

<https://github.com/amartinsson/DukeMicroCourse>

Sampling Algorithms “Microcource”

Session 2

Ben Leimkuhler (U of Edinburgh)
Tutorials by: Matthias Sachs (SAMSI) and
Anton Martinssen (U of Edinburgh)

Duke University, Spring 2018

Overview

Lecture 1: Monte Carlo, Hamiltonian systems, splitting methods, microcanonical sampling, SDEs, Euler-Maruyama & Leimkuhler-Matthews, Langevin splitting, superconvergence

Tutorial 1: methods for Langevin (and overdamped) dynamics

Lecture 2: Noisy gradients, then:
**increasing the timestep, multiple timestepping,
holonomic constraints, isokinetic constraints & SIN,
ensemble preconditioning**

Tutorial 2: enhancing the convergence to equilibrium/
decreasing the integrated autocorrelation time

Langevin Dynamics

$$dq = M^{-1} pdt$$

$$dp = F(q)dt - \gamma M^{-1} pdt + \sqrt{2\beta^{-1}\gamma} dW$$

Newton's Equations

Stochastic Perturbation

With Periodic Boundary Conditions
and smooth potential, ergodic sampling
of the canonical distribution with density

$$\rho_{\text{can}} \propto e^{-\beta(p^T M^{-1} p / 2 + U(q))}$$

Mattingly, Stuart, Higham, Hairer, Nier...

Splitting Methods

$$\begin{aligned} dq &= pdt \quad (\text{A}) \quad (M=I) \\ dp &= F(q)dt \quad (\text{B}) \quad -\gamma pdt + \sqrt{2\beta^{-1}\gamma} dW \quad (\text{O}) \end{aligned}$$

$$\mathbf{A} : \quad \dot{q} = p, \quad \dot{p} = 0$$

$$\mathbf{B} : \quad \dot{q} = 0, \quad \dot{p} = F(q)$$

$$\mathbf{O} : \quad \dot{q} = 0, \quad \dot{p} = -\gamma p + \sqrt{2\beta^{-1}\gamma} \dot{W}$$

Time stepsize δt

$$p := e^{-\gamma\delta t} p + \sqrt{(1 - e^{-2\gamma\delta t})\beta^{-1}} \mathcal{N}(0, 1)$$

$$\mathbf{ABO} \quad p := \delta t F(q)$$

$$q := q + \delta t p$$

Splitting Methods

$$\mathcal{L} = \mathcal{A} + \mathcal{B} + \mathcal{O}$$

$$\mathcal{A} = p^T M^{-1} \nabla_q \quad \mathcal{B} = -\nabla U(q)^T \nabla_p \quad \mathcal{O} = -\gamma p^T M^{-1} \nabla_p + \beta^{-1} \Delta_p$$

Propagator:

$$\mathcal{P}_t = e^{t\mathcal{L}}$$

Splitting Method:

$$\mathcal{P}_t \approx e^{t\mathcal{A}} e^{t\mathcal{B}} e^{t\mathcal{O}}$$

Drift Kick Shuffle “OBA”

“ABO” “OAB” “ABOBA” “OBABO” ...

Expansion of the invariant distribution

(Talay-Tubaro expansion¹ in the ergodic limit)

$$[\mathcal{L}^\dagger + \delta t^2 \mathcal{L}_2^\dagger + \dots] e^{-\beta(H + \delta t^2 f_2 + \dots)} = 0$$

Leading order:

$$\mathcal{L}^\dagger(\rho_{\text{can}} f_2) = \beta^{-1} \mathcal{L}_2^\dagger \rho_{\text{can}}$$

L. & Matthews, AMRX, 2013

L., Matthews, & Stoltz, IMA J. Num. Anal. 2015

- detailed treatment of all 1st and 2nd order splittings
- estimates for the operator inverse discrete inv. measure
- treatment of nonequilibrium (e.g. transport coefficients)

¹Denis Talay and Luciano Tubaro. Expansion of the global error for numerical schemes solving stochastic differential equations. *Stochastic Analysis and Applications* 4 (1990)

Two level expansion

For each of **ABOBA** and **BAOAB**, we find the first terms of the two-level expansion

$$\hat{\rho} = \exp(-\beta[H + \delta t^2(f_{2,0} + f_{2,1}\varepsilon + O(\varepsilon^2)) + O(\delta t^4)])$$

$$\varepsilon = 1/\gamma$$

$$f_{2,0} \equiv f_{2,0}^{\text{BAOAB}} = \frac{1}{8} (p^T U''(x) p - \beta^{-1} \Delta U(x)),$$

$$f_{2,1} \equiv f_{2,1}^{\text{BAOAB}} = \frac{1}{24} \beta^{-1} p^T \nabla_x \Delta_x U(x) - \frac{1}{72} p^T \nabla_x p^T U''(x) p,$$

$$f_{2,2} \equiv f_{2,2}^{\text{BAOAB}} = \frac{1}{296} p^T \nabla_x p^T \nabla_x p^T U''(x) p - \frac{1}{48} \nabla U(x) \cdot \nabla_x p^T U''(x) p.$$

Configurational Sampling

Integrate out with respect to momenta...and
discover a **surprise**:

Proposition:

The right order of Langevin building blocks provides substantial improvements in accuracy (low sampling bias).

In the high friction limit: 4th order, and with just one force evaluation per timestep.

2.1 Noisy gradients

Problem: use stochastic dynamics to accurately sample a distribution with given positive smooth density

$$\rho \propto \exp(-U)$$

in case the force $-\nabla U$ can only be computed approximately

Examples:

Multiscale models

several flavors of hybrid **ab initio MD Methods**

QM/MM methods

...Many applications in **Bayesian Inference & Big Data Analytics**

What to do about the force error?

$$\tilde{F}(x) = -\nabla U(x) + \eta(x)$$

a sampling error... it seems natural to take

$$\eta(x) \sim \mathcal{N}(0, \sigma(x))$$

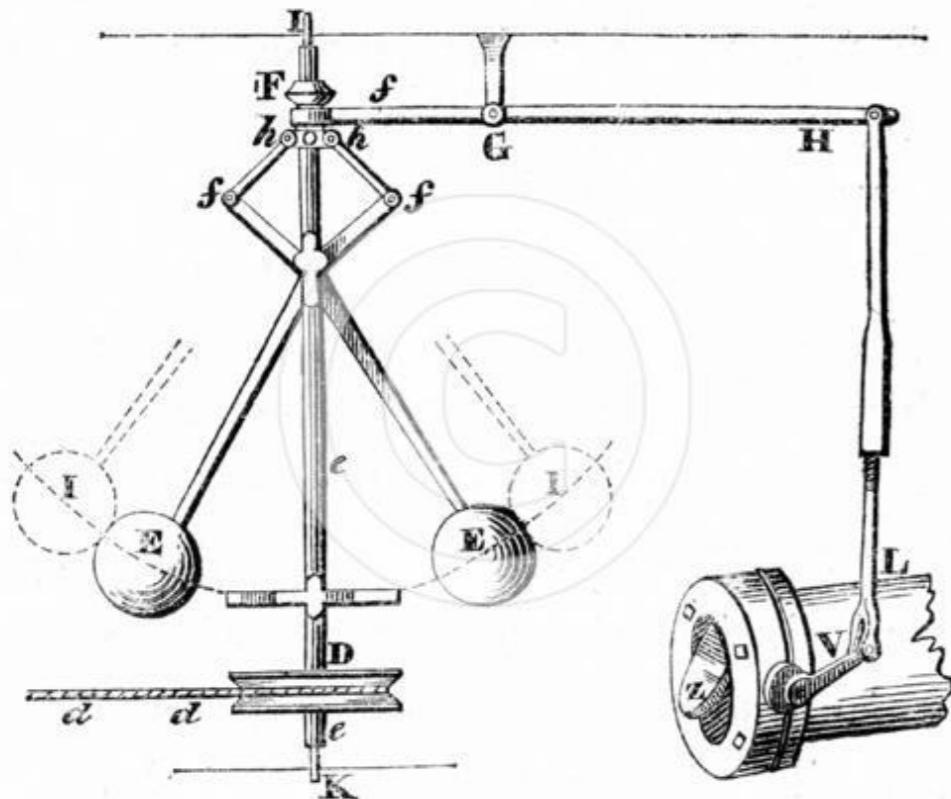
and also, at least in the first stage, to assume $\sigma(x) \approx \sigma$

$$\begin{aligned} x_{n+1} &= x_n + hF(x_n) + h\sigma \tilde{R}_n + \sqrt{2h} R_n \\ &= x_n + hF(x_n) + \sqrt{h} \sqrt{h\sigma^2 + 2} \hat{R}_n \end{aligned}$$

Like Euler-Maruyama discretization of

$$dx = F(x)dt + \sqrt{2 + \sigma h} dW$$

James Watt's Engine



Too fast: balls move to outside, opening valve, releasing steam, reducing pressure, reducing speed

Too slow: balls fall to inside, closing valve, leading to an increase in pressure, increasing speed

Nose-Hoover dynamics - a “Gibbs Governor”

$$\dot{q} = p$$

$$\dot{p} = -\nabla U(q) - \xi p$$

$$\dot{\xi} = p^2 - kT$$

Preserves

$$e^{-\beta[p^2/2+U(q)]} \times e^{-\beta\xi^2/2}$$

Adaptive Langevin, later “SGNHT”

Jones & L. J. Chem. Phys. 2011

Applying Nosé-Hoover Dynamics to a system with noisy gradient corrects the canonical distribution.

Adaptive (Automatic) Langevin

$$dx = M^{-1} pdt$$

$$dp = -\nabla U dt + \sqrt{h}\sigma dW - \xi pdt + \sigma_A dW_A$$

$$d\xi = \mu^{-1} [p^T M^{-1} p - n\beta^{-1}] dt$$

$$\tilde{\rho} = e^{-\beta[p^T M^{-1} p/2 + U(x)]} \times e^{-\beta\mu(\xi - \gamma)^2/2}$$

Shift in auxiliary variable by $\gamma = \frac{\beta(h\sigma^2 + \sigma_A^2)}{2\text{Tr}(M)}$

Discretization

[L. & Shang, SISC, 2016]

generator: $\mathcal{L} = \mathcal{L}_A + \mathcal{L}_B + \mathcal{L}_O + \mathcal{L}_D$

$$\mathcal{L}_A = (M^{-1}p) \cdot \nabla_x$$

$$\mathcal{L}_B = -\nabla U(x) \cdot \nabla_p + \frac{h\sigma^2}{2} \Delta_p$$

$$\mathcal{L}_O = -\xi p \cdot \nabla_p + \frac{\sigma_A^2}{2} \Delta_p$$

$$\mathcal{L}_D = G(p) \frac{\partial}{\partial \xi}$$

define related operator by composition, e.g. **BADODAB**

$$e^{h\hat{\mathcal{L}}} = e^{\frac{h}{2}\mathcal{L}_B} e^{\frac{h}{2}\mathcal{L}_A} e^{\frac{h}{2}\mathcal{L}_D} e^{h\mathcal{L}_O} e^{\frac{h}{2}\mathcal{L}_D} e^{\frac{h}{2}\mathcal{L}_A} e^{\frac{h}{2}\mathcal{L}_B}$$

Mimic properties of BAOAB in the high σ_A regime!

Superconvergence

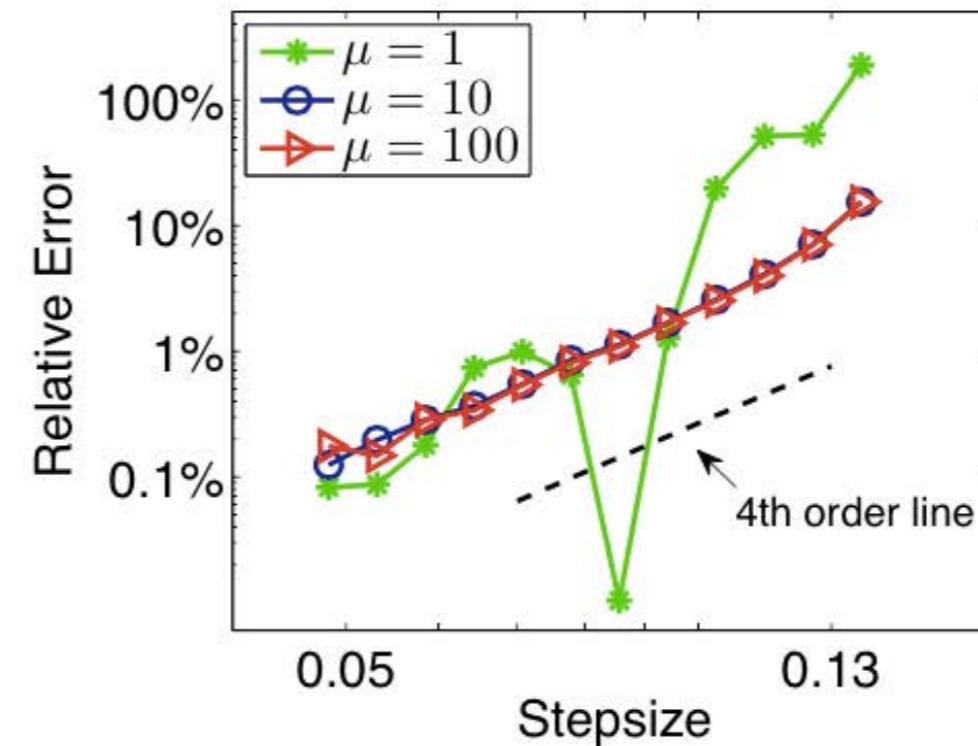
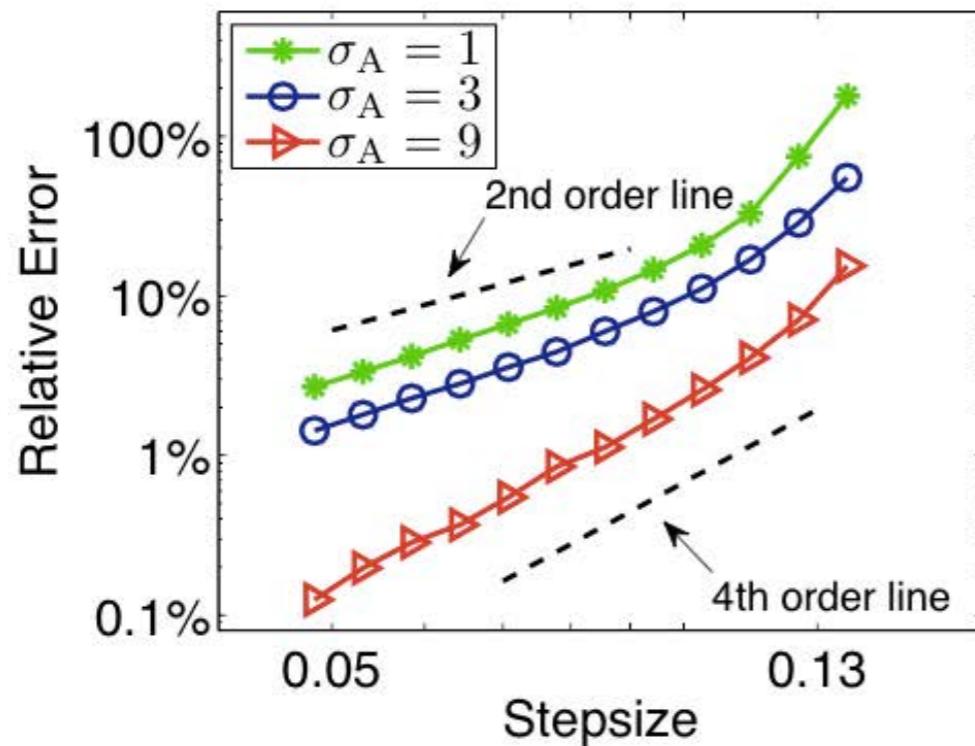
BAOAB, in the high friction limit, gives a ***superconvergence*** property for configurational quantities.

By taking large $\gamma \propto \sigma_A^2$ and $\mu \propto \sigma_A^2$ we can make BADODAB behave like BAOAB in the high friction limit after averaging over the auxiliary variable.

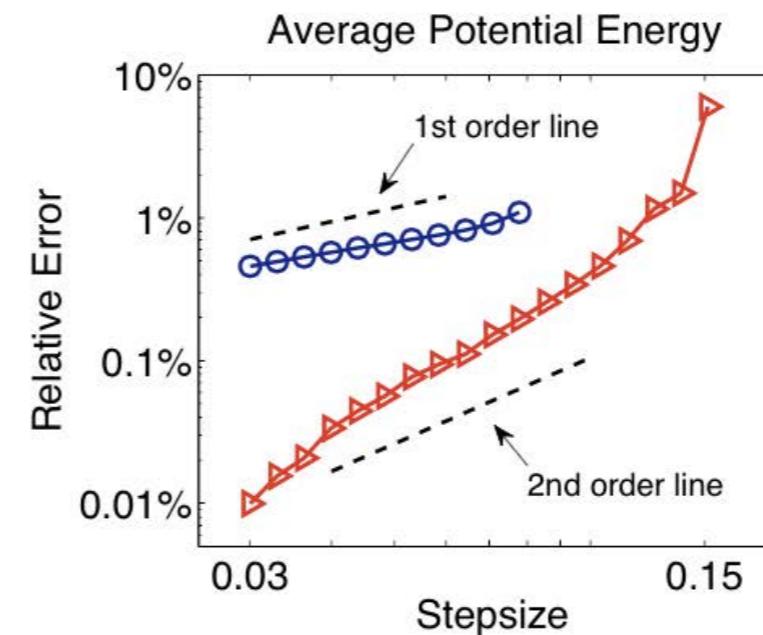
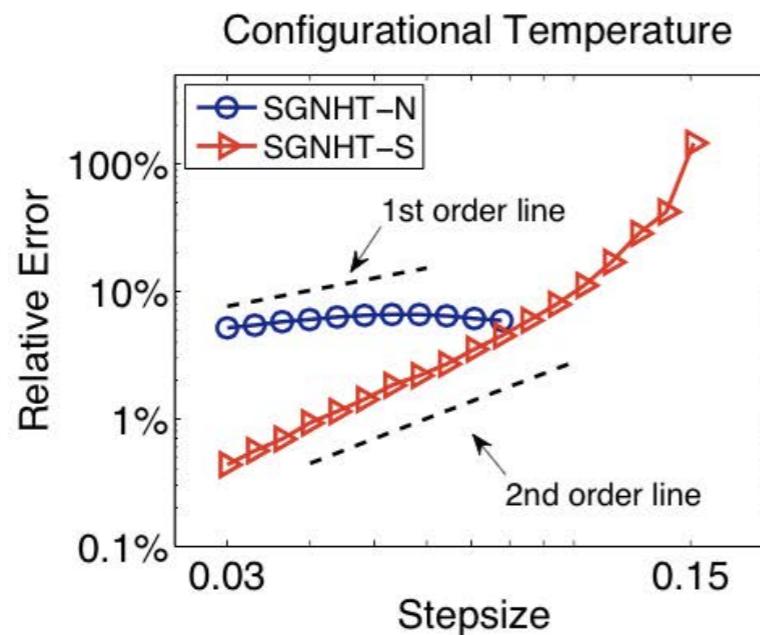
Effectively the extra driving noise implements a projection to the case of Langevin dynamics, ***but large driving noise means large effective friction, restricts phase space exploration.***

500 LJ particles, clean gradient

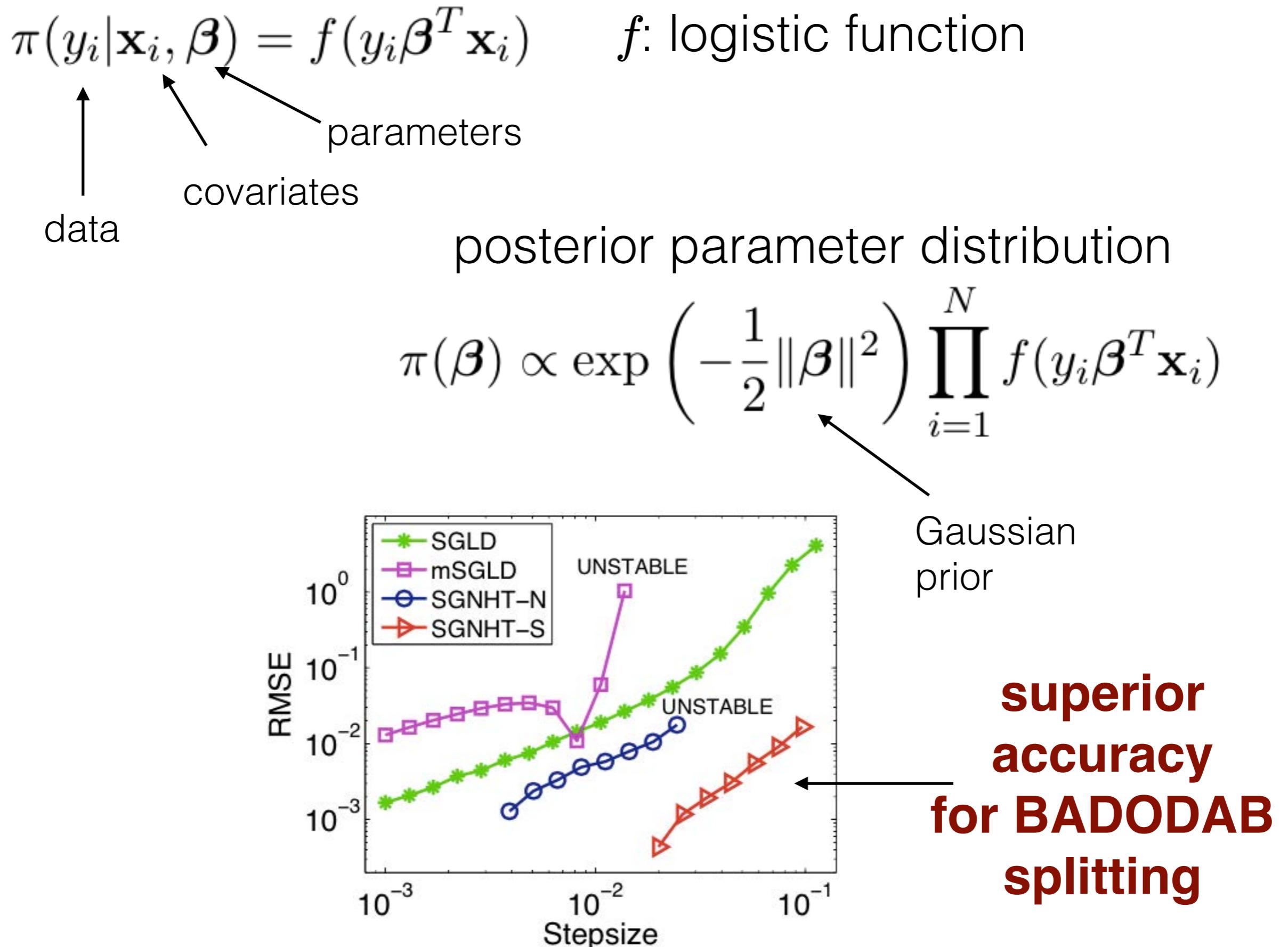
configurational temperature



Comparison with Chen et al. (SGNHT, NIPS 2015)



Bayesian Logistic Regression



$L^2(\mu)$ Convergence Analysis

[w./Matthias Sachs & Gabriel Stoltz, 2018]

HYPOCOERCIVITY FOR LINEAR KINETIC EQUATIONS CONSERVING MASS,
J. Dolbeault, C. Mouhot, & C. Schmeiser, Trans. AMS 367, 3807-3828, 2015

$$\partial_t f + \mathbf{T} f = \mathbf{L} f$$

\mathbf{T} : anti-symmetric; \mathbf{L} symmetric part of operator in $L^2(\mu)$

Π : projection operator based on p -marginalization

Define an entropy functional,

$$H[f] := \frac{1}{2} \|f\|^2 + \varepsilon \langle A f, f \rangle, \quad \text{with } A := (1 + (\mathbf{T} \Pi)^*(\mathbf{T} \Pi))^{-1} (\mathbf{T} \Pi)^*$$

and show it is decreasing along the flow: $\frac{d}{dt} H[f] \leq -\frac{2\kappa}{1+\varepsilon} H[f],$

Consequence: $L^2(\mu)$ hypocoercivity + spectral gap

DMS Convergence Analysis of AdL

$$dq_t = p_t dt,$$

$$dp_t = \left(-\nabla U(q_t) - \epsilon^{-1} \xi p_t - \gamma p_t \right) dt + \sqrt{\frac{2\gamma}{\beta}} dW_t,$$

$$d\xi_t = \epsilon^{-1} \left(|p|^2 - n\beta^{-1} \right) dt,$$

Assume a Poincaré inequality of the form

$$\left\| \varphi - \int_{\Omega} \varphi d\mu_q \right\|_{L^2(\mu_q)}^2 \leq \frac{1}{C} \|\nabla \varphi\|_{L^2(\mu_q)}^2, \quad \varphi \in H^1(\mu_q)$$

Then

$$\forall t \geq 0, \quad \left\| e^{t\mathcal{L}_{\gamma,\epsilon}} \left(\varphi - \int \varphi d\mu \right) \right\|_{L^2(\mu)} \leq K_{\epsilon,\gamma} e^{-\kappa_{\epsilon,\gamma} t} \left\| \varphi - \int \varphi d\mu \right\|_{L^2(\mu)}$$

with

$$\kappa_{\epsilon,\gamma} \in O \left(\min(1, \epsilon^{-2}) \cdot \min(\gamma, \gamma^{-1}) \right)$$

CLT for AdL

Bhattacharya et al 1982: “On the functional central limit theorem...”

$$\text{CLT: } \sqrt{t} (\mu_t[\varphi] - \mu[\varphi]) \sim \mathcal{N}(0, \sigma_\varphi^2), \text{ as } t \rightarrow \infty.$$

Take Home Message #2.

Adaptive Langevin dynamics offers a simple method for removing bias in noisy gradient simulations while maintaining high accuracy (low bias).

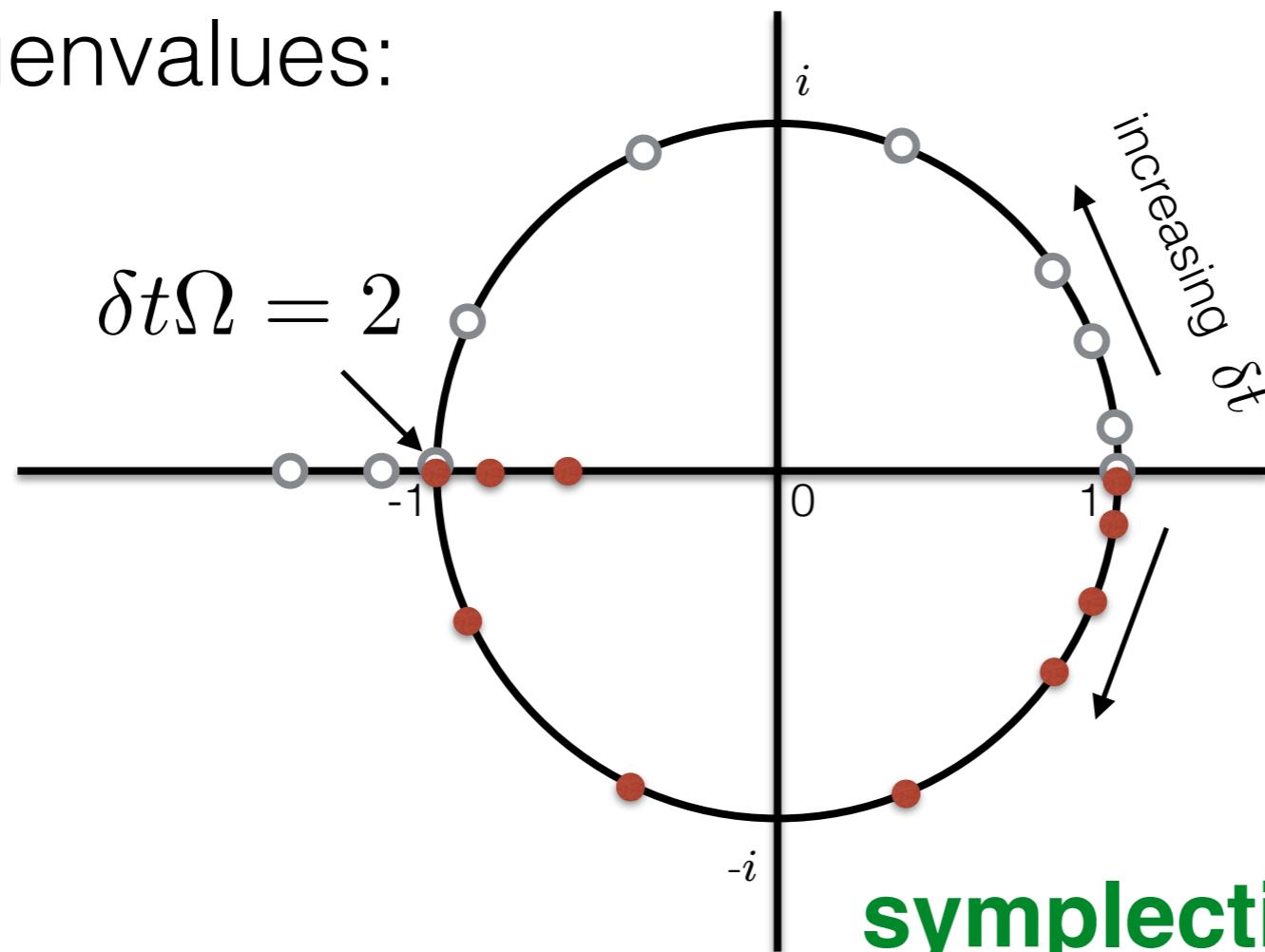
2.2 The problem of the timestep

Harmonic Oscillator

$$\dot{q} = p$$

$$\dot{p} = -\Omega^2 q$$

eigenvalues:



leapfrog/Verlet

$$q_{n+1} = q_n + \delta t p_{n+1/2}$$

$$p_{n+1/2} = p_n - \frac{\delta t}{2} \Omega^2 q_n$$

$$p_{n+1} = p_{n+1/2} - \frac{\delta t}{2} \Omega^2 q_{n+1}$$

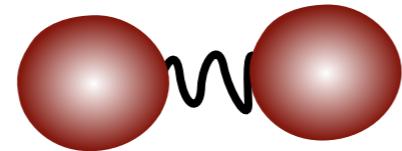
stable for
 $\delta t \Omega \leq 2$

symplectic for all δt
but only **useful** for small δt

Stepsize Restriction

$$\text{H.O. } H = \frac{p^2}{2} + \frac{\Omega^2 q^2}{2}$$

think bonds



stability threshold

$$\delta t < 2/\Omega$$

in molecular dynamics
bonds to H atoms have
high frequencies

$$\delta t < 3\text{fs} = 3 \times 10^{-15}\text{s}$$

The Timestep Problem in MD

Find timestepping methods that allow

$$\delta t \gg \delta t_{\text{Verlet}}$$

Three possibilities:

~~Use an **implicit method**~~

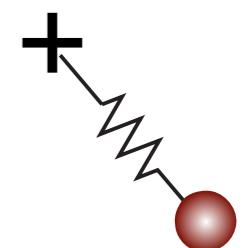
Eliminate the motion of fast components
constraints

Isolate the stiff terms for efficient treatment
multiple timestepping

Stiff Spring Oscillator

energy:

$$H = \frac{|p|^2}{2m} + \frac{k(|q| - L)^2}{2} = E$$

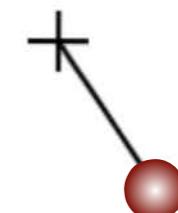


equations of motion:

$$\frac{dq}{dt} = p/m$$

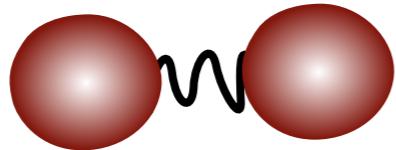
$$\frac{dp}{dt} = -k \left(1 - \frac{L}{|q|}\right) q$$

If energy is fixed, then in the limit of large k
we must have $|q| \sim L$

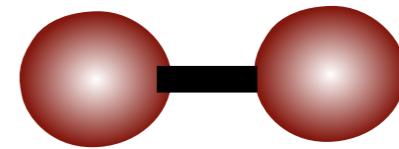


Constrained dynamics: $|q| = L$

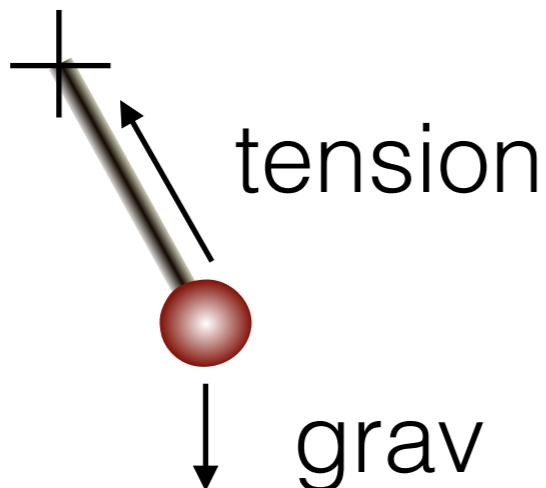
Pendulum



stiff harmonic
spring with rest length



rigid rod
(holonomic constraint)



$$\dot{q} = m^{-1} p$$

$$\dot{p} = - \begin{bmatrix} 0 \\ g \end{bmatrix} - q\lambda$$

$$0 = q \cdot q - 1$$

Constrained Hamiltonian Systems

$$\frac{d}{dt}q = \mathbf{M}^{-1}p$$

$$\frac{d}{dt}p = F - \sum_{i=1}^m \lambda_i \nabla g_i(q)$$

$$0 = g_j(q), \quad j = 1, 2, \dots, m$$

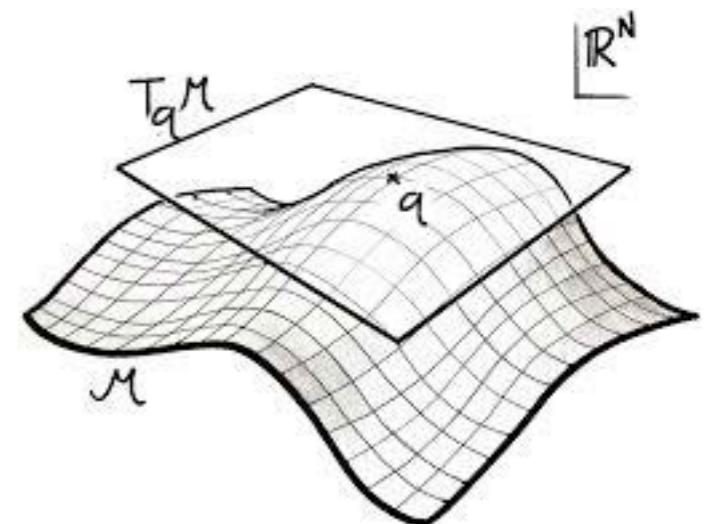
$$0 = \nabla g_j(q)^T \mathbf{M}^{-1} p, \quad j = 1, 2, \dots, m$$

“hidden constraints”

“cotangent bundle”

$$T^*\mathcal{M} = \{(q, p) | g_j(q) = 0, \nabla g_j(q)^T \mathbf{M}^{-1} p = 0, j = 1, 2, \dots, m\}$$

symplectic: dynamics preserves $[dq \wedge dp]_{T^*\mathcal{M}}$



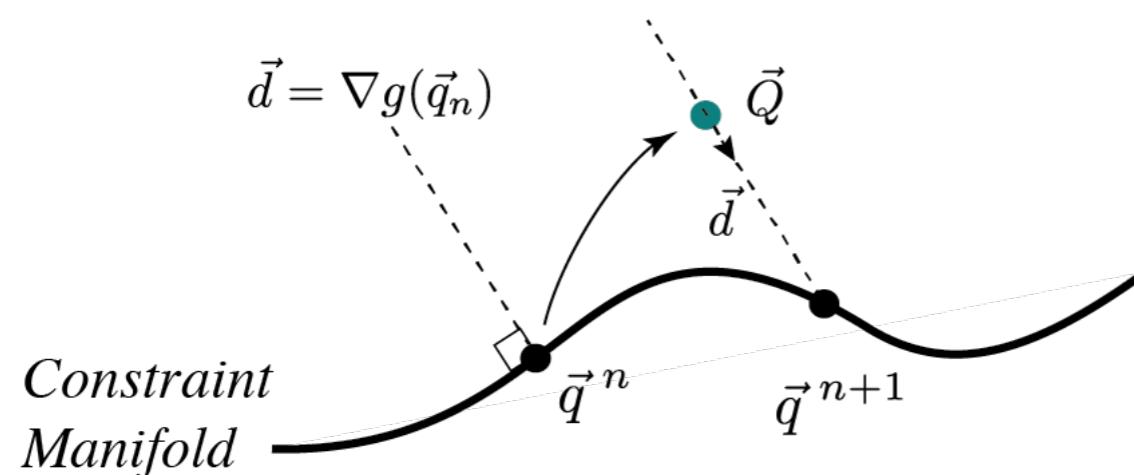
SHAKE

[Berendsen, Ciccotti, Ryckaert 1977]

RATTLE

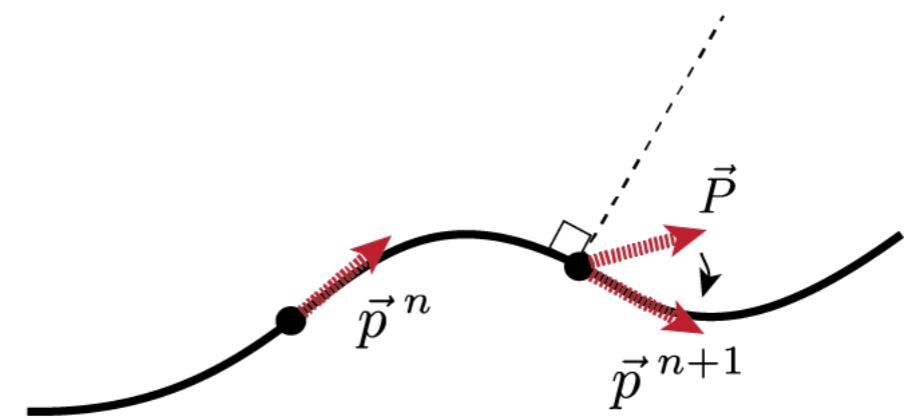
[Andersen 1983]

“SHAKE” Projection



$$\mathcal{M} = \{\vec{q} \mid g(\vec{q}) = 0\}$$

“RATTLE” Projection



L. & Skeel, J.Comput. Phys, 1994

SHAKE and **RATTLE** are **symplectic** methods
and are actually the same method (**conjugate** methods).

Constraint preserving methods

$$H(q, p) = |p|^2/2 + U(q)$$

$$\begin{aligned}\dot{q} &= p \\ \dot{p} &= -\nabla_q U(q) - \lambda \nabla g \\ 0 &= g(q) \\ 0 &= \nabla g(q) \cdot p\end{aligned}$$

$$(q, p) \mapsto (Q, P)$$

approximating a
step in time of size h

$$Q := q + h\tilde{P}$$

$$\tilde{P} := p - h\nabla U(q) - \lambda \nabla g(q)$$

$$0 = g(Q)$$

$$P := \tilde{P} - \mu \nabla g(Q)$$

$$0 = \nabla g(Q) \cdot P$$

first order

$$Q := q + (h/2)\tilde{P}$$

$$\tilde{P} := p - (h/2)\nabla U(q) - \lambda \nabla g(q)$$

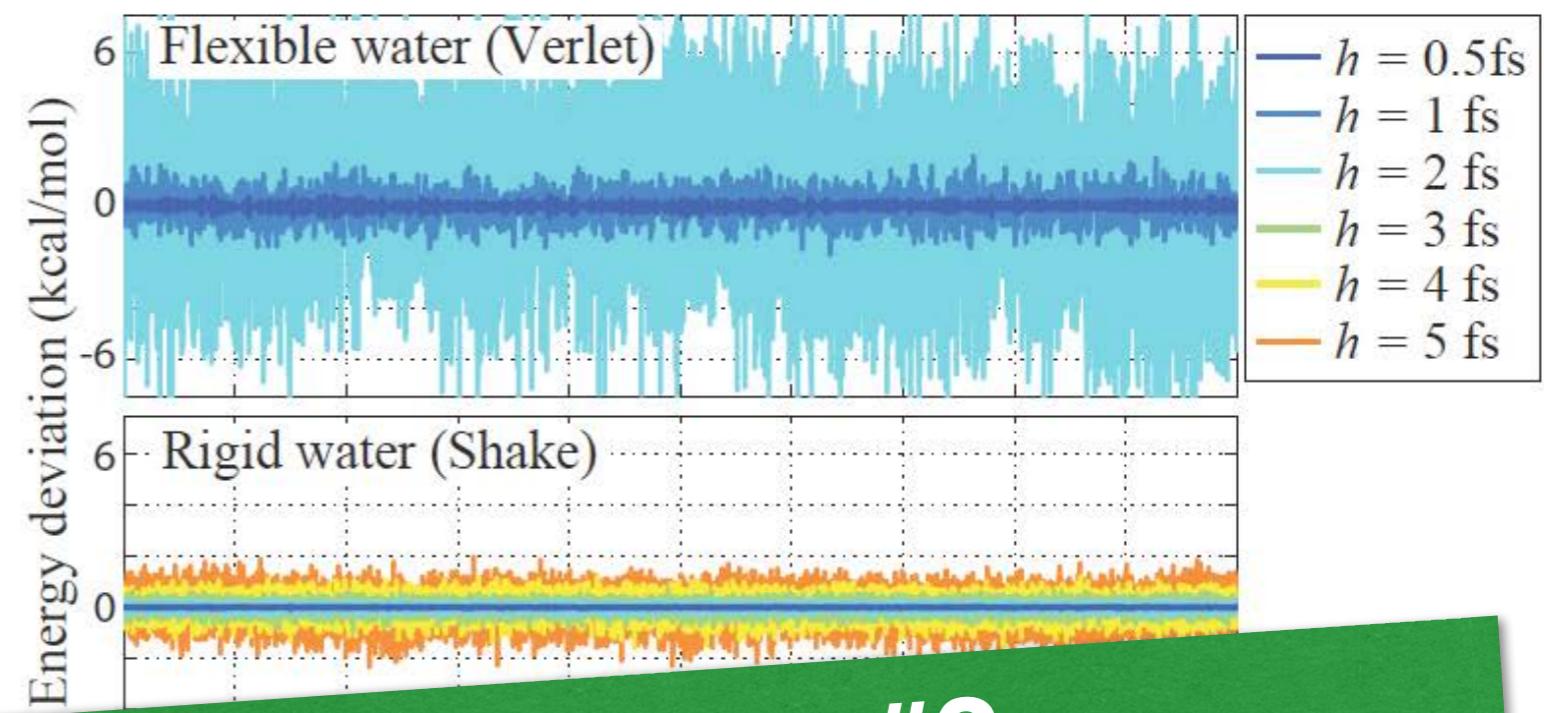
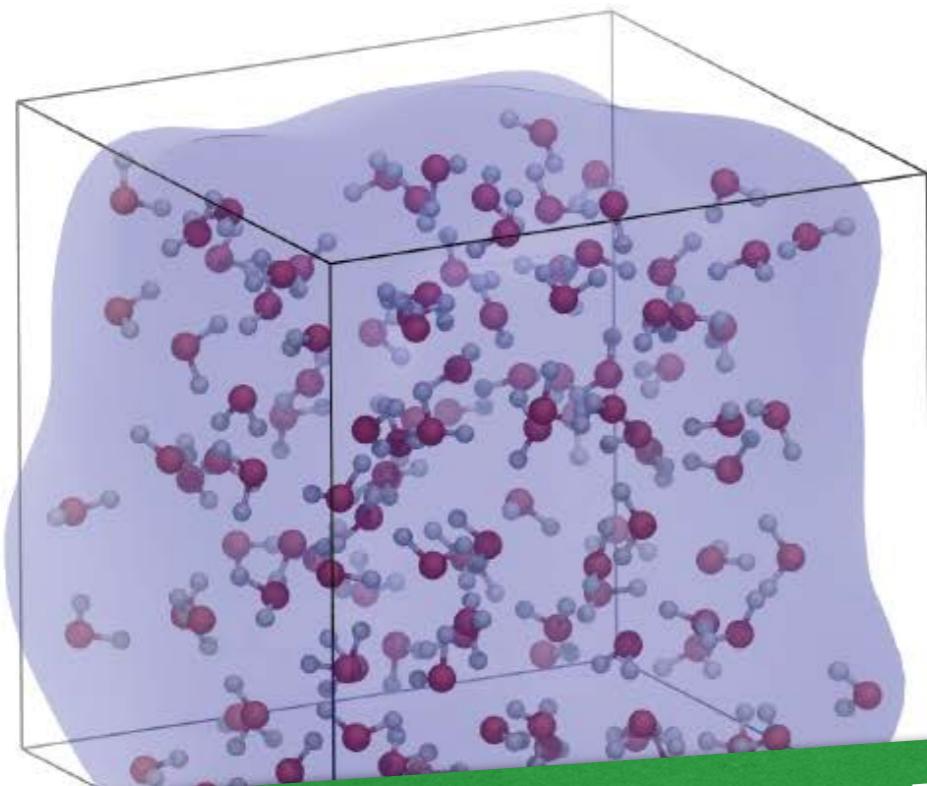
$$0 = g(Q)$$

$$P := \tilde{P} - (h/2)\nabla U(Q) - \mu \nabla g(Q)$$

$$0 = \nabla g(Q) \cdot P$$

2nd order “RATTLE”

Microcanonical Sampling



Take Home Message #3
Constraints don't buy much stability in
deterministic MD, but they do improve
accuracy!
1 fs with 1 kcal/mol energy accuracy

SHAKE (RATTLE) OK up to as much as
3-4fs with energy accuracy

Multiple timestepping

Multiple Timestepping (RESPA)

Tuckerman, Martyna, Berne '90

also Grubmüller, Heller, Windemuth, Schulten '91

Some forces are **fast-changing** (e.g. short-ranged)
but **cheap** to compute

Some forces are **slow-changing** (e.g. long-ranged)
and **costly** to compute

$$\begin{aligned} & \exp\left(\frac{\delta t}{2}\mathcal{L}_{U^{\text{slow}}}\right) \times \\ & \left[\exp\left(\frac{\delta t}{2r}\mathcal{L}_{U^{\text{fast}}}\right) \exp\left(\frac{\delta t}{r}\mathcal{L}_K\right) \exp\left(\frac{\delta t}{2r}\mathcal{L}_{U^{\text{fast}}}\right) \right]^r \\ & \times \exp\left(\frac{\delta t}{2}\mathcal{L}_{U^{\text{slow}}}\right) \end{aligned}$$

Many fast cheap evaluations
Few slow costly ones

Linear Model Problem

model problem

$$H = \frac{p^2}{2} + \frac{(1 + \Omega^2)q^2}{2} \quad \Omega \gg 1$$

idealized *multiple timestepping (fast solve is exact)*

$$H(q, p) = p^2/2 + U_{\text{S}}(q) + U_{\text{F}}(q)$$

$$U_{\text{S}}(q) = q^2/2$$

$$U_{\text{F}}(q) = \Omega^2 q^2/2$$

Kick with U_{S}

$$\delta t/2$$

Solve $H_{\text{F}} = p^2/2 + U_{\text{F}}(q)$

$$\delta t$$

Kick with U_{S}

$$\delta t/2$$

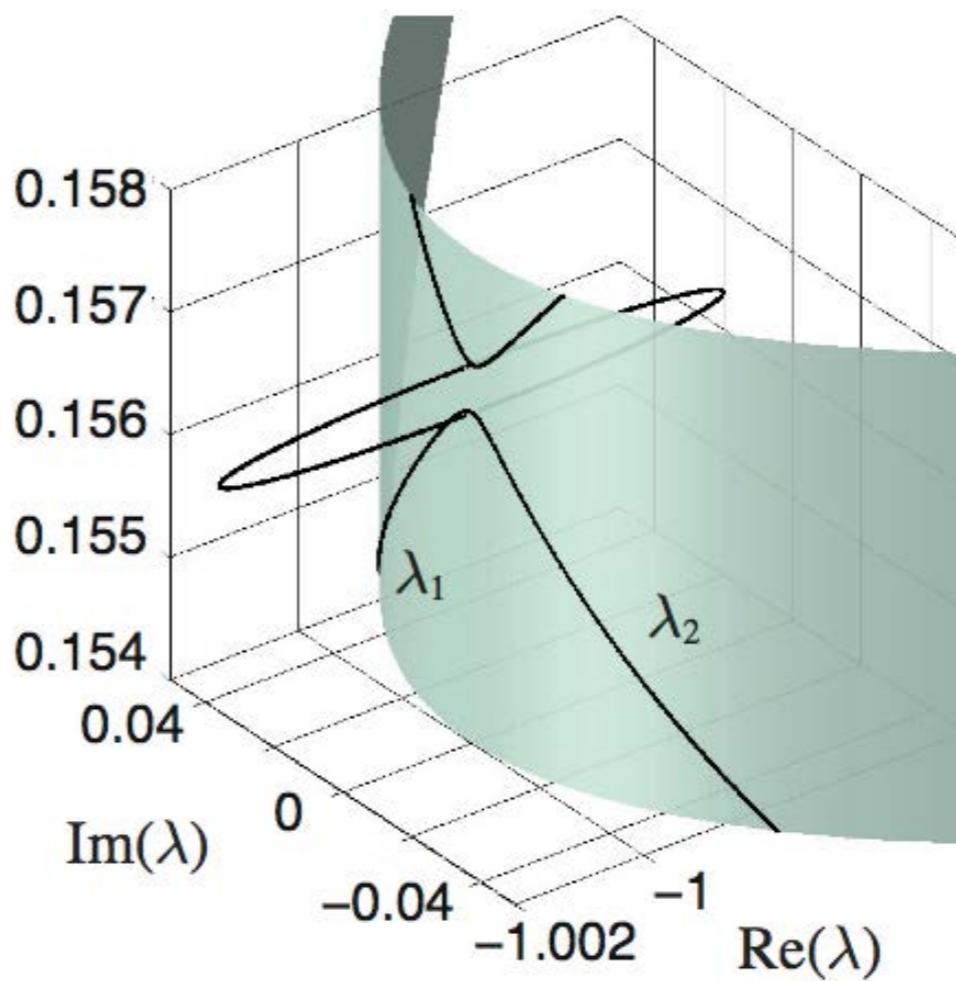
Resonance

Eigenvalues of timestep map:

$$\lambda_1 \lambda_2 = 1; \quad \lambda_1 + \lambda_2 = 2 \cos(\delta t \Omega) - \frac{\delta}{\Omega} \sin(\delta t \Omega)$$
$$\varepsilon \propto \delta t - \pi/\Omega < 0 \quad \lambda = -1 \pm (|\varepsilon|/2)^{1/2} + O(\varepsilon)$$

instabilities...

δt



Resonance

For the harmonic model, Verlet introduces a stability restriction

$$\delta t < 2/\Omega$$

Multiple timestepping, in the idealized form described here has a stability restriction of about

$$\delta t < \pi/\Omega$$

i.e., not very dramatic improvement. In practice, we are limited to around **3.5fs**

Mollified Impulse Method

Garcia-Archilla, Sanz-Serna and Skeel 1998

To stabilize multiple timestepping, one idea is to average out over an associated fast dynamics to ``mollify'' the impulse in RESPA

starting from e.g $(q, 0)$, solve

$$H^{\text{fast}}(\mathbf{q}, \mathbf{p}) = \mathbf{p}^T \mathbf{M}^{-1} \mathbf{p}/2 + U_F(\mathbf{q})$$

to produce a sequence $\theta_0 = q, \theta_1, \theta_2, \dots, \theta_K$
then define:

$$\mathcal{A}(\mathbf{q}) = \frac{1}{K+1} \sum_{i=0}^K \phi_i \theta_i(\mathbf{q})$$

Advance the step using a *mollified impulse*:

$$\tilde{U}(q) := U_s(\mathcal{A}(q))$$

Pushing the limits

SHAKE-MOLLY [*Izaguirre, Reich, Skeel 1999*]

constrains the slow force evaluations to lie exactly at the bond stretch minima (multiple timestepping + constraints).

With the very best **deterministic** schemes, however, we find for a biological molecule:

Take Home Message #4
Combining constraints and multiple
timestepping can improve accuracy and
stability

Stochastic Methods

Introducing random perturbations might seem to complicate the numerical integrator.

In fact, if done correctly, it is possible to gain in two ways:

- 1) a significant **stability improvement** is possible
- 2) substantially **higher accuracy** is possible by shifting emphasis to the **invariant distribution**

These benefits carry over to **multiple timestepping** and **constraints**, and **the combination can be even more powerful than expected!!**

Constrained Langevin Dynamics

Constrained Langevin Dynamics

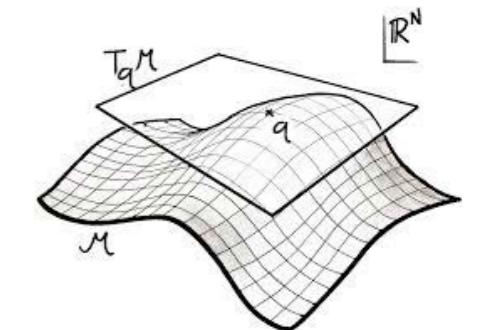


$$\frac{d}{dt}q = \mathbf{M}^{-1}p$$

$$\frac{d}{dt}p = F - \gamma p + \sqrt{2k_B T \gamma} \mathbf{M}^{1/2} \eta(t) - \sum_{i=1}^m \lambda_i \nabla g_i(q)$$

$$0 = g_j(q), \quad j = 1, 2, \dots, m$$

$$0 = \nabla g_j(q)^T \mathbf{M}^{-1} p, \quad j = 1, 2, \dots, m$$



Lagrange multipliers

$$T^* \mathcal{M} = \{(q, p) | g_j(q) = 0, \nabla g_j(q)^T \mathbf{M}^{-1} p = 0, j = 1, 2, \dots, m\}$$

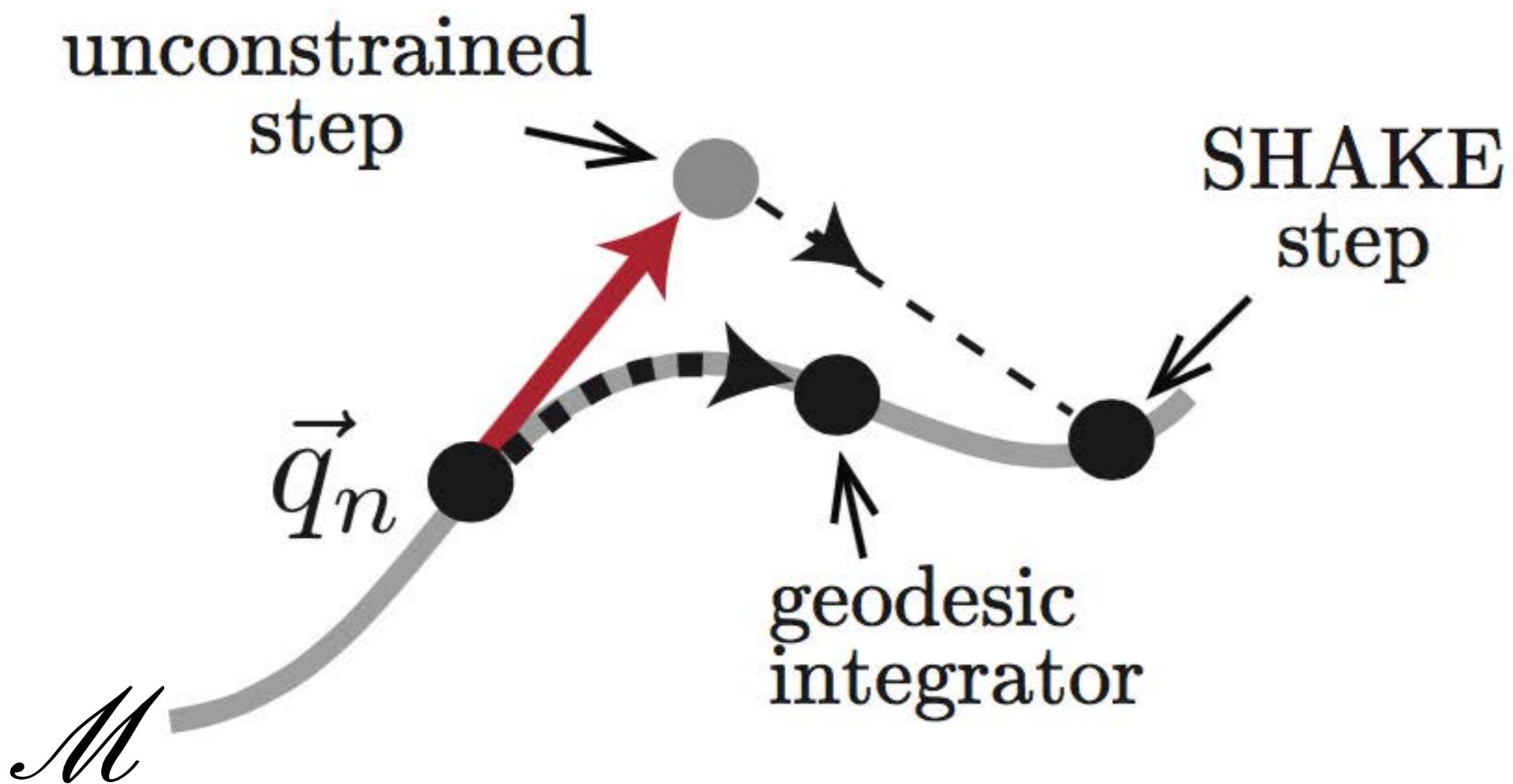
Some technical aspects:

Lelievre, Rousset and Stoltz, Math. Comp., 2010

Geodesic Integrator

B.L. and G. Patrick, J. Nonlin. Sci. 1996 (deterministic)

B.L. and C. Matthews, Proc Roy Soc A 2016 (stochastic)



Geodesic Integrator

B.L. and G. Patrick, J. Nonlin. Sci. 1996 (deterministic)

B.L. and C. Matthews, Proc Roy Soc A 2016 (stochastic)

An alternative to SHAKE discretization

Idea: preserve the configuration manifold during position moves and the cotangent space during impulse.

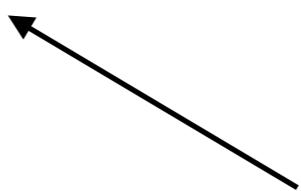
The natural constrained analogue of **Verlet (BAB)** or **BAOAB**

$$g\text{-Verlet: } \Phi_{h,H_{T^*\mathcal{M}}} \approx \Phi_{h/2,U_{T^*\mathcal{M}}} \circ \Phi_{h,T_{T^*\mathcal{M}}} \circ \Phi_{h/2,U_{T^*\mathcal{M}}}$$

Combines: **geodesic flow projected “kicks”
projective OU**

Constrained OU

$$dp = -\gamma \Pi p dt + \sqrt{2\gamma k_B T} \Pi M^{1/2} dW$$



Projector onto the co-tangent space

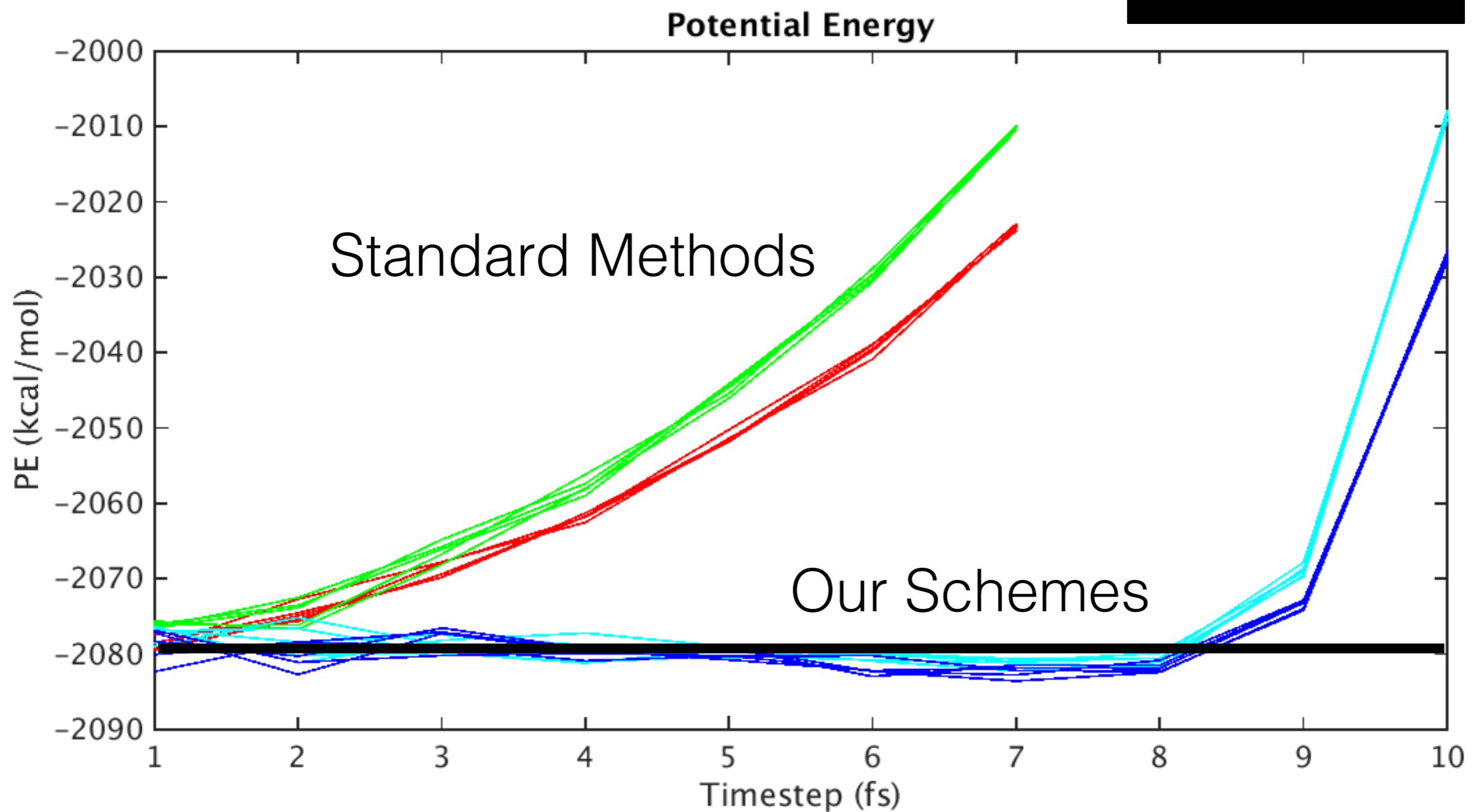
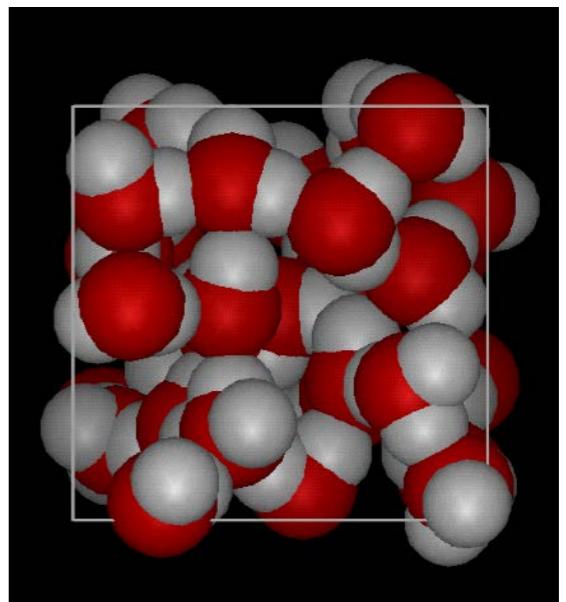
Easily solved using Rodrigues' formula:

$$p(t) = e^{-\gamma t} p(0) + \sqrt{k_B T (1 - e^{-2\gamma t})} \Pi M^{1/2} R(t).$$

assuming $\Pi p(0) = p(0)$

Box of H₂O

For rigid body water, implement the geodesic integrator using
SETTLE: *exact geodesic steps*



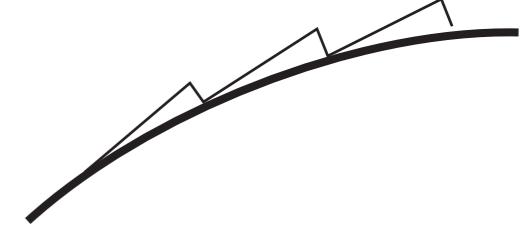
Multiple-Timestepping Implementation

To implement the geodesic integrator, we use
a sequence of SHAKE/RATTLE steps for the “**A**” step

These steps do not require re-evaluation of the force
field, so each iteration is rel. cheap (in the MD universe).

The **B** and **O** (for Langevin) steps are simple RATTLE
projections, so there is no significant added cost for these.

***No need for Hessians/normal modes,
or tuning/parameterization***



Geodesic Integrator

$$p \leftarrow p + \frac{\delta t}{2} F(q) + \sum_j \mu_j G_j(q) \in T_q^* \mathcal{M},$$

**Cotangent
space projection**

For k from 1 to K_r do:

$$(q, p) \leftarrow A \left(q, p, \frac{\delta t}{2K_r} \right)$$

**Several Rattle
(geodesic flow) Steps**

end do

$$p \leftarrow a_2 p + b_2 \mathbf{M}^{1/2} \mathbf{R} + \sum_j \mu_j G_j(q) \in T_q^* \mathcal{M},$$

O Step

For k from 1 to K_r do:

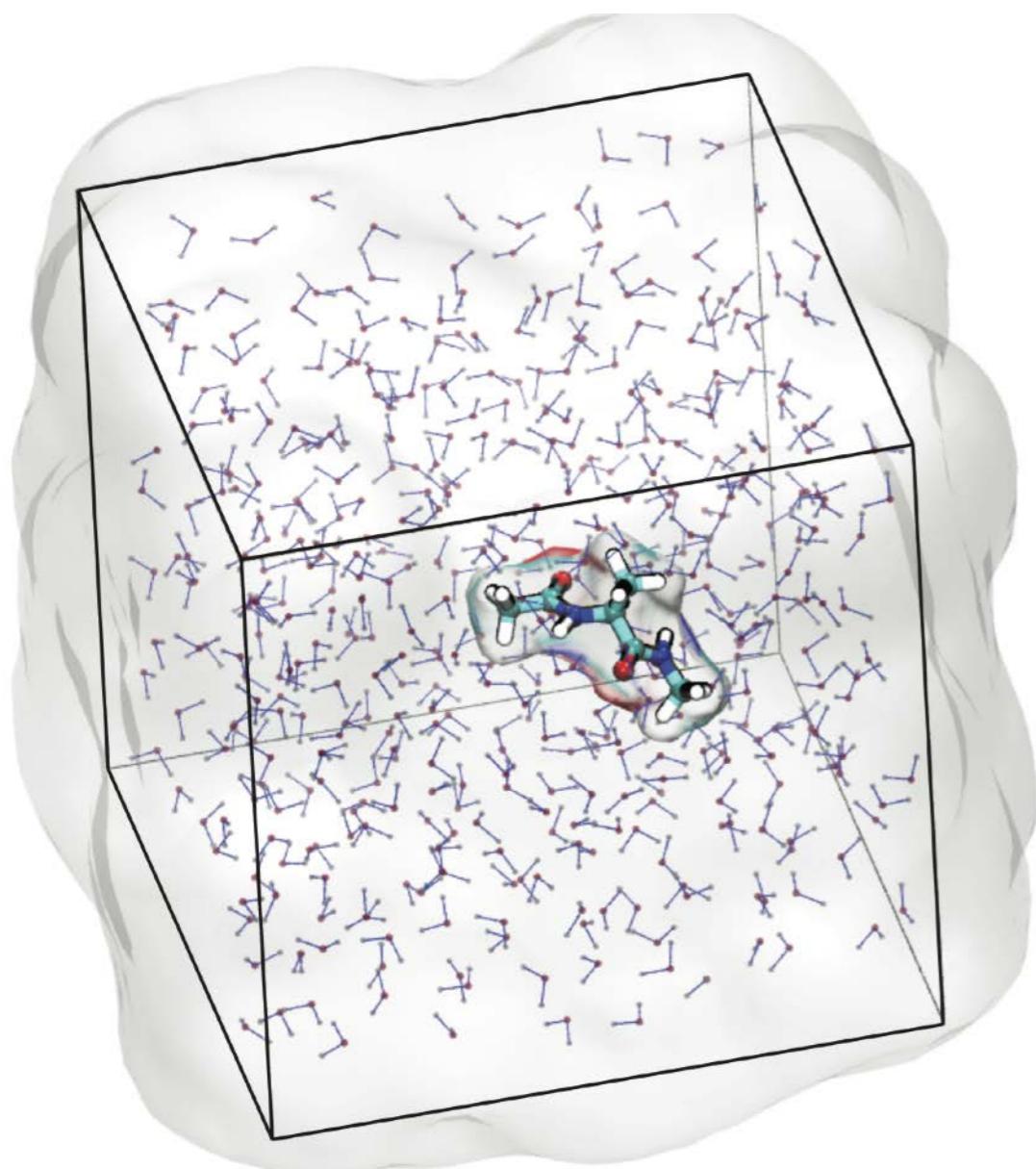
$$(q, p) \leftarrow A \left(q, p, \frac{\delta t}{2K_r} \right)$$

end do

$$p \leftarrow p + \frac{\delta t}{2} F(q) + \sum_j \mu_j \nabla g_j(q) \in T_q^* \mathcal{M},$$

Solute-Solvent Splitting

The object of **bio-MD simulation** is virtually always a protein or nucleic acid fragment + solvent (water) bath.



Once the bonds and selected (H-X-H, X-Y-H) angles of the solute and waters are constrained, **the next fastest modes are due to flexible angle bonds of the solute.**

These would limit the stepsize to around **5fs**, even with the geodesic integrator.

Solute-Solvent Splitting

The obvious solution is to break the interaction forces into three pieces, denoted **PP** (protein-protein), **PS** (protein-solvent), and **SS** (solvent-solvent)

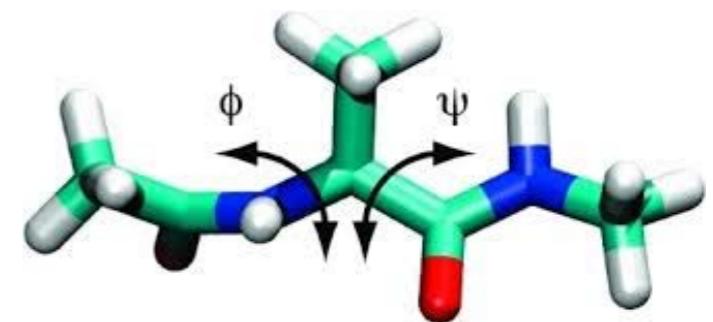
SS dominates the computational cost. **PP**, **PS** determine the stepsize. Therefore, consider
stochastic two-level multiple timestepping:

$$\exp\left(-\frac{h}{2}U_{\text{SS}}(q)\right) \boxed{\exp\left(-\frac{h}{2}[T(p) + U_{\text{PP}}(q) + U_{\text{PS}}(q)]\right)} \exp\left(-\frac{h}{2}U_{\text{SS}}(q)\right)$$

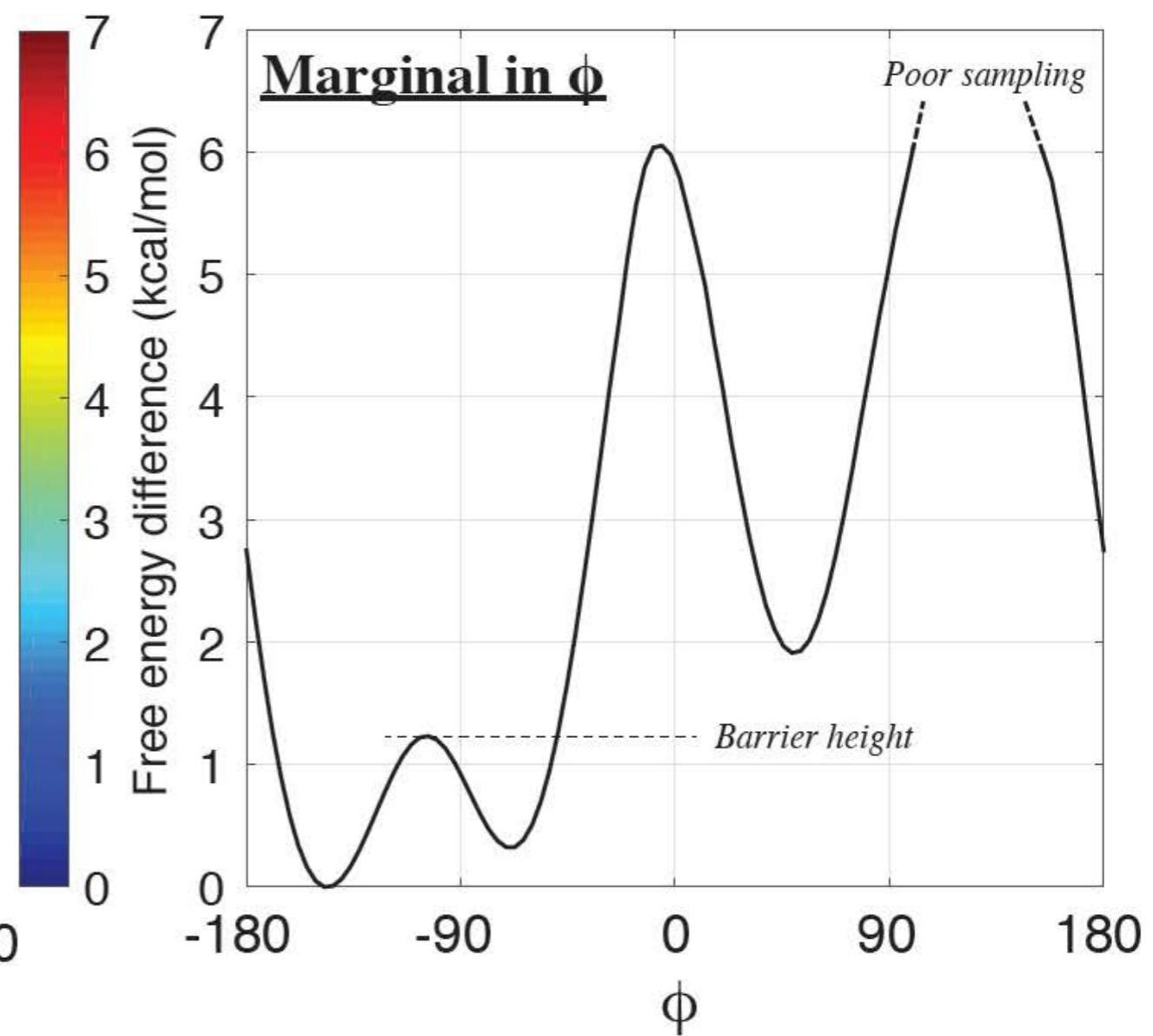
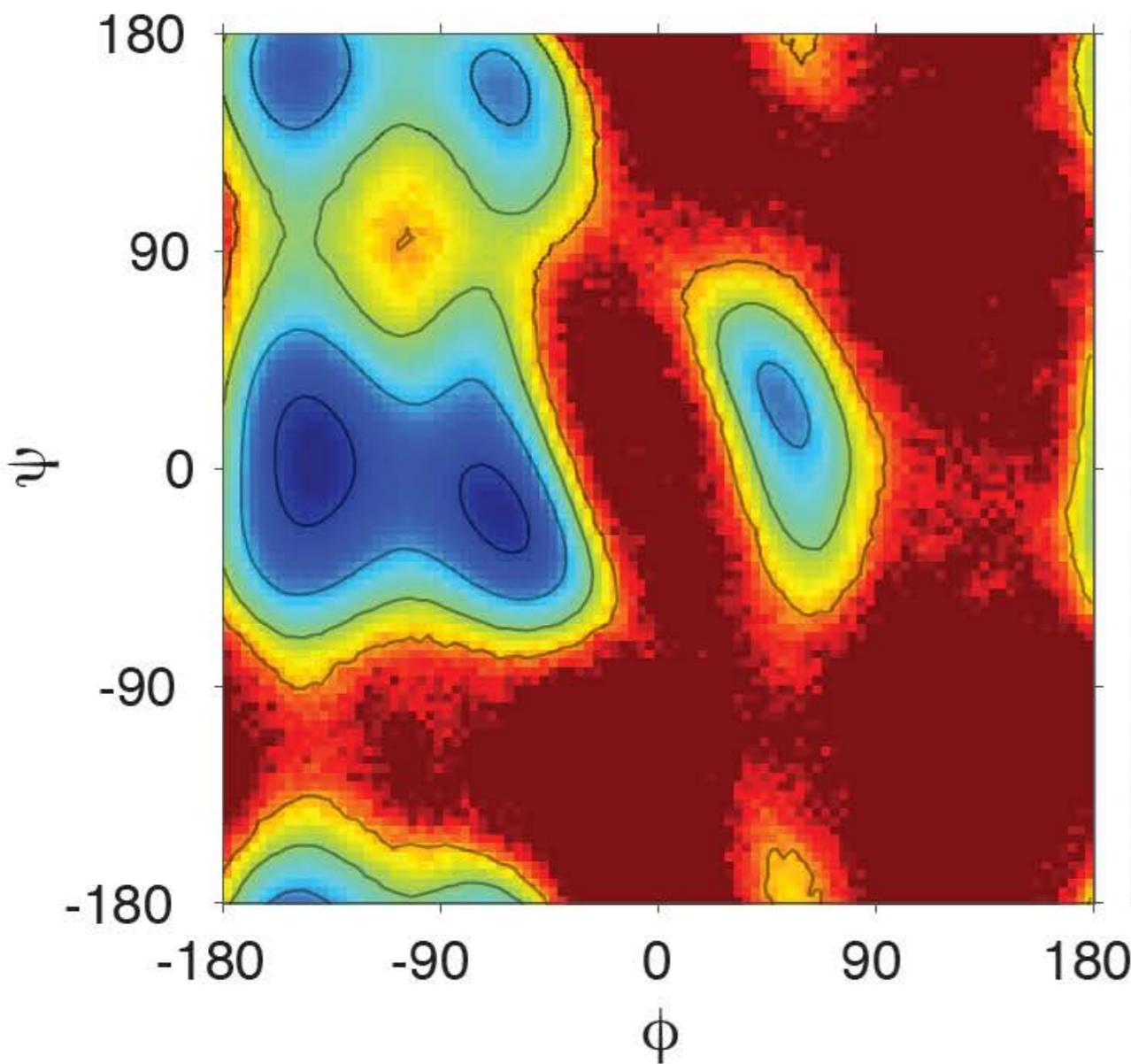
**m=2 or m=3 iterations
of a geodesic Langevin integrator
using several RATTLE substeps**

1. **PP+PS** is viewed as a many-body, stochastic system
2. Water motion can be implemented using **SETTLE**.

Alanine Dipeptide, Solvated, 300K, Tinker Code



illustrative FE surface



Comparisons

RATTLE: Standard RATTLE + Langevin

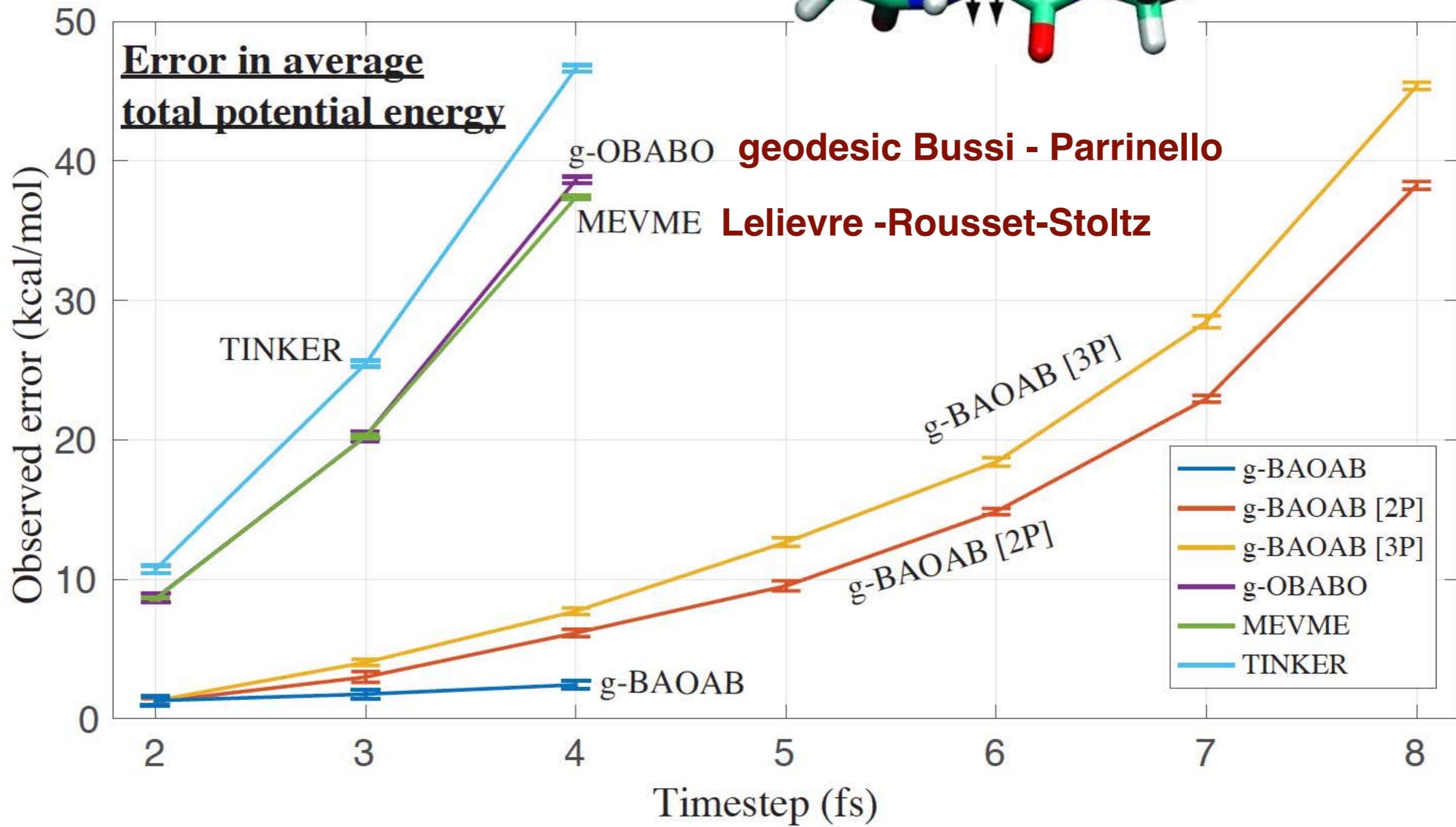
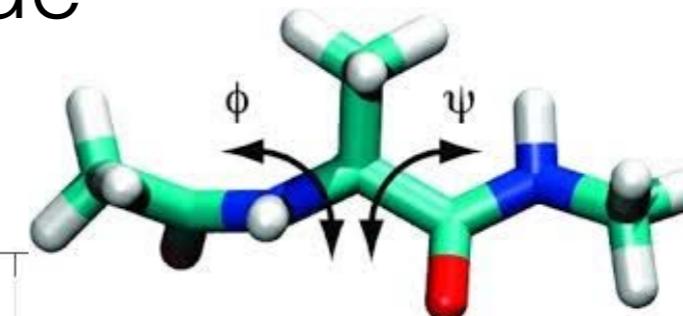
TINKER: Tinker's internal scheme

OBABO with constraint projections

MEVME: Scheme of Lelievre, Rousset, Stoltz 2010

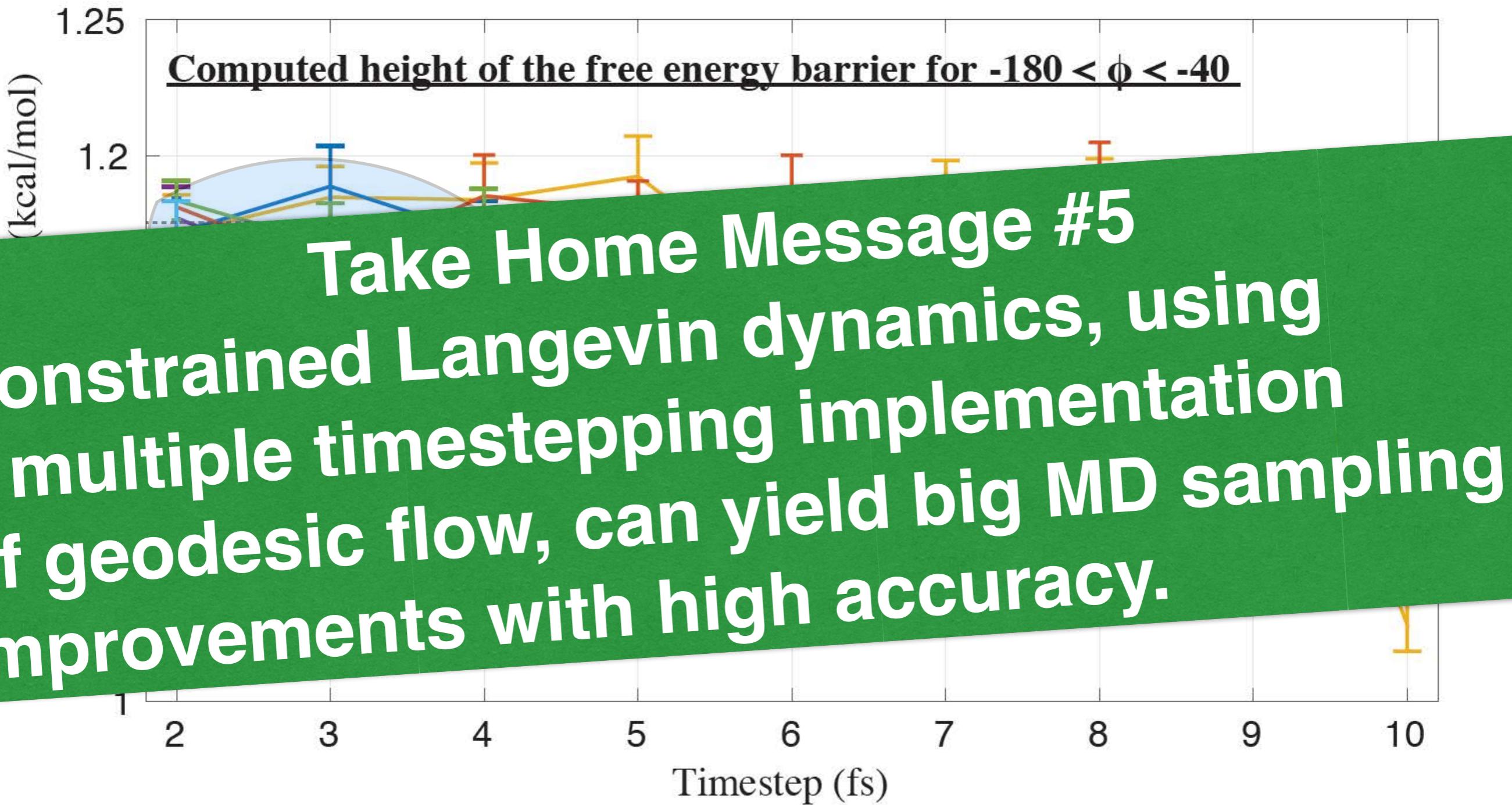
MTS-BAOAB: geodesic-MTS method, using 20 steps to compute the geodesics and m=3 (PP+PS) steps per step

Alanine Dipeptide



WHY do we get this extreme accuracy?
not explained by current BAOAB analysis

Alanine Dipeptide Effective Free Energy Barrier Height *using geodesic solvent-solute splitting*



Stochastic Backpropagation and Approximate Inference in Deep Generative Models

Danilo J. Rezende, Shakir Mohamed, Daan Wierstra

{danilor, shakir, daanw}@google.com

Google DeepMind, London

Abstract

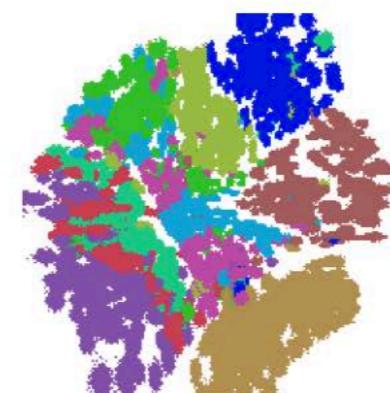
We marry ideas from deep neural networks and approximate Bayesian inference to derive a generalised class of deep, directed generative models, endowed with a new algorithm for scalable inference and learning. Our algorithm introduces a recognition model to represent an approximate posterior distribution and uses this for optimisation of a variational

Uria et al., 2014; Gregor et al., 2014) can be easily sampled from, but in most cases, efficient inference algorithms have remained elusive. These efforts, combined with the demand for accurate probabilistic inferences and fast simulation, lead us to seek generative models that are i) *deep*, since hierarchical architectures allow us to capture complex structure in the data, ii) allow for *fast sampling* of fantasy data from the inferred model, and iii) are computationally *tractable and scalable* to high-dimensional data.

$$U \mapsto Y = [\theta, X] = g(U) = [g_1(U), g_2(U)]$$

high D inputs parameters “code” features low D

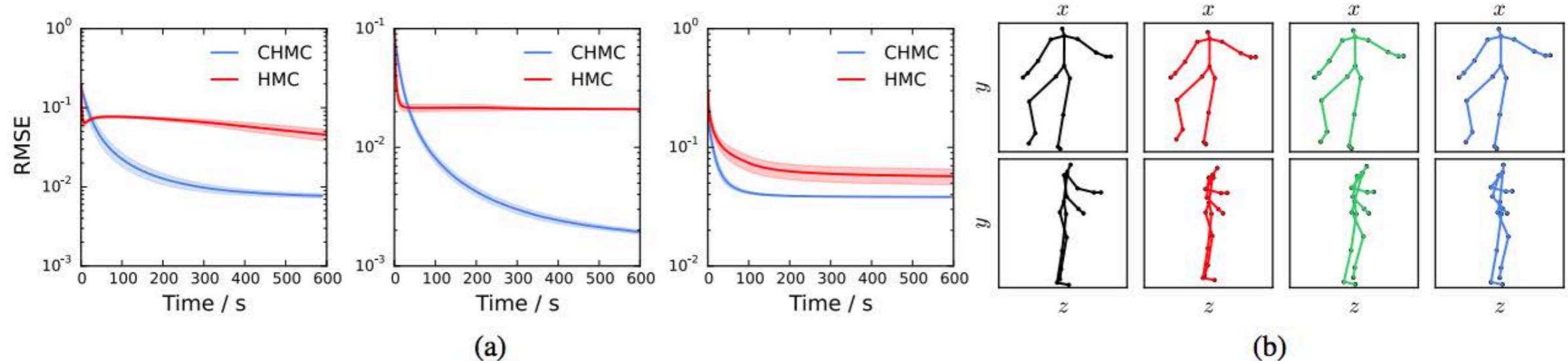
0 2 2 3 8 6 7 3 8 8	0 6 7 0 5 7 8 9 0 9	6 6 9 4 5 7 9 8 0 9
9 0 5 5 0 9 7 6 4 8	3 5 7 6 6 9 5 1 3 0	3 5 7 6 8 9 5 1 3 0
4 6 3 2 4 1 7 1 7 7	4 1 4 5 5 4 0 6 4 9	4 4 4 5 5 4 0 6 4 9
5 1 8 4 8 6 6 5 4 9	8 9 7 7 2 9 0 7 4 8	8 9 7 7 2 9 0 7 4 8
3 3 0 6 1 3 2 6 2 3	6 5 4 0 0 9 4 2 2 8	6 5 4 0 0 9 4 2 2 8
6 4 5 0 1 1 4 5 8 1	8 9 5 6 1 5 0 7 7 6	8 9 5 6 1 5 0 7 7 6
7 8 3 7 9 7 1 6 7 9	5 6 2 9 7 6 9 4 0 9	5 6 2 9 7 6 9 4 0 9
0 0 1 7 3 3 1 3 2 1	2 3 1 3 4 1 5 6 4 0	2 3 1 3 4 1 5 6 4 0
3 3 9 3 6 9 8 7 8 6	1 2 5 7 6 9 9 5 3 7	1 2 5 7 6 9 9 5 3 7
2 4 8 4 9 5 1 6 8 8	6 2 3 8 7 4 0 9 4 3	6 2 3 8 7 4 0 9 4 3



Constraints in Statistical Inference

From **A. Storkey and M. Graham, AISTATS 2017**

Ex. Human Pose Reconstruction



→ **Stochastic, constrained model within HMC**

“We use a generalisation of the RATTLE scheme to simulate the dynamic. The inner updates of the state to solve for the geodesic motion on the constraint manifold are split into multiple smaller steps. This is a special case of the scheme described in [L. & Matthews, 2016] and allows more flexibility in choosing an appropriately small step-size to ensure convergence of the iterative solution of the equations projecting on to the constraint manifold while still allowing a more efficient larger step size for updates to the momentum”

Stochastic multiple timestepping

Langevin-RESPA

L. & Mathews, Molecular Dynamics, Springer 2015

It is natural to think that stochastic perturbations may have a stabilizing influence on multiple timestepping.

Consider the combination of Langevin (stochastic) dynamics with RESPA,

$$\begin{bmatrix} dq \\ dp \end{bmatrix} = \underbrace{\begin{bmatrix} p dt \\ -\Omega^2 q dt - \gamma p dt + \sqrt{2\gamma/\beta} dW \end{bmatrix}}_{\text{Fast}} + \underbrace{\begin{bmatrix} 0 \\ -q dt \end{bmatrix}}_{\text{Slow}}$$

Configurational Sampling

$$H(q, p) = p^T M^{-1} p / 2 + U(q)$$

$$e^{-\beta H} = e^{-\beta T(p)} e^{-\beta U(q)}$$

In MD we typically are interested in ***q-dependent quantities***. The Hamiltonian is used to enhance exploration but ultimately we don't much care about momenta.

Question: can we use this freedom, together with stochastic dynamics, to **control resonance in multiple timestepping** while taking advantage of a known force field decomposition?

Fast Dynamics:

$$\mathbf{F} = \begin{bmatrix} 0 & 1 \\ -\Omega^2 & -\gamma \end{bmatrix} \quad \begin{bmatrix} q(t) \\ p(t) \end{bmatrix} = \exp(\mathbf{F}t) \begin{bmatrix} q(0) \\ p(0) \end{bmatrix} + \mathbf{B}_t \mathbf{R}$$

$$\mathbf{B}_t \mathbf{B}_t^T = \frac{2\gamma}{\beta} \int_0^t \exp(\mathbf{F}(t-s)) \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \exp(\mathbf{F}^T(t-s)) ds.$$

Combine with kicks by the slow force

$$z_{n+1} = X z_n + Y R_n$$

$$\mathbf{X} = \begin{bmatrix} 1 & 0 \\ -h & 1 \end{bmatrix} \exp(h\mathbf{F}), \quad \mathbf{Y} = \begin{bmatrix} 1 & 0 \\ -h & 1 \end{bmatrix} \mathbf{B}_h$$

$$z_{n+1} = X z_n + Y R_n$$

Stability if

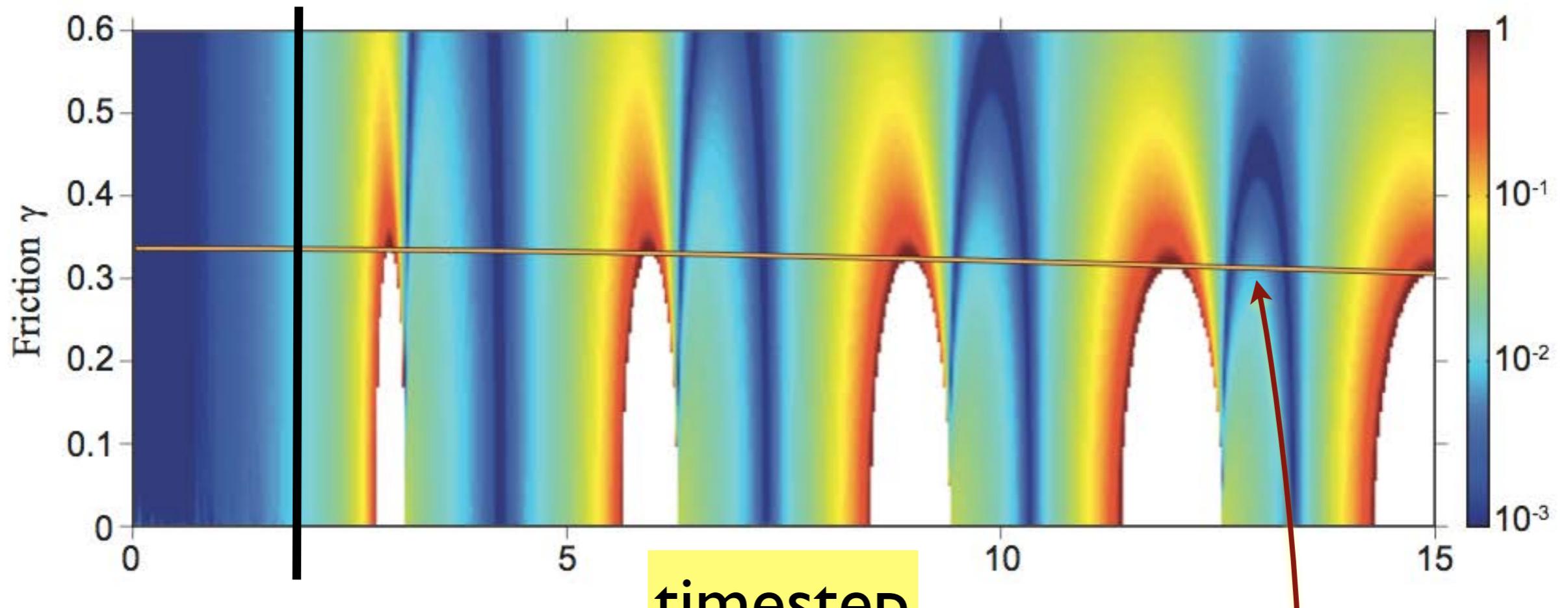
$$\rho(X) \leq 1$$

$$\left| \frac{\eta - h}{\eta} e^{(\eta - \gamma)h/2} + \frac{\eta + h}{\eta} e^{(-\eta - \gamma)h/2} \right| \leq 1 + e^{-\gamma h}$$

$$\frac{h}{\sqrt{4\Omega^2 - \gamma^2}} \leq \sinh\left(\frac{\gamma h}{2}\right)$$

sampling error

noise strength



Verlet stability threshold

timestep

$$\frac{h}{\sqrt{4\Omega^2 - \gamma^2}} \leq \sinh\left(\frac{\gamma h}{2}\right)$$

Langevin can stabilize RESPA but only at high friction

Stochastic Isokinetic Nosé-Hoover (SIN)

L., Margul and Tuckerman, Mol. Phys. 2013

view as 1 dof model

$$\dot{q} = p$$

$$\dot{p} = F(q) - \lambda p$$

$$\dot{\xi}_1 = -\lambda\xi_1 - \xi_1\xi_2$$

$$\dot{\xi}_2 = \xi_1^2 - kT - \gamma\xi_2 + \sqrt{2\gamma kT}\eta(t)$$

$$\lambda = \frac{2p \cdot F - \xi_1^2\xi_2}{2K(p, \xi_1)}$$
 to make $K(p, \xi_1) := p \cdot p + \frac{\xi_1^2}{2} = kT$
isokinetic constraint

Compressible statistical mechanics

M. E. Tuckerman, Y. Liu, G. Ciccotti, and G. J. Martyna, *J. Chem. Phys.* **115**, 1678 (2001).

$$\dot{x} = f(x)$$

conservation laws:

$$C_j(x) = c_j, j = 1, 2, \dots, k$$

phase space compressibility:

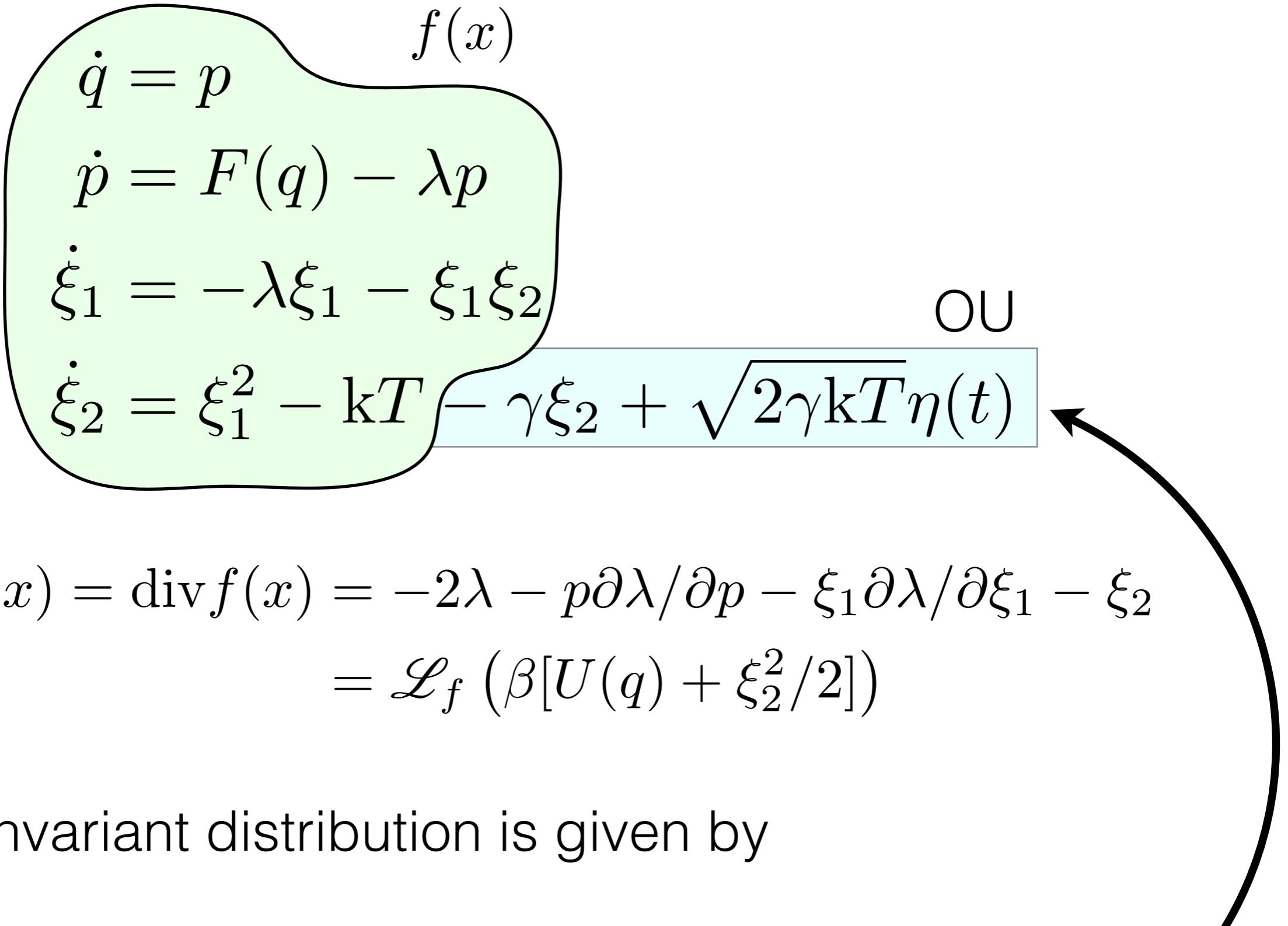
$$\kappa(x) = \operatorname{div} f(x)$$

and find:

$$\kappa = \mathcal{L}_f w$$

Then, if ergodic, the partition function is:

$$\Omega(N, V, \beta^{-1}) = \int e^{-w(x)} \prod_{j=1}^k \delta[C_j(x) - c_j] dx$$



$$\rho(x) = e^{-\beta U(x)} e^{-\beta \xi_2^2/2} \delta[K(p, \xi_1) - kT]$$

$$\rho(x) = e^{-\beta U(x)} e^{-\beta \xi_2^2/2} \delta[K(p, \xi_1) - kT]$$

Dynamics evolves on the **generalized semi-cylinder** defined by periodic boundary conditions (q), and the isokinetic constraint on (p, ξ_1) , with the restriction $\xi_1 > 0$

We need to demonstrate the **Hörmander condition** on the manifold, and find a **Lyapunov function**.

Lyapunov function: 😊

$$\Phi(\xi_2) = 1 + \xi_2^{2s}$$

Ergodicity Property of SIN

$$f = p\partial_q - (q + \lambda p)\partial_p - (\lambda\xi_1 + \xi_2\xi_1)\partial_{\xi_1} + [(\xi_1^2 - k_B T) - \gamma\xi_2]\partial_{\xi_2},$$

$$g = \sigma\partial_{\xi_2}.$$

$$\text{Span}\{\mathbf{f}, \mathbf{g}, [\mathbf{f}, \mathbf{g}], [\mathbf{f}, [\mathbf{f}, \mathbf{g}]], [\mathbf{f}, [\mathbf{f}, [\mathbf{f}, \mathbf{g}]]]\} = T\mathcal{M}$$

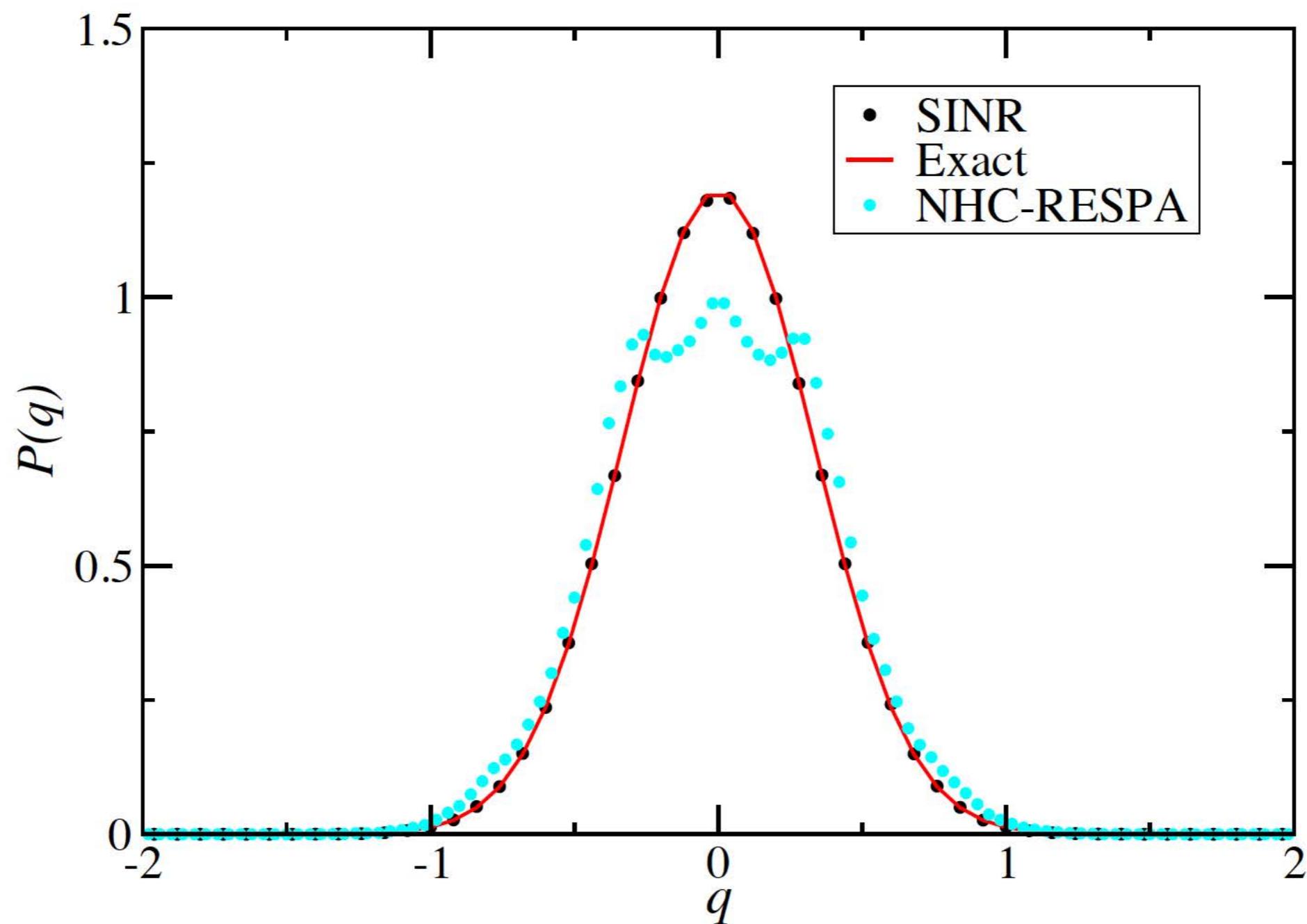
The Alexander Davie commutator



Thus conclude (subject to a minorization condition, left as an exercise) that **SIN** is **ergodic** on \mathcal{M} with a unique invariant distribution which, after marginalization reduces to:

$$\rho(x) = e^{-\beta U(x)}$$

sampling $U(q) = \frac{1}{2}\omega^2 q^2 + \frac{1}{4}gq^4$



SIN(R)

Stochastic Isokinetic Nosé-Hoover (RESPA)

$$dq = pdt,$$

$$dp = Fdt - \left(\frac{p(F^F + F^S) - \frac{1}{2}\mu_1\xi_1^2\xi_2}{\Lambda} \right) pdt,$$

$$d\xi_1 = - \left(\frac{p(F^F + F^S) - \frac{1}{2}\mu_1\xi_1^2\xi_2}{\Lambda} \right) \xi_1 dt - \xi_2 \xi_1 dt,$$

$$d\xi_2 = \mu_2^{-1}(\mu_1\xi_1^2 - \beta^{-1})dt - \gamma\xi_2 dt + \sigma dW,$$

Kinetic part:

$$\mathcal{L}_K = p \frac{\partial}{\partial q},$$

isokinetic Fast force:

$$\mathcal{L}_F = [F^F - (p^2/\Lambda)F^F] \frac{\partial}{\partial p} - \frac{p\xi_1 F^F}{\Lambda} \frac{\partial}{\partial \xi_1},$$

isokinetic Slow force:

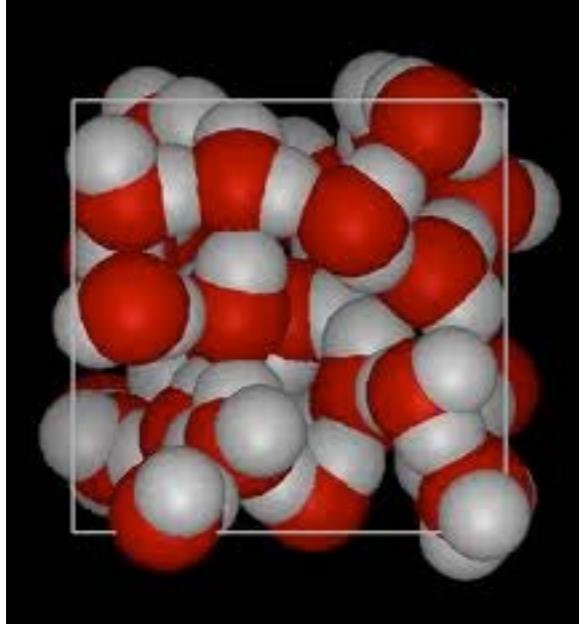
$$\mathcal{L}_S = [F^S - (p^2/\Lambda)F^S] \frac{\partial}{\partial p} - \frac{p\xi_1 F^S}{\Lambda} \frac{\partial}{\partial \xi_1},$$

isokinetic Nosé terms:

$$\mathcal{L}_N = \frac{1}{2} \frac{\mu_1 \xi_1^2 \xi_2 p}{\Lambda} \frac{\partial}{\partial p} + \frac{1}{2} \frac{\mu_1 \xi_1^3 \xi_2}{\Lambda} \frac{\partial}{\partial \xi_1} - \xi_2 \xi_1 \frac{\partial}{\partial \xi_1} + \mu_2^{-1} [\mu_1 \xi_1^2 - \beta^{-1}] \frac{\partial}{\partial \xi_2},$$

Ornstein-Uhlenbeck term:

$$\mathcal{L}_O = -\gamma \xi_2 \frac{\partial}{\partial \xi_2} + \frac{\sigma^2}{2} \frac{\partial^2}{\partial \xi_2^2}.$$



periodic box, 25Å sides,
512 fully flexible
SPC water molecules (all atom)
periodic boundary conditions
three level splitting

$$U = U_{\text{bond}} + U_{\text{s.r.}} + U_{\text{l.r.}}$$

Coulombic forces: Smooth Particle Mesh Ewald

Long ranged forces includes ‘reciprocal space’ part + screened coulomb interactions within the simulation cell

Short-ranged Regime: e.g. Lennard-Jones cutoffs to 6Å, ‘real space’ part of Ewald.

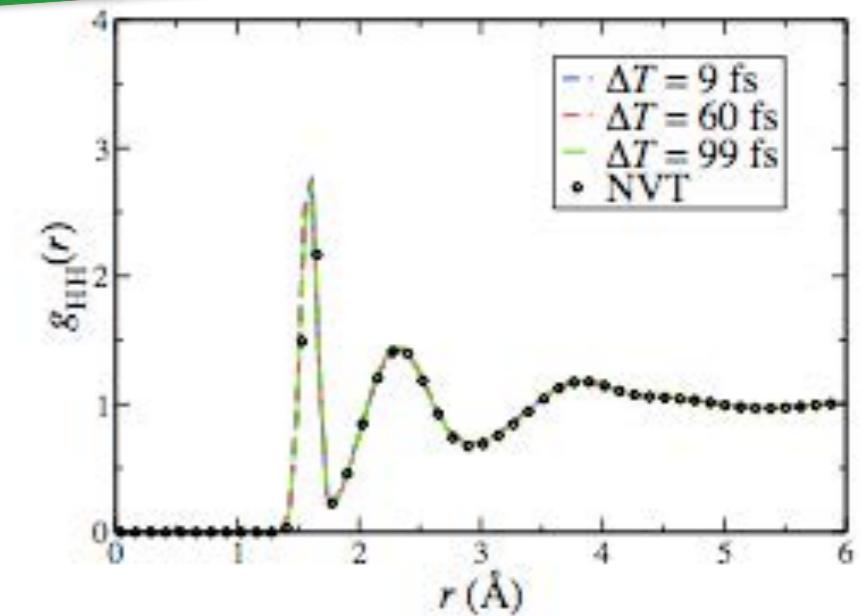
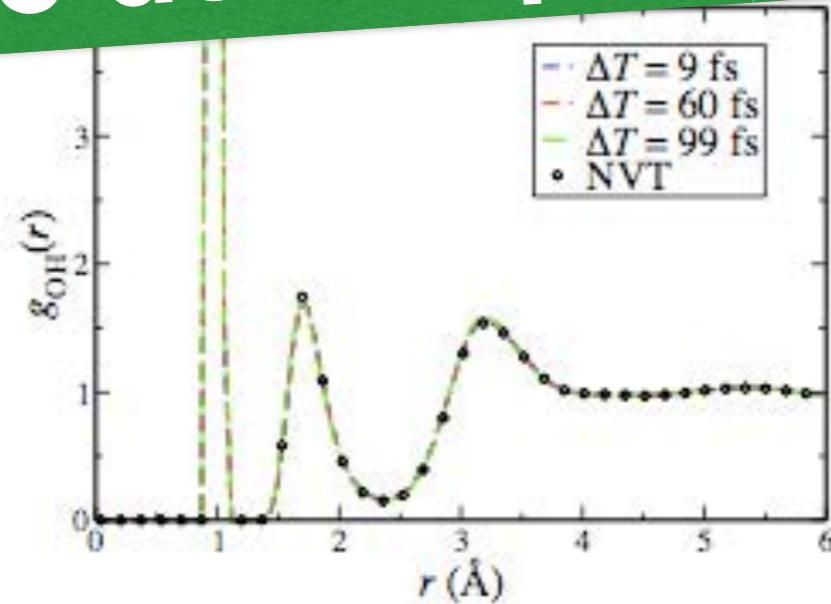
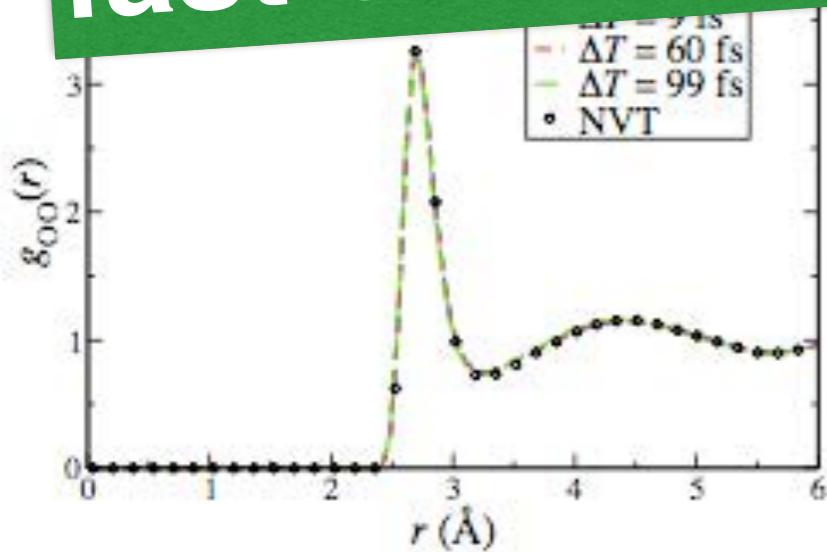
Bond part: Intramolecular interactions, e.g. angle and length bonds.

Flexible H₂O simulations @ $\delta t = 99\text{fs}$

$$\delta t_{\text{inner}} = 0.5\text{fs}, \quad \delta t_{\text{mid}} = 3\text{fs}$$

Take Home Message #6

Isokinetic constraints can stabilize
stochastic multiple timestepping based on
fast-slow force decomposition.



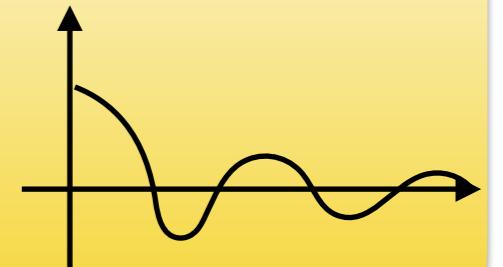
Ensemble Preconditioning

Goal: redesign the dynamics (and integrator) to enhance the rate of convergence for typical observables f

$$\langle f \rangle = \int f(x) \rho(x) dx = \lim_{N \rightarrow \infty} N^{-1} \sum_{t=1}^N f(x_t)$$

figure of merit = **Integrated Autocorrelation Time (IAT)**

$$\tau_f = 1 + 2 \sum_{t=1}^{\infty} \text{cor}(f(x_t), f(x_0))$$

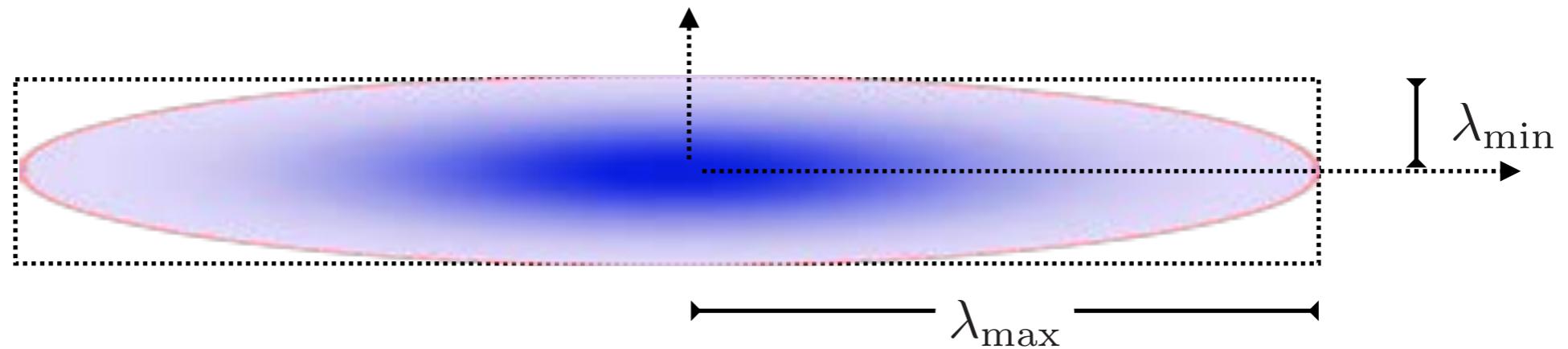


We would like to have τ_f as small as possible

Motivating Example:

$$\pi(x) = \exp\left(-\frac{1}{2}x^T \Sigma^{-1} x\right)$$

Eigenvalues: $0 < \lambda_{\min} < \dots < \lambda_{\max} = \rho(\Sigma)$



For MCMC schemes like Euler-Maruyama or Leimkuhler-Matthews, **stability requires** $h = O(\lambda_{\min})$

But for $f(x) = x \cdot \mathbf{e}_{\max}$ $\tau_f = O(\lambda_{\max}/\lambda_{\min})$

Poor Scaling \Rightarrow **Slow Convergence**

Ensemble Preconditioning

More generally, we wish to sample problems with complicated energy functions, where each basin or local approximation may be very poorly scaled.



Related Concepts

Stochastic Newton schemes

BFGS Method

MC Hammer (Goodman and Weare)

Compare to: Riemannian Manifold HMC (Girolami et al)

Ensemble Preconditioning

Use local information to estimate inverse Hessian matrix; precondition (rescale) dynamics to enhance convergence (reduce IAT)

Wishlist:

- Increase efficiency by reducing the IAT
- Compute the preconditioning based on a local ensemble approximation
- Allow for inertial effects (underdamped Langevin/HMC)

Idea: Use a collection of “walkers” to generate local covariance information and use this to estimate the inverse Hessian adaptively

Procedure

Use an ensemble of L walkers:

$$Q = (q_1, q_2, \dots, q_L) \in \mathbb{R}^{dL}, \quad P = (p_1, p_2, \dots, p_L) \in \mathbb{R}^{dL},$$

$$\bar{\pi}(Q, P) = \prod_{i=1}^L \hat{\pi}(q_i, p_i), \quad \int \bar{\pi}(Q, P) \text{d}P = \prod_{i=1}^L \pi(q_i).$$

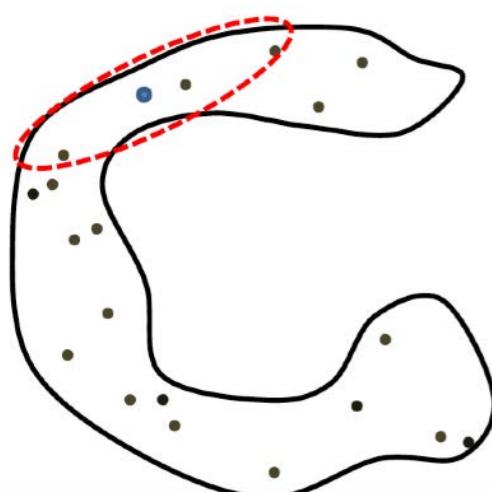
*each walker samples the
same target distribution $\pi(q)$*

We construct dynamics in the extended space and compute ensemble averages by marginalisation over the individual walkers.

$$\dot{Q} = B(Q)P,$$

$$\dot{P} = B(Q)^T \nabla \log(\pi(Q)) + \operatorname{div}(B(Q)^T) - \gamma P + \sqrt{2\gamma} \eta(t).$$

$$B(Q) = \operatorname{diag}(B_1(Q), B_2(Q), \dots, B_L(Q))$$



$$B_i(Q) = \sqrt{I_d + \eta \operatorname{wcov}(Q_{[i]}, \omega_{\lambda(Q_{[i]}, q_i)})}$$

*blend with identity
(robustness)*

*collects
covariance
info
of nearby
walkers*

$$Q_{[i]} = (q_1, q_2, \dots, q_{i-1}, q_{i+1}, \dots, q_L)$$

Basing B_i on other walkers only
eliminates the problems of multiplicative noise

Discretization of SDEs: similar to **BAOAB**

Discretization

BAOABish:

$$p^{(n+1/2)} = p^{(n)} + \frac{\delta t}{2} F(q^{(n)}),$$

$$q^{(n+1/2)} = q^{(n)} + \frac{\delta t}{2} B\left(q^{(n+1/2)}\right) p^{(n+1/2)}$$

$$\begin{aligned}\hat{p}^{(n+1/2)} &= \alpha p^{(n+1/2)} + \frac{(\alpha + 1)\delta t}{2} \text{div}\left(B\left(q^{(n+1/2)}\right)^T\right) \\ &\quad + \sqrt{1 - \alpha^2} \mathbf{R}^{(n)}\end{aligned}$$

$$q^{(n+1)} = q^{(n+1/2)} + \frac{\delta t}{2} B\left(q^{(n+1/2)}\right) \hat{p}^{(n+1/2)}$$

$$p^{(n+1)} = \hat{p}^{(n+1/2)} + \frac{\delta t}{2} F(q^{(n+1)})$$

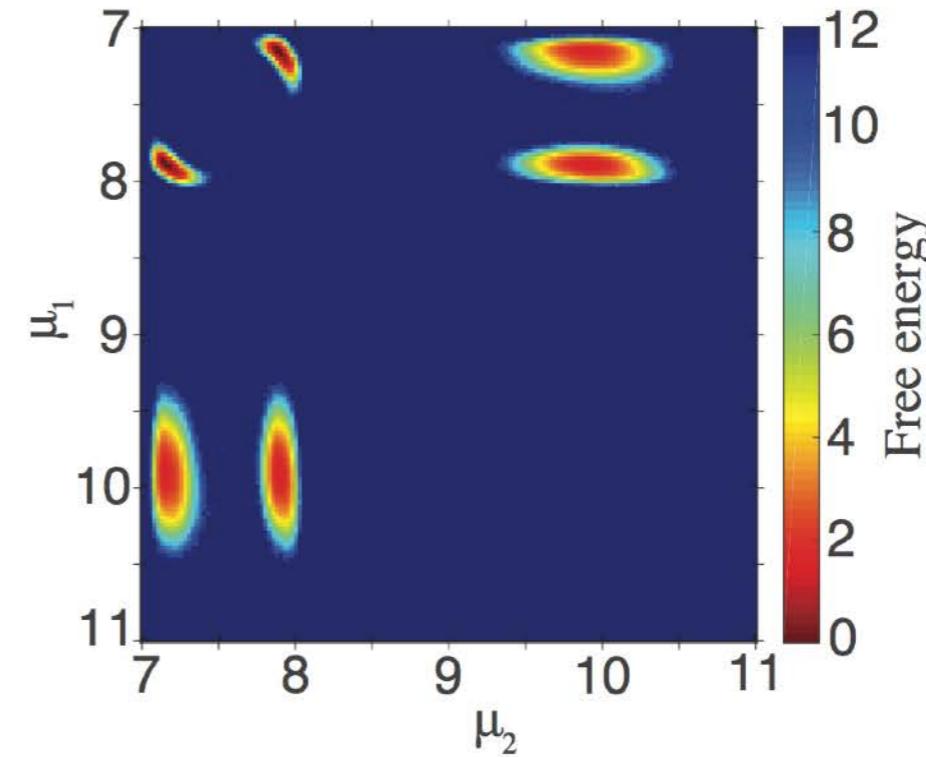
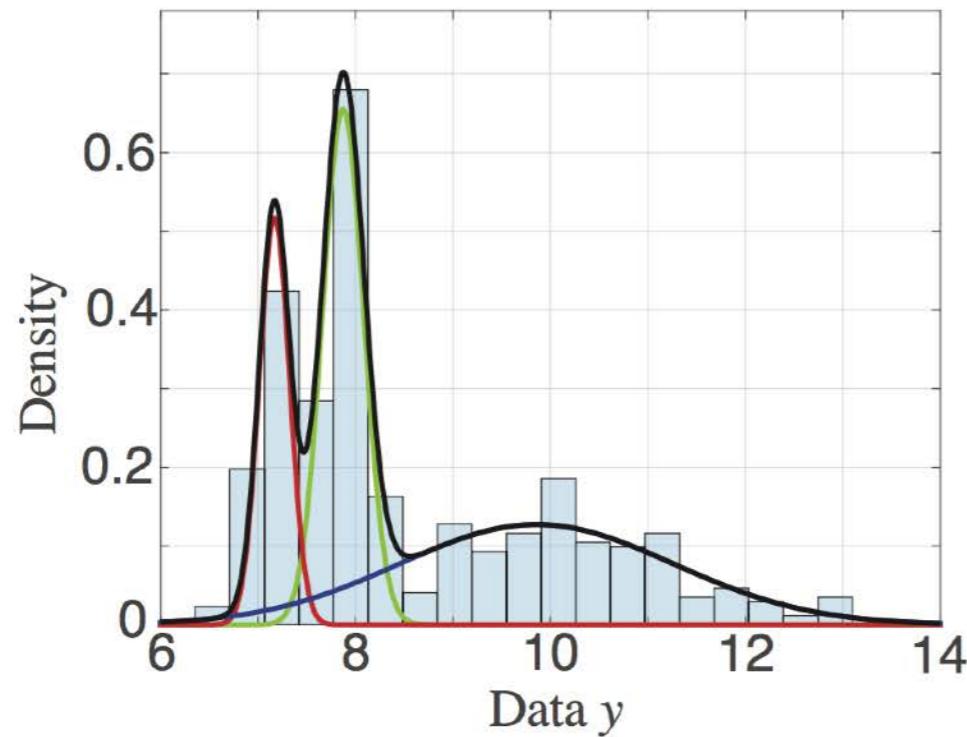
$$\alpha = e^{-\gamma h}$$

Gaussian Mixture Model: Hidalgo Stamps

“Adventures in Stamp Collecting”

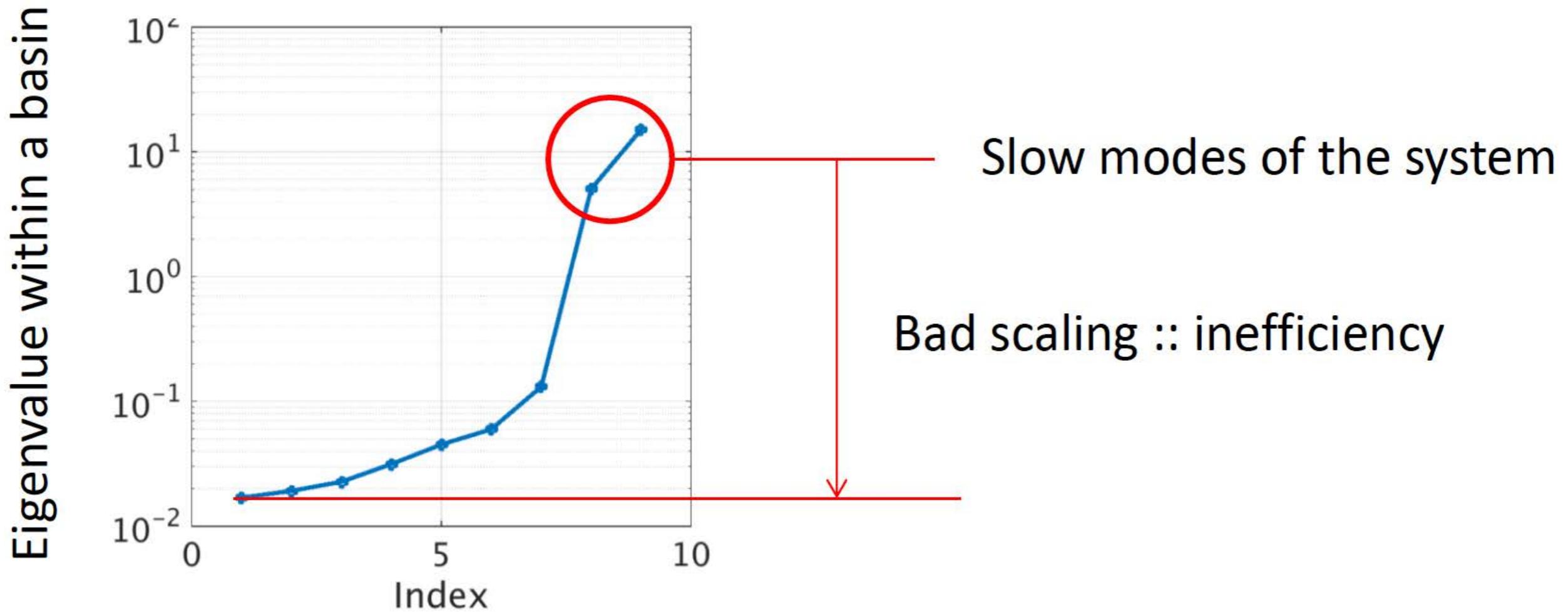
$$\sum_{n=1}^3 z_n \mathcal{N}(\mu_n, \lambda_n^{-2})$$

Dataset: Thickness of
485 stamps from
Mexico in 1872.



- (moderately) poorly scaled basins
- multimodal due to “label switching” symmetry

Within one basin, we have bad scaling:



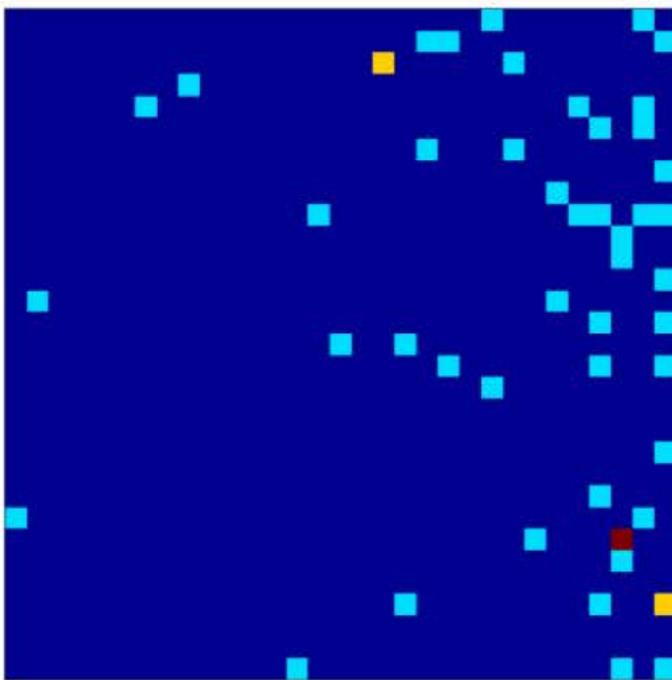
But the eigensystems are **different in different basins**, so the localized covariance is needed...

Gaussian Mixture Model: Hidalgo Stamps

Integrated Autocorrelation Times of Different Schemes

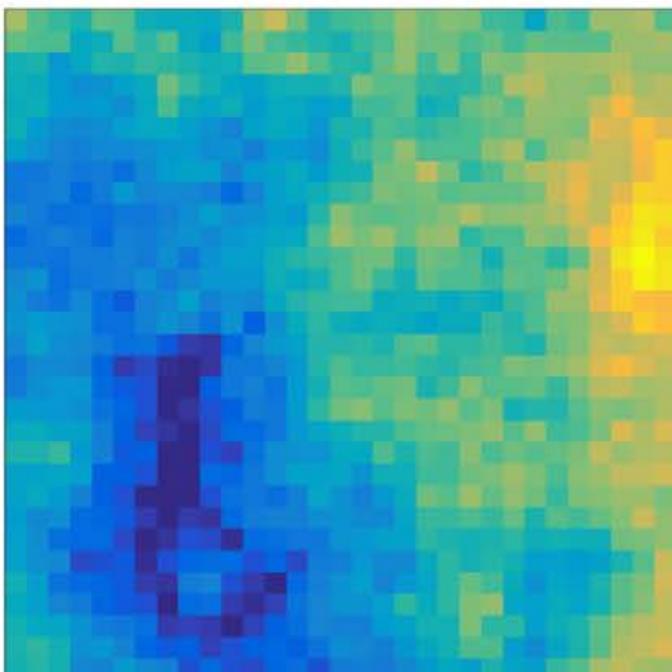
Scheme	$\min(z)$	$\max(\lambda)$	$\min(\mu)$	β
HMC	21495	42935	27452	7148
Langevin Dynamics	6825	13279	8384	4641
Ensemble Q-N	69	83	98	115

Numerical test: Log-Gaussian Cox model



Observations X

Means $\exp(Y_{i,j})$



Break $[0,1]^2$ into a 32×32 grid.

Observed intensity in box (i,j) is $X_{i,j}$,
Poisson distributed with mean

$$\Lambda(i,j)/32^2, \quad \Lambda(i,j) = \exp(Y_{i,j})$$

where $Y \sim N(\mu, \Sigma)$, with

$$\Sigma_{(i,j),(i',j')} = \sigma^2 \exp[-\sqrt{(i-i')^2 + (j-j')^2}/(32\beta)]$$

We generate synthetic data X using

$$\sigma^2 = 1.91, \quad \beta = 1/33, \quad \mu = \log(126) - \sigma^2/2$$

We fix μ and aim to infer likely Y , using
hyperparameters σ^2, β , with prior $\text{gamma}(2, \frac{1}{2})$

Log Gaussian Cox Model

Scheme	x	σ^2	β	Efficiency
HMC	800.7	1041.6	1318.7	1.0
RMHMC	2158.9	34.0	1502.0	0.15
LD	1051.1	1051.1	1051.1	1.0

Take Home Message #7

Using a collection of walker particles can improve conditioning of poorly scaled basins and increase sampling efficiency.

Ensemble Quasi-Newton python package

http://bitbucket.org/c_matthews/ensembleqn

7 Messages of this Lecture

1. The right order of Langevin building blocks provides substantial improvements in accuracy (low sampling bias).
2. Adaptive Langevin dynamics offers a simple method for removing bias in noisy gradient simulations while maintaining high accuracy (low bias).
3. Constraints don't buy much stability in deterministic MD, but they do improve accuracy!
4. Combining constraints and multiple timestepping can improve accuracy and stability
5. Constrained Langevin dynamics, using a multiple timestepping implementation of geodesic flow, can yield big MD sampling improvements with high accuracy.
6. Isokinetic constraints can stabilize stochastic multiple timestepping based on fast-slow force decomposition.
7. Using a collection of walker particles can improve conditioning of poorly scaled basins and increase sampling efficiency.