_1^*, x_1, \dots, x_T\}$  **trajectory**

$\tilde{x}$  **is momentum-reversed**  $x$   
 $\alpha_t(x, y) > 0 \Leftrightarrow \alpha_t(y, x) > 0$   
 $K_t(x, y) > 0 \Leftrightarrow K_t(y, x) > 0$   
 $\tilde{\Lambda}$  **is time-reversed**  
 $\tilde{X}$  **is time-reversed**

**forward process**  $x_0 \xrightarrow{\alpha_1} x_1^* \xrightarrow{K_1} x_1 \longrightarrow \dots \longrightarrow x_{T-1} \xrightarrow{\alpha_T} x_T^* \xrightarrow{K_T} x_T.$   
**reverse process**  $\tilde{x}_T \xrightarrow{K_T} \tilde{x}_T^* \xrightarrow{\alpha_T} \tilde{x}_{T-1} \longrightarrow \dots \longrightarrow \tilde{x}_1 \xrightarrow{K_1} \tilde{x}_1^* \xrightarrow{\alpha_1} \tilde{x}_0.$

$P(\tilde{\Lambda}|\tilde{x}_T) > 0 \Leftrightarrow P(\Lambda|x_0 > 0)$   
**both must be selected with nonzero probability!**

Enforce strict “pathwise” form of detailed balance (which also ensures detailed balance is satisfied):

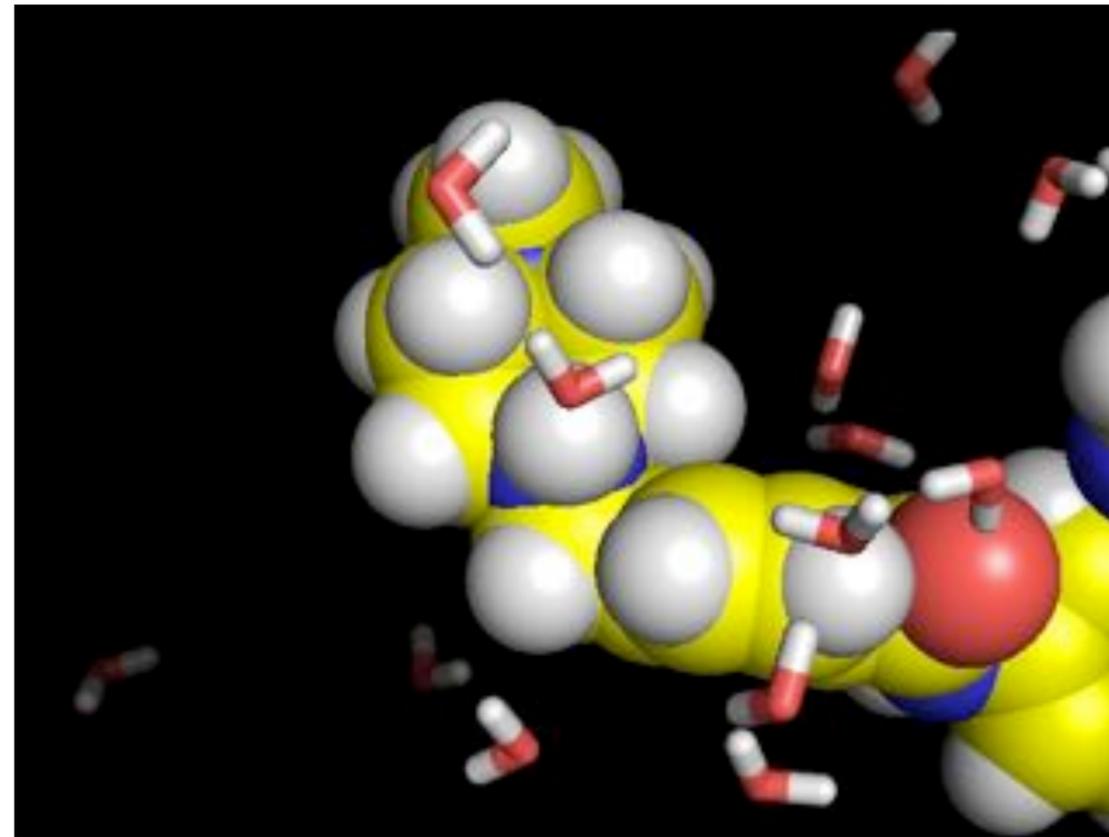
$$\underbrace{A(X|\Lambda)}_{\text{acceptance probability}} \underbrace{\Pi(X|x_0, \Lambda)}_{\text{path generation probability}} \underbrace{P(\Lambda|x_0, \lambda_0)}_{\text{protocol selection}} \underbrace{\pi(x_0, \lambda_0)}_{\text{equilibrium}} = \underbrace{A(\tilde{X}|\tilde{\Lambda})}_{\text{acceptance probability}} \underbrace{\Pi(\tilde{X}|\tilde{x}_T, \tilde{\Lambda})}_{\text{path generation probability}} \underbrace{P(\tilde{\Lambda}|\tilde{x}_T, \lambda_T)}_{\text{protocol selection}} \underbrace{\pi(\tilde{x}_T, \lambda_T)}_{\text{equilibrium}}$$

Result is a general acceptance criteria for any nonequilibrium perturbation:

$$A(X|\Lambda) = \min \left\{ 1, \frac{\pi(x_T, \lambda_T)}{\pi(x_0, \lambda_0)} \frac{P(\tilde{\Lambda}|\tilde{x}_T, \lambda_T)}{P(\Lambda|x_0, \lambda_0)} \frac{\tilde{\alpha}(\tilde{X})}{\alpha(X)} e^{-\Delta S(X)} \right\}$$

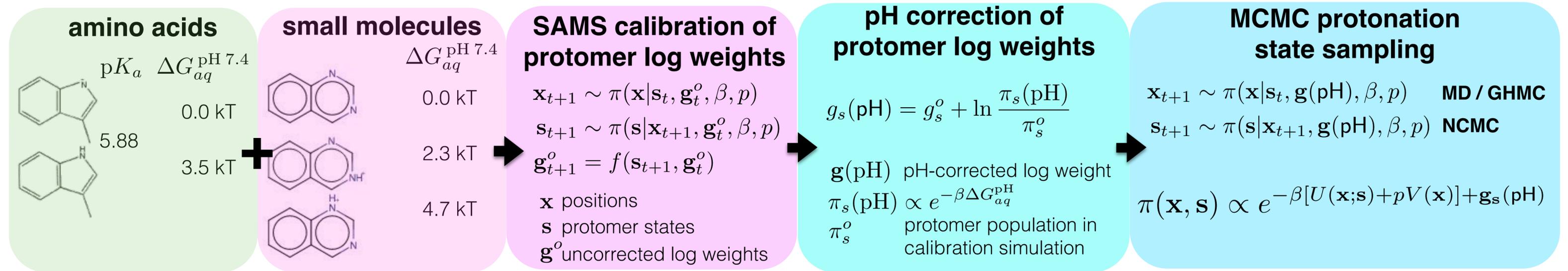
equilibrium
selection
perturbation action

# ANNIHILATING A PROTON IN EXPLICIT SOLVENT OVER 20 PS NCMC SWITCHING TRAJECTORY



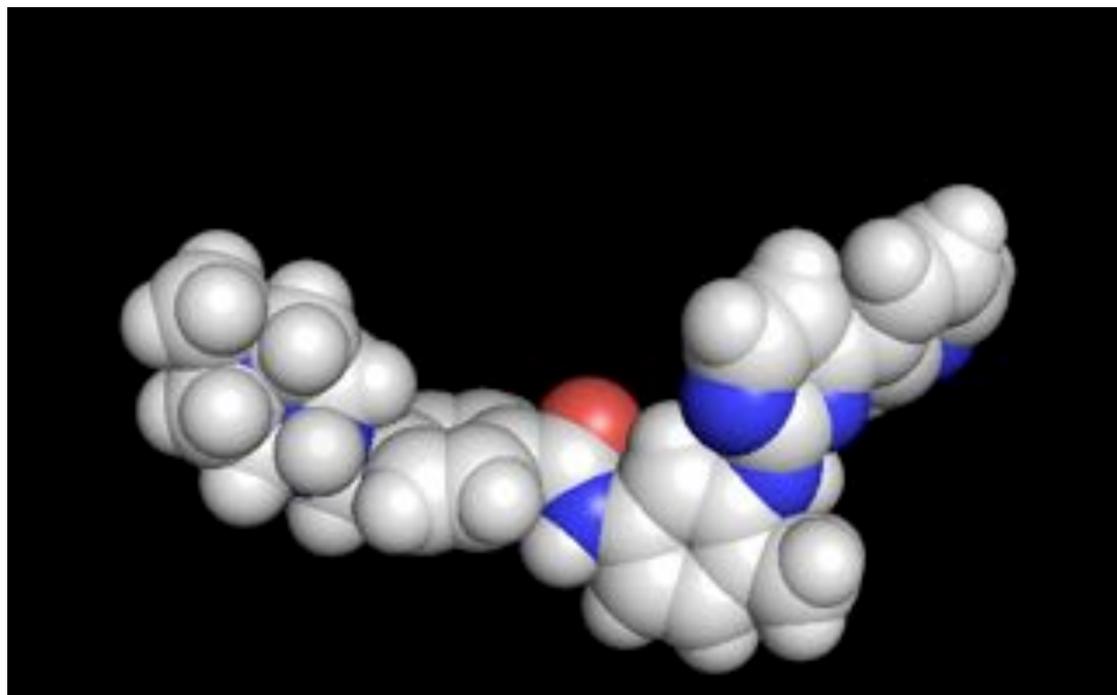
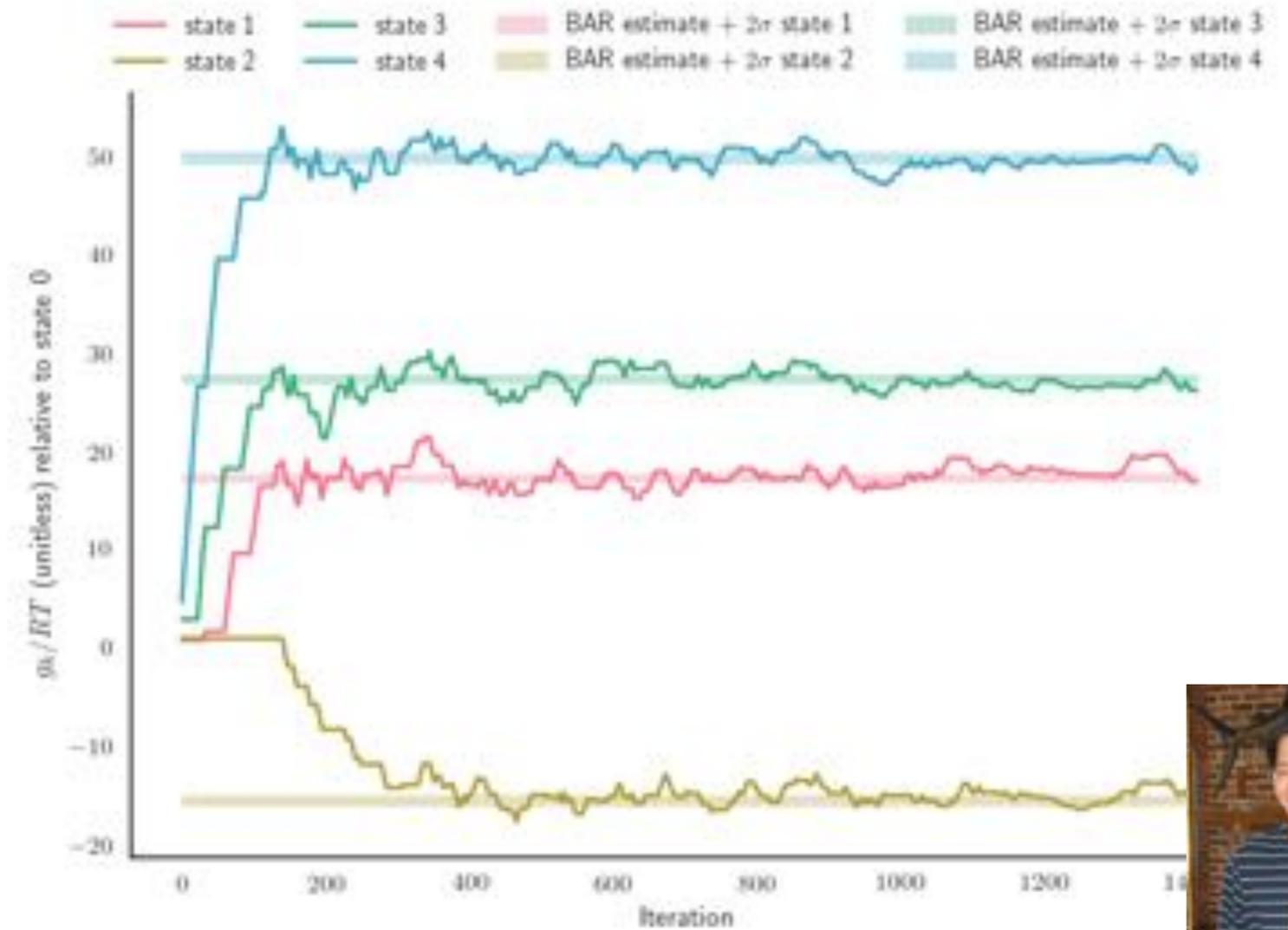
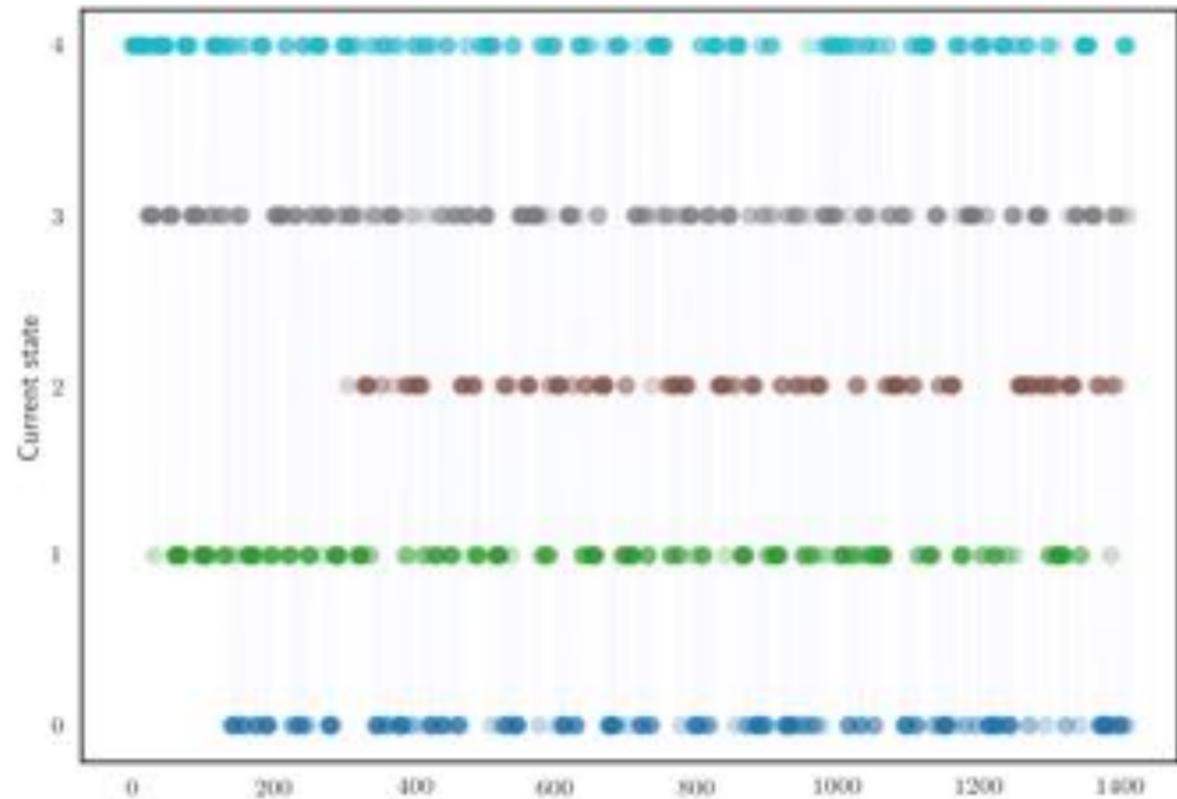
**BAS RUSTENBURG**

# DYNAMIC LIGAND PROTONATION STATES REQUIRE AN EXPLICIT CALIBRATION STEP



**BAS RUSTENBURG**

# CALIBRATION USES SAMs + BAR TO ESTIMATE RELATIVE PROTONATION STATE FREE ENERGIES

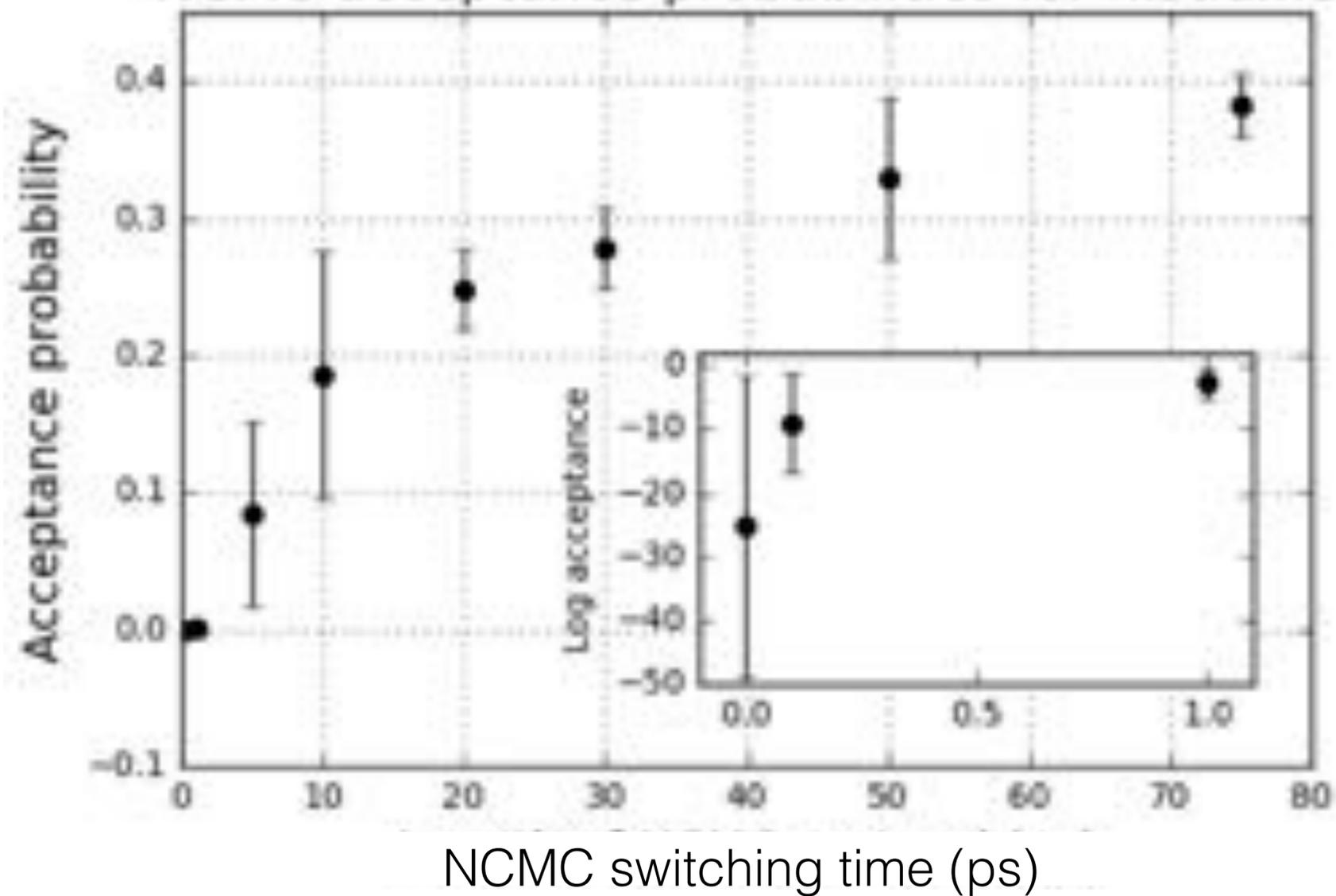


**ZHIQIANG TAN**  
Rutgers



# NCMC GIVES USEFUL ACCEPTANCE RATES

**A** NCMC acceptance probability for histidine in explicit solvent



**BAS RUSTENBURG**

# PROTON-DRIVE

## Protons: Protonation states and tautomers for OpenMM

### Note:

This module is undergoing heavy development. None of the API calls are final.

### Introduction

This python module implements a constant-pH MD scheme for sampling protonation states and tautomers of amino acids and small molecules in OpenMM.

### Installation

Use the command

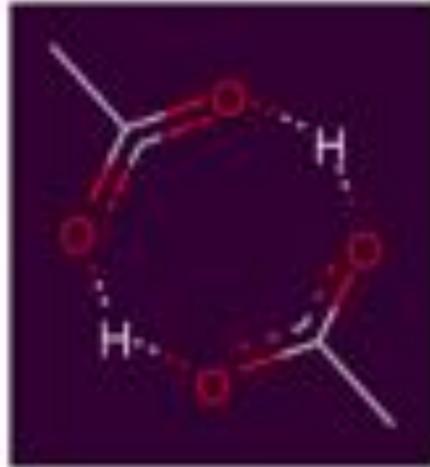
```
python setup.py install
```

to install the package. The installation does not automatically check for requirements.

To test the installation, run

```
nosetests protons
```

### Requirements



### Protons

Protonation states and tautomers for OpenMM

Watch

### Navigation

Setting up a constant-pH MD simulation

Advanced calibration options

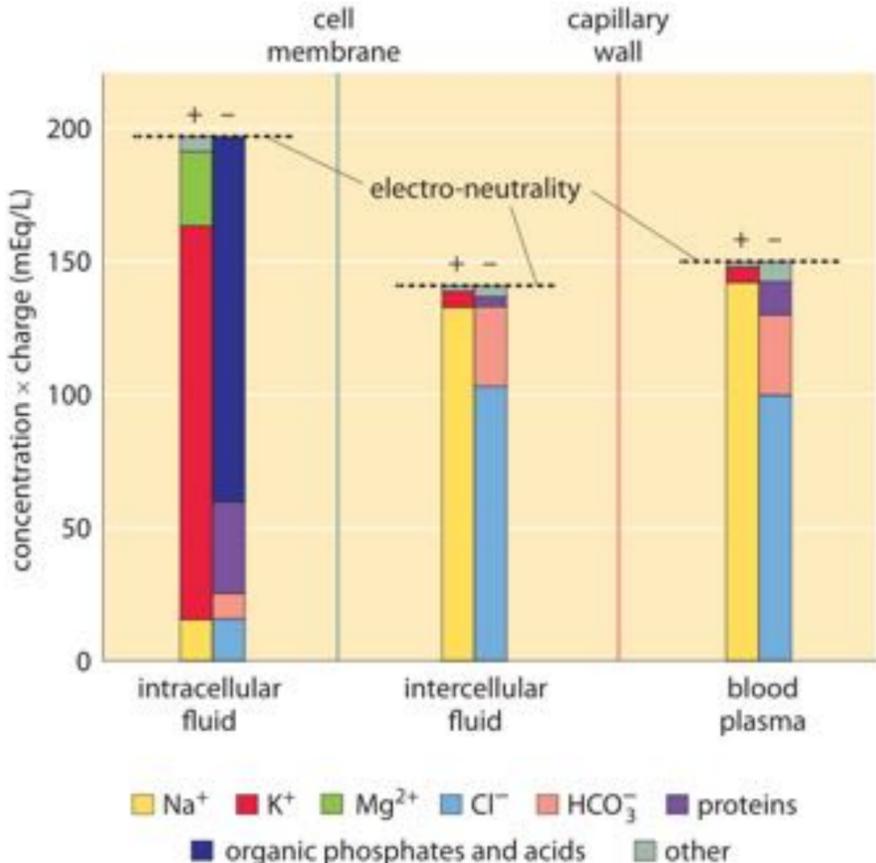
The ligutils submodule



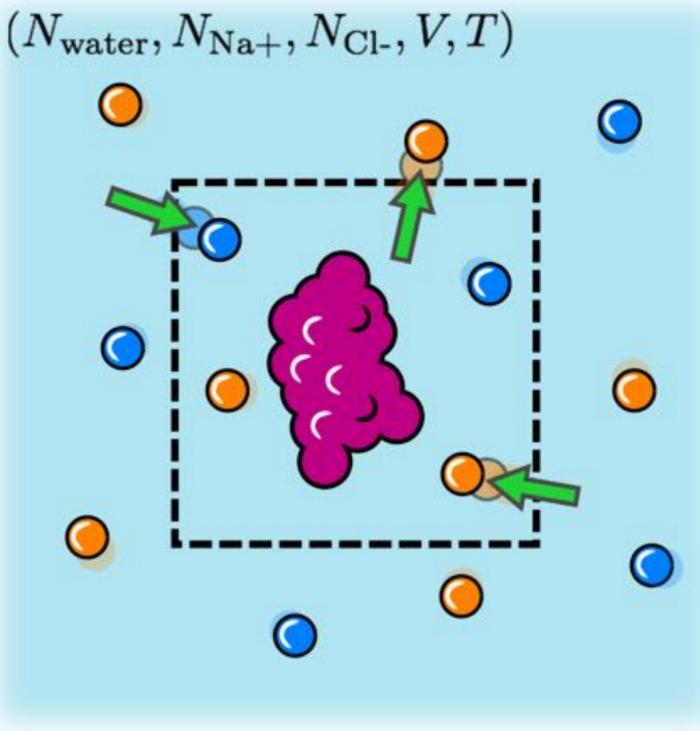
**BAS RUSTENBURG**

<https://github.com/choderalab/protons>

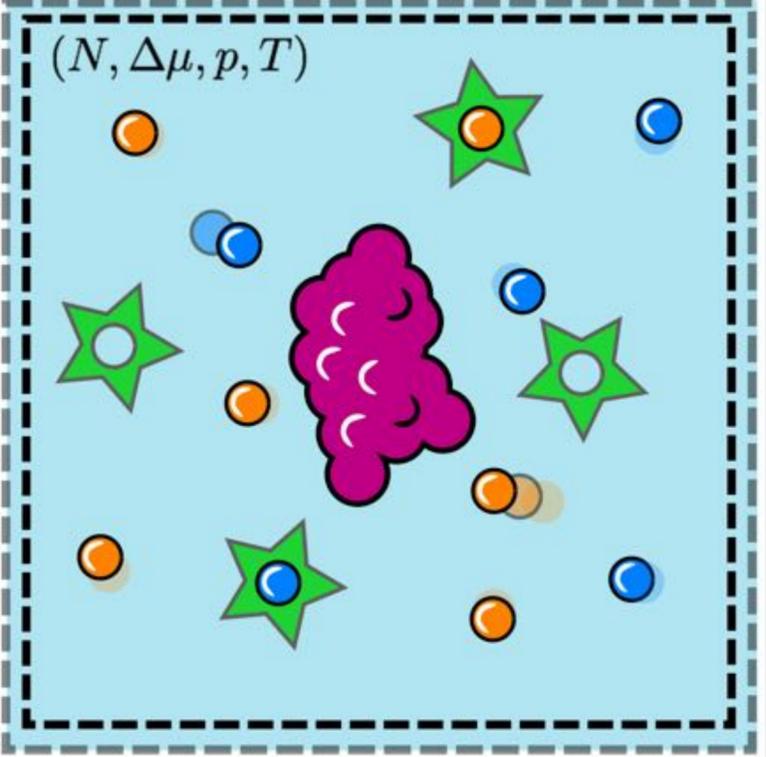
# AN NCMC OSMOSTAT CAN MODEL REALISTIC SALT CONCENTRATION FLUCTUATIONS



real system



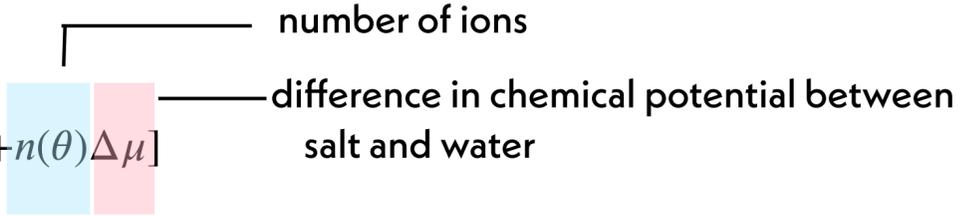
semigrand ensemble



semigrand ensemble

$$\pi(x, \theta; \Delta\mu, N, p, T) = \frac{1}{\Xi(\Delta\mu, N, p, T)} e^{-\beta[U(x, \theta) + pV(x) + n(\theta)\Delta\mu]}$$

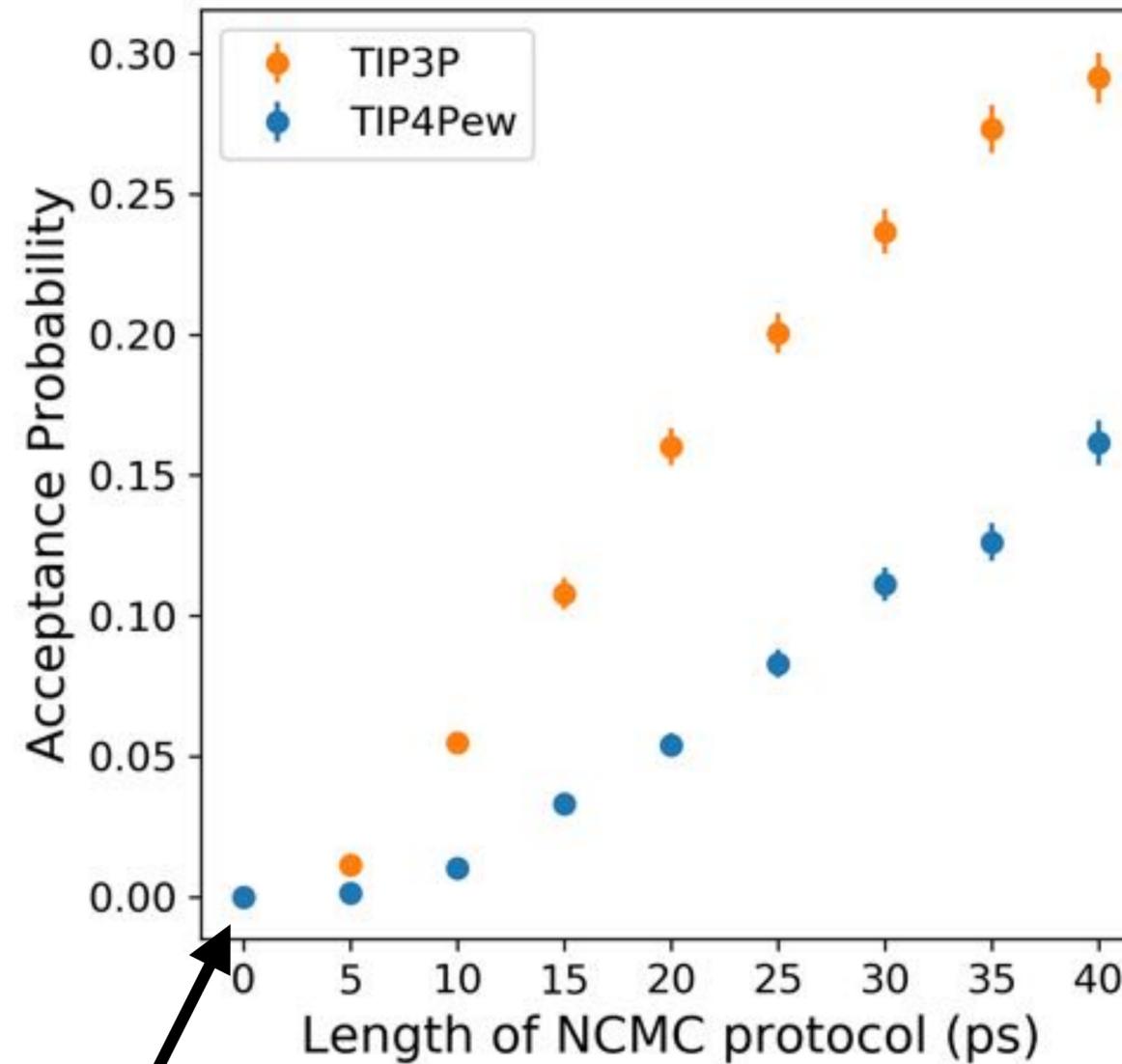
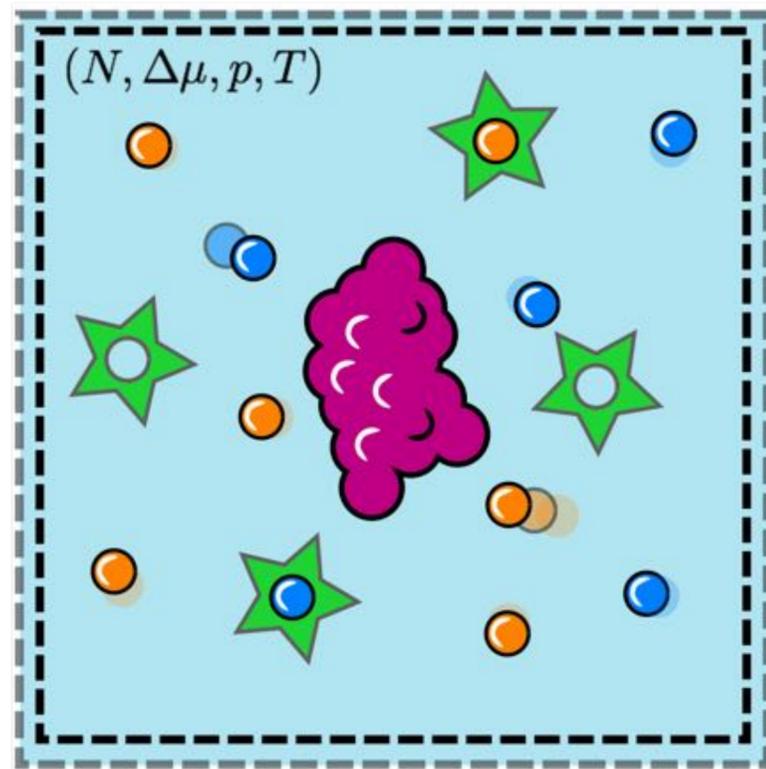
$$\Xi(\Delta\mu, N, p, T) = \sum_{\theta, \text{subject to } \sum_i^N \theta_i = -z} \int dx \pi(x, \theta, \Delta\mu, N, p, T),$$



# AN NCMC OSMOSTAT CAN MODEL REALISTIC SALT CONCENTRATION FLUCTUATIONS

NCMC acceptance probability  
for water-salt exchange

semigrand ensemble

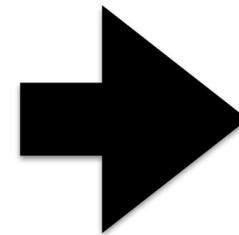
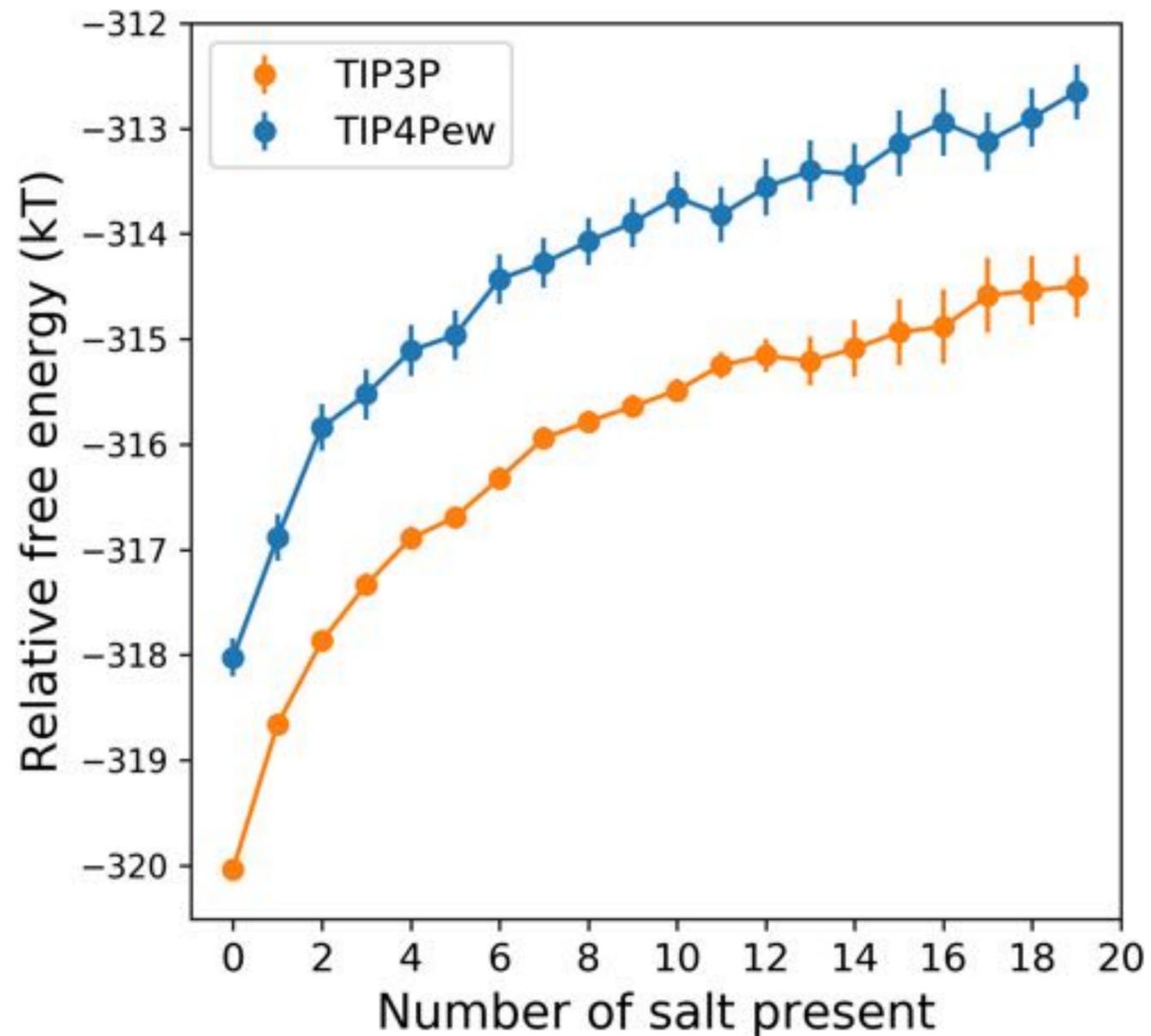


**GREGORY ROSS**

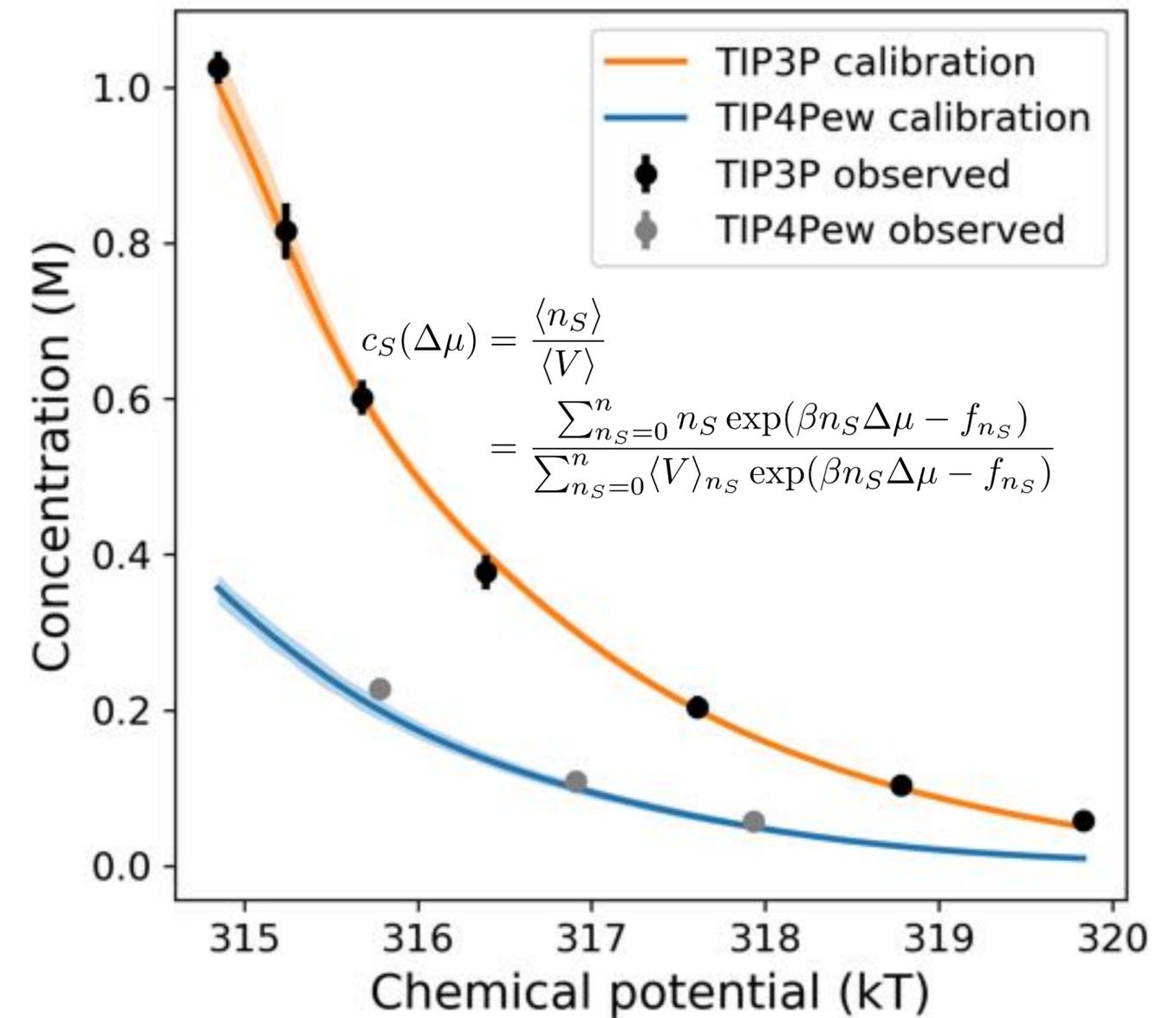
$10^{-22}$  if instantaneous MC used

# AN NCMC OSMOSTAT CAN MODEL REALISTIC SALT CONCENTRATION FLUCTUATIONS

SAMS+BAR estimates free energy change for inserting salt pair

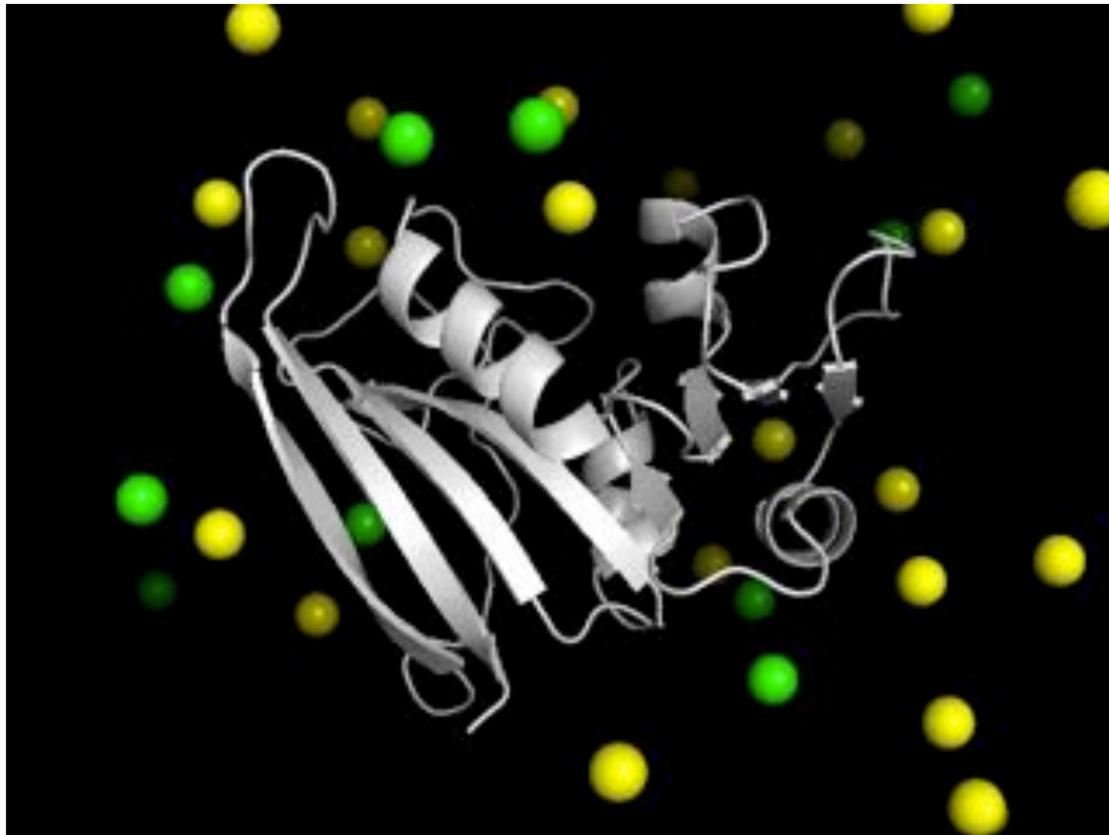


Macroscopic salt concentration vs chemical potential

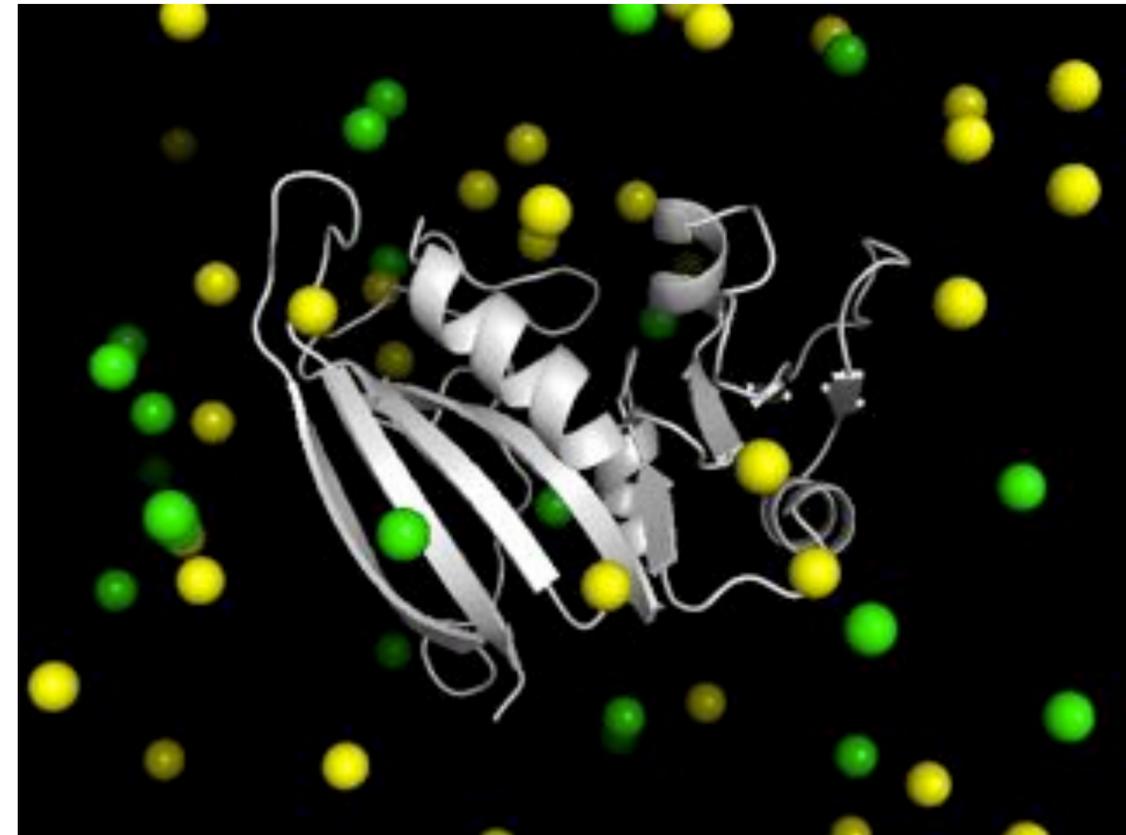


# COUNTERION ENVIRONMENTS NEAR BIOMOLECULES ARE HIGHLY DYNAMIC

100 mM NaCl



200 mM NaCl



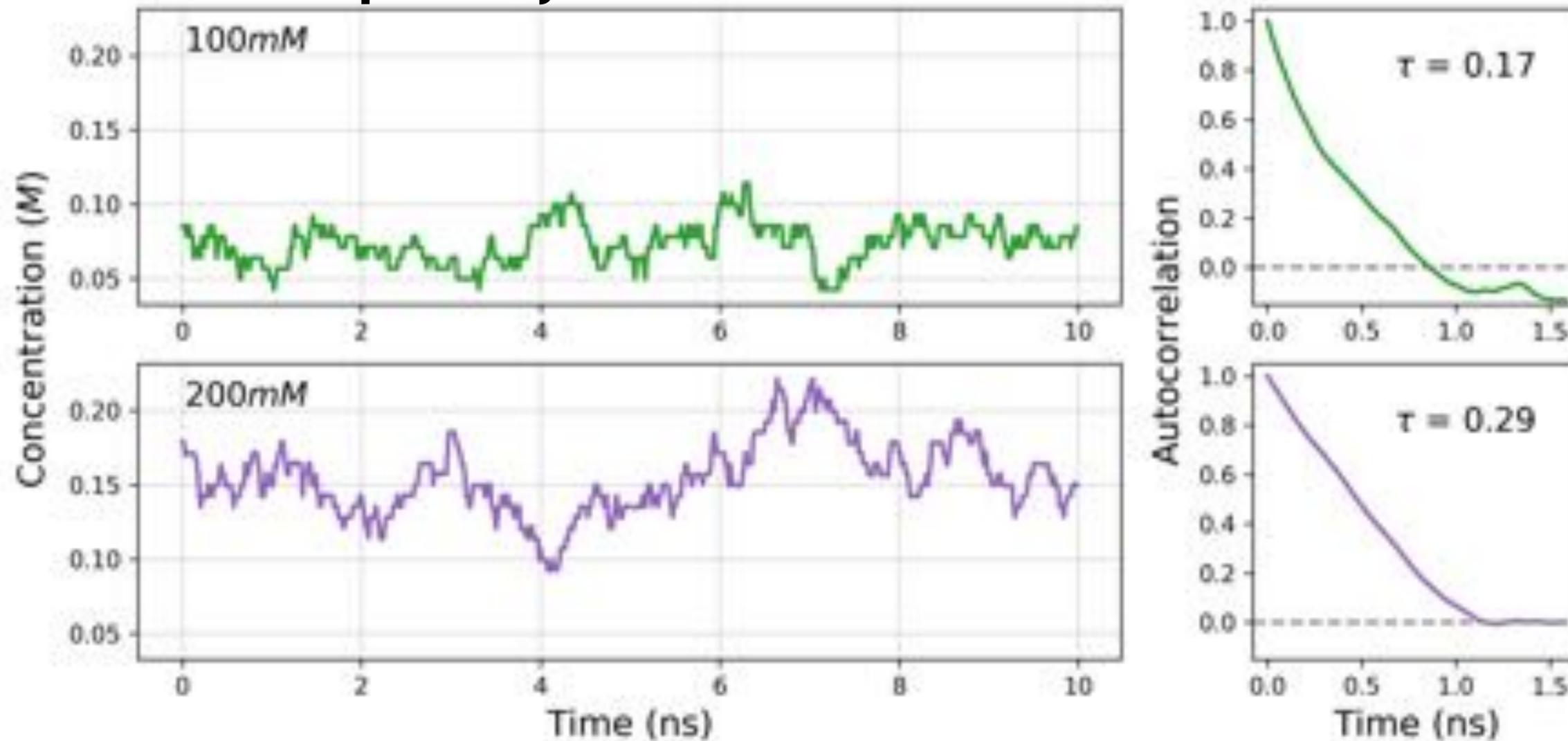
DHFR (from JAC benchmark) in TIP3P with PME  
AMBER99SB-ILDN with Cheatham-Joung ion parameters

**GREGORY ROSS**



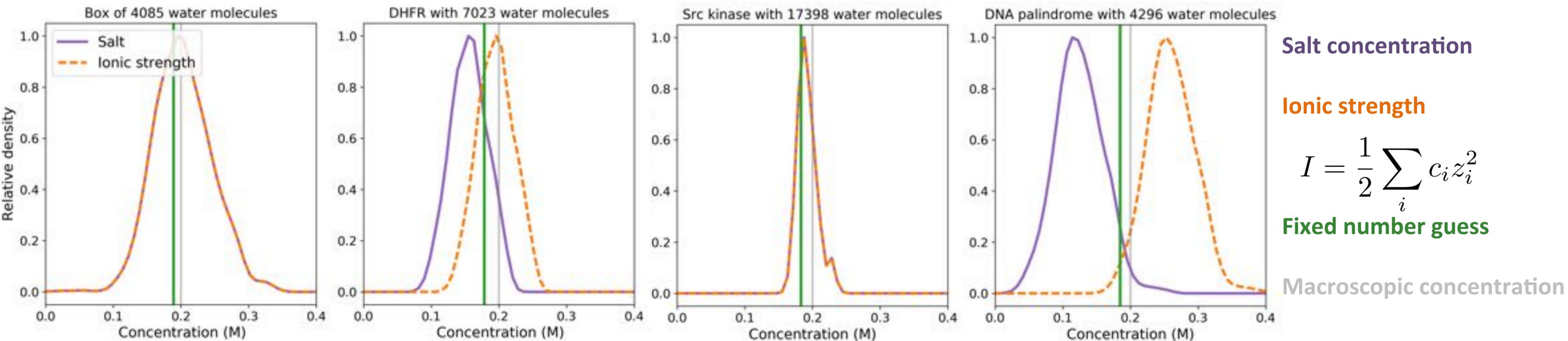
# COUNTERION ENVIRONMENTS NEAR BIOMOLECULES ARE HIGHLY DYNAMIC

salt pair trajectories for DHFR



**GREGORY ROSS**

# COUNTERION ENVIRONMENTS NEAR BIOMOLECULES ARE HIGHLY DYNAMIC

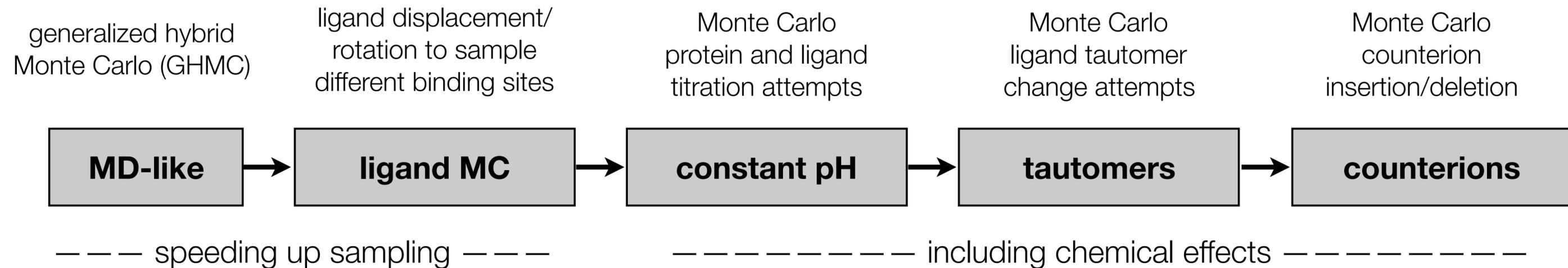


Real salt environments can differ substantially from our typical fixed-number guess



**GREGORY ROSS**

# MARKOV CHAIN MONTE CARLO (MCMC) PROVIDES A FLEXIBLE FRAMEWORK FOR ENHANCEMENTS



**WE CAN USE FREE ENERGY CALCULATIONS AND EXPERIMENTS TO **QUANTIFY**  
**THE ERROR** IN NEGLECTING OF PROTOMERS AND TAUTOMERS**

# SAMPL6 IS COMING SOON

Model systems of **intermediate complexity** to focus community on challenges in blind tests

## Model protein-ligand systems

Isolate individual physical challenges (e.g. binding of charged ligands)

## Physical properties

Tests of forcefield accuracy in hydrated or protein-like environments

Isolate chemical effects (protonation states, ligand conformations) **without** slow protein timescales

## Host-guest systems

Binding of small drug-like molecules with protein-like affinities, **without** slow protein timescales

**SAMPL0**  
2007

JNK3 kinase inhibitors  
hydration free energies

**SAMPL1**  
2008

CDK2 kinase inhibitors  
hydration free energies

**SAMPL2**  
2009

hydration free energies  
tautomer ratios

**SAMPL3**  
2011

trypsin inhibitors  
hydration free energies

**SAMPL4**  
2013

HIV-1 integrase inhibitors  
hydration free energies  
octaacid host-guest  
CB7 host-guest

**SAMPL5**  
2016

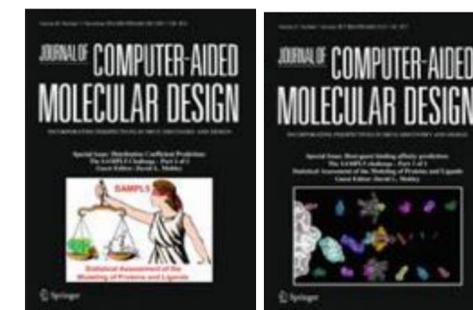
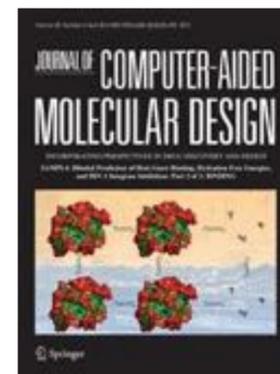
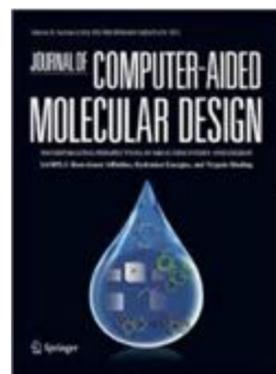
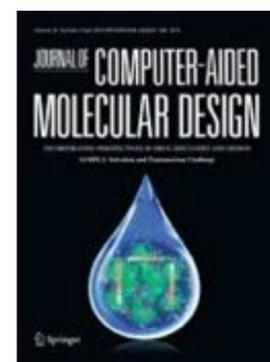
distribution coefficients  
CBClip host-guest  
CB7 host-guest

**SAMPL6**  
2017

distribution coefficients  
solubilities (CheqSol)  
octaacid host-guest  
CB8 host-guest

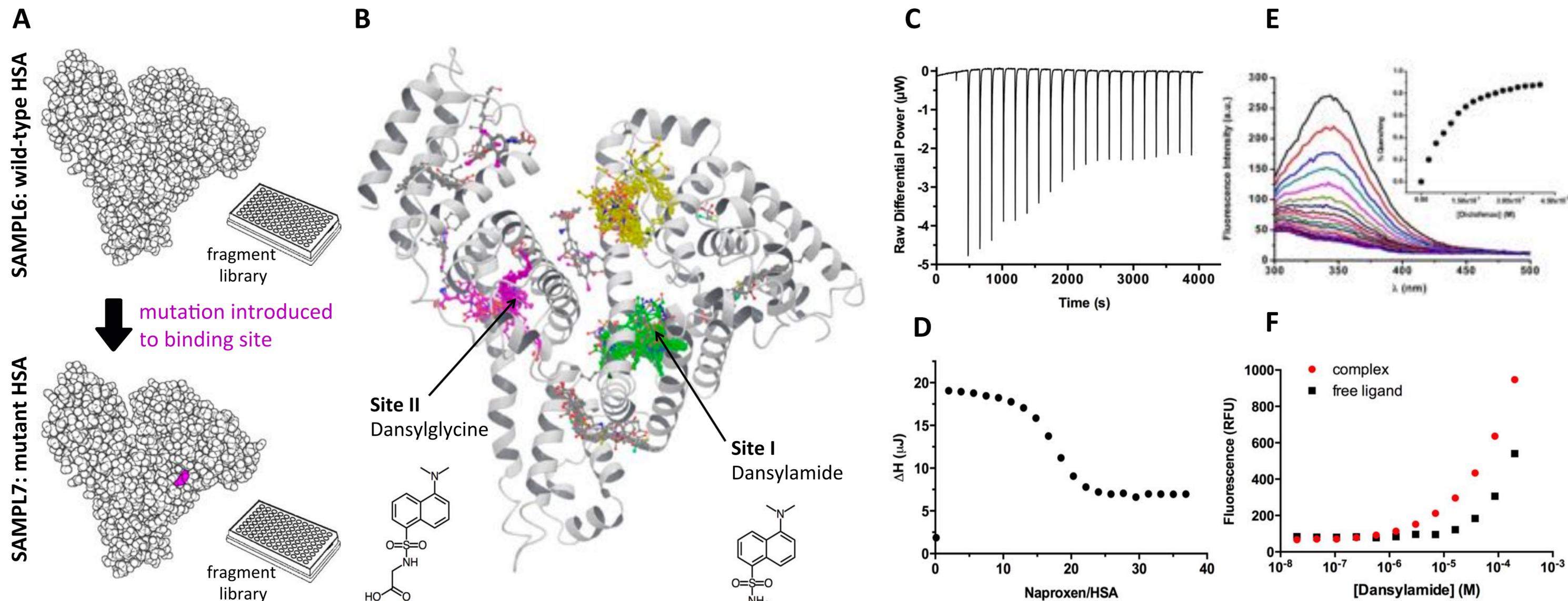
**SAMPL7**  
2018

human serum albumin  
pKas  
distribution coefficients  
solubilities (CheqSol)  
octaacid host-guest  
CB8 host-guest



<https://drugdesigndata.org/about/sampl>

# SAMPL7/8: HUMAN SERUM ALBUMIN



# THE CHODERA LAB @ MSKCC



**Josh Fass**  
**Bas Rustenburg**  
**Andrea Rizzi**  
**Patrick Grinaway**  
**Greg Ross**  
**Levi Naden**

Code and data available at <http://www.choderalab.org>

Start Folding at <http://folding.stanford.edu>

# WE MAKE TOOLS FOR THE COMMUNITY

OpenMM (GPU-accelerated MD with Python API) : <https://openmm.org>

openmmtools (integrators, alchemy, samplers) : <http://openmmtools.readthedocs.io>

YANK (absolute alchemical free energy calculations) : <http://getyank.org>

Software best practices : <https://github.com/choderalab/software-development>

Omnia consortium : <http://omnia.md>

Open Forcefield Group : <https://github.com/open-forcefield-group>