

# Molecular Dynamics

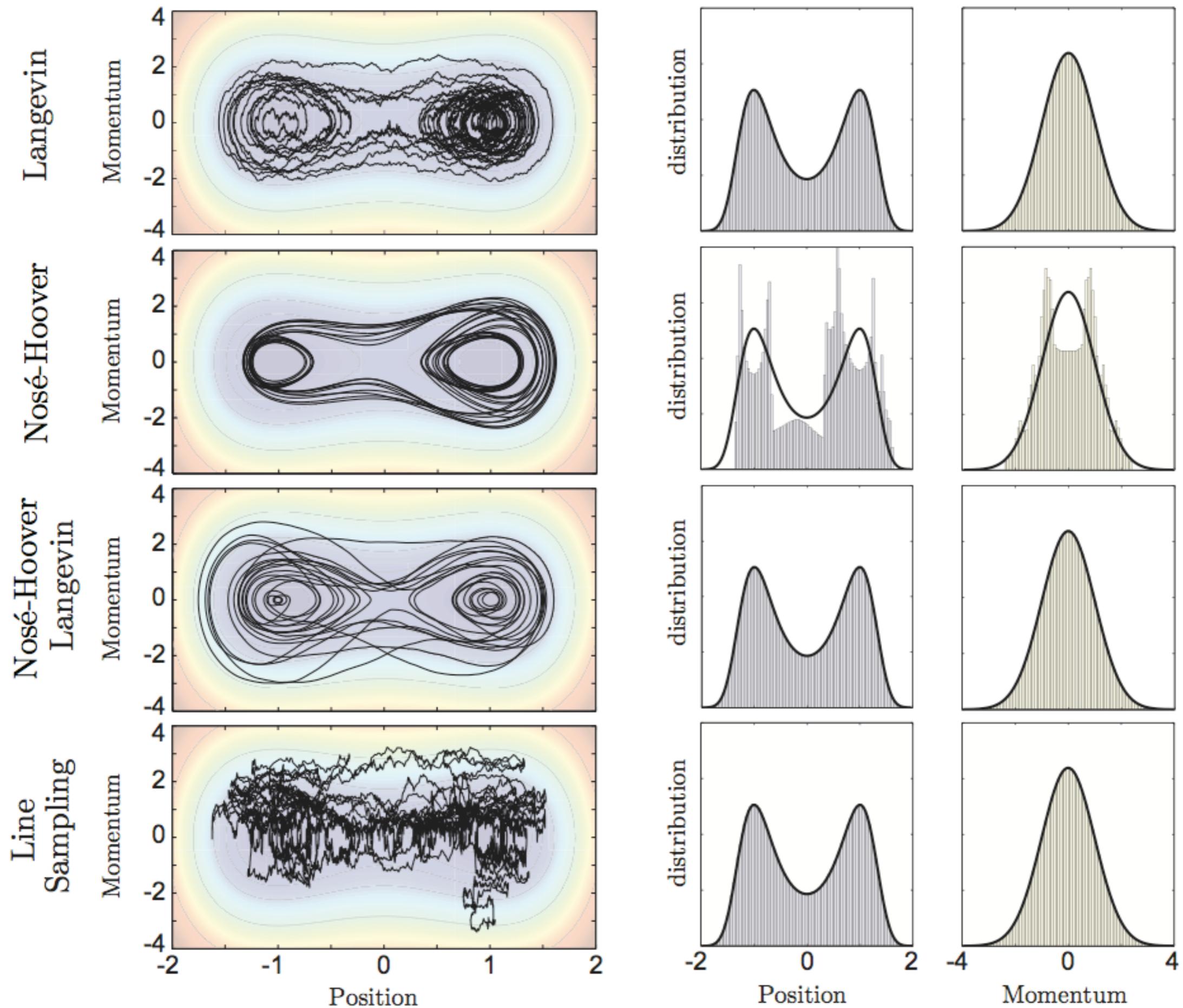
## Day 5

Ben Leimkuhler

**general thermostats and SDEs  
Nosé-Hoover methods  
adaptive thermostats  
noisy gradients  
ensemble preconditioning**

Peking 2018

From L., Generalized Bulgak-Kusnezov Thermostats, PRE, 2010



# Finding the “Right” Dynamics for the Job

There are many different stochastic models that can be used in MD, but they can have very different efficiencies for a particular task.

## Overdamped Langevin Dynamics

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

**great** for sampling well scaled multivariate Gaussian distribution,

**awful** for a highly corrugated landscape

## Nosé-Hoover

**gentle** - good for autocorrelation functions in systems with strong internal mixing properties...

**not ergodic** - lousy for nucleic acid simulations in implicit solvent

# Thermostats

Gibbs distribution

**Overdamped Langevin**

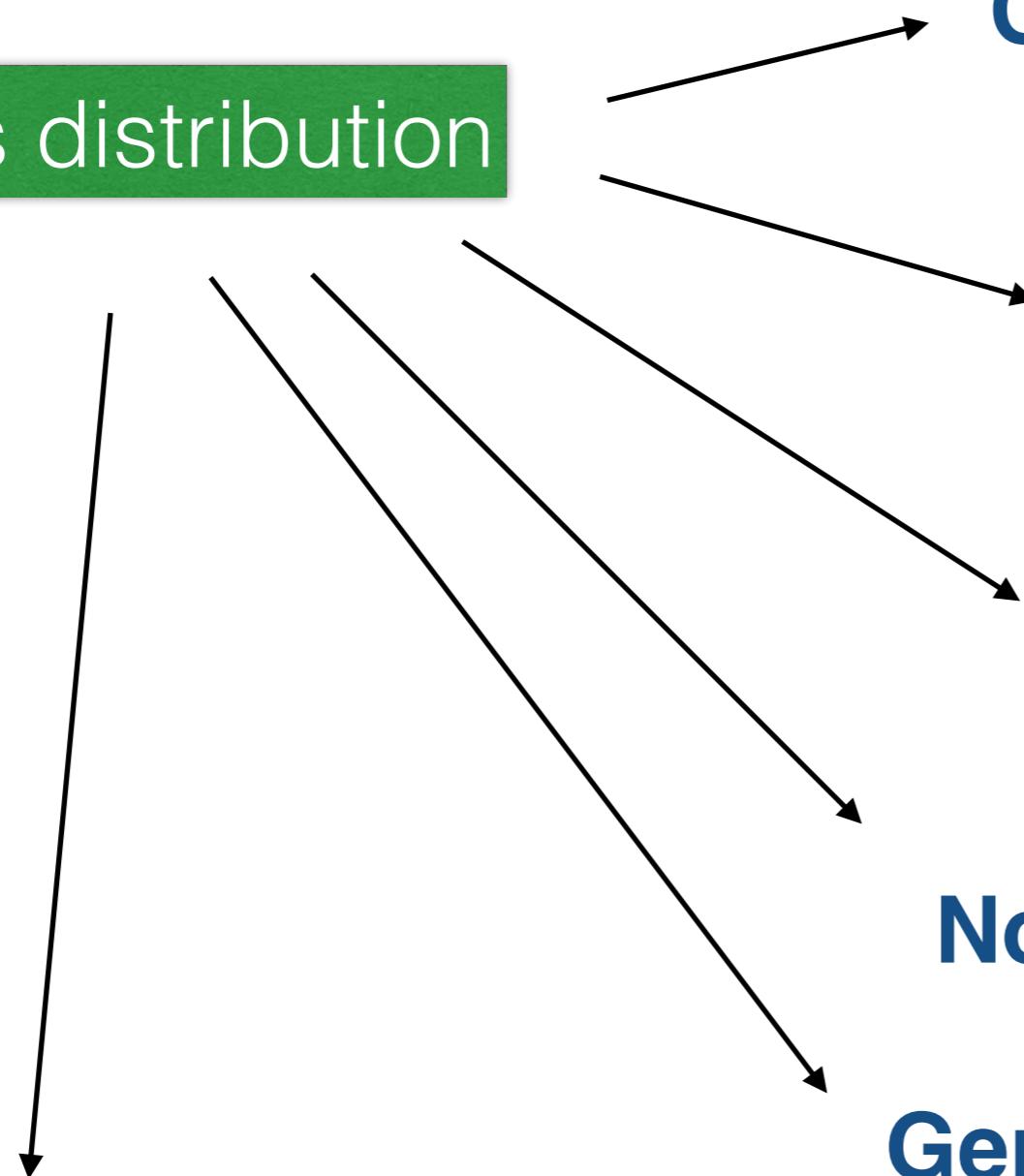
**Langevin Dynamics**

**Nosé-Hoover Langevin**

**Nosé-Hoover Dynamics**

**Generalized  
Bulgac-Kusnezov**

**Preconditioned Methods  
Ensemble Quasi-Newton**



# A generic sampling dynamics

$$dx = [J(x) + S(x)]\nabla \log \pi(x) + \nabla \cdot [J(x) + S(x)] + \sqrt{2S(x)}dW$$

**$J(x)$  antisymmetric**  
 **$S(x)$  symmetric**

Includes, e.g.,

SDEs like **Brownian** and **Langevin** dynamics  
**non-reversible perturbation** methods  
various **ensemble sampling** schemes

Questions:

- Which approach **converges most rapidly?** (small IAT)
- What is the **sampling bias under discretization?**
- How to effectively **combine with extension?**

# Generalized Sampling

Up to know we have assumed the situation of a known distribution with invariant density

$$\rho \propto e^{-\beta U(q)}$$

What if we don't know  $U$  or cannot exactly resolve the force?

**Multiscale models, e.g. ab initio MD Methods and QM/MM methods** (heating due to force mismatch)

**Nonequilibrium MD** (e.g Shear Flows)

Applications in **Bayesian Inference & Big Data Analytics**

# Problems for today

1. How to **gently perturb Hamiltonian dynamics** in order to achieve thermal equilibration.
2. How to handle **noisy gradient systems** and driven systems efficiently, in particular with momentum constraints for shear flow applications.
3. How to **accelerate convergence** to equilibrium by use of an **ensemble of “particles”** (walkers).

# Additivity

The thermostats can be **combined** in most cases without altering their effectiveness (often improving it).

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

$$\mathcal{L}_{f+g} = \mathcal{L}_f + \mathcal{L}_g$$

$$\begin{aligned}\mathcal{L}_f^\dagger \rho &= 0 \\ \mathcal{L}_g^\dagger \rho &= 0\end{aligned}\Rightarrow \mathcal{L}_{f+g}^\dagger \rho = 0$$

Works for **SDEs** too...

# Extension

Many schemes make good use of the concept of **extension**

$$\int \pi(x) \tilde{\pi}(y) dy \propto \pi(x)$$

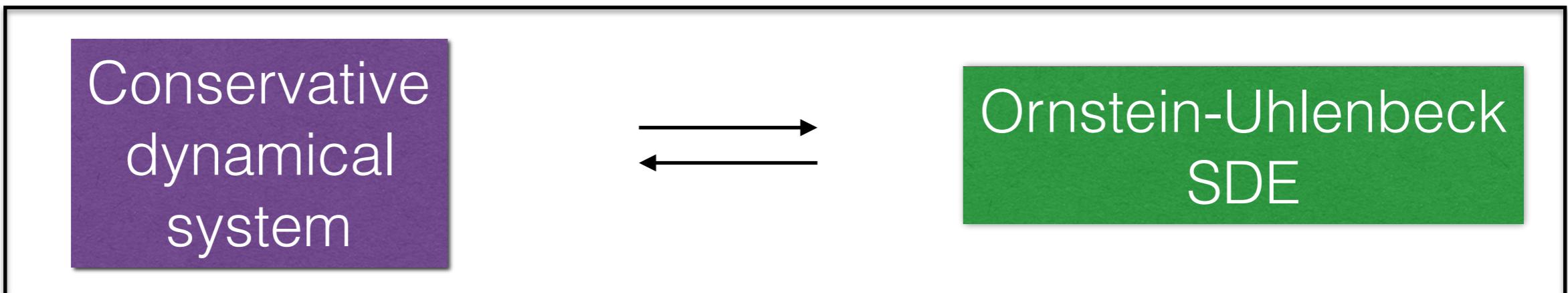
This looks banal but the key point is that although  $x$  and  $y$  decouple in the invariant distribution, they may be tightly coupled in the associated SDEs.

Example: **Langevin dynamics**

$$\int e^{-\beta p^2/2} e^{-\beta U(x)} dp \propto e^{-\beta U(x)}$$

# Remote control of thermal equilibration

Ex: Langevin dynamics



Preserves

$$\rho_\beta = e^{-\beta p^T M^{-1} p} e^{-\beta U}$$

Ergodic for

$$\bar{\rho} = e^{-\beta p^T M^{-1} p}$$

- The two systems are both **compatible** with  $\rho_\beta$
- Sufficient mixing

The ergodicity of the OU process implies ergodicity of the full system

# **Nosé-Hoover and gentle equilibration**

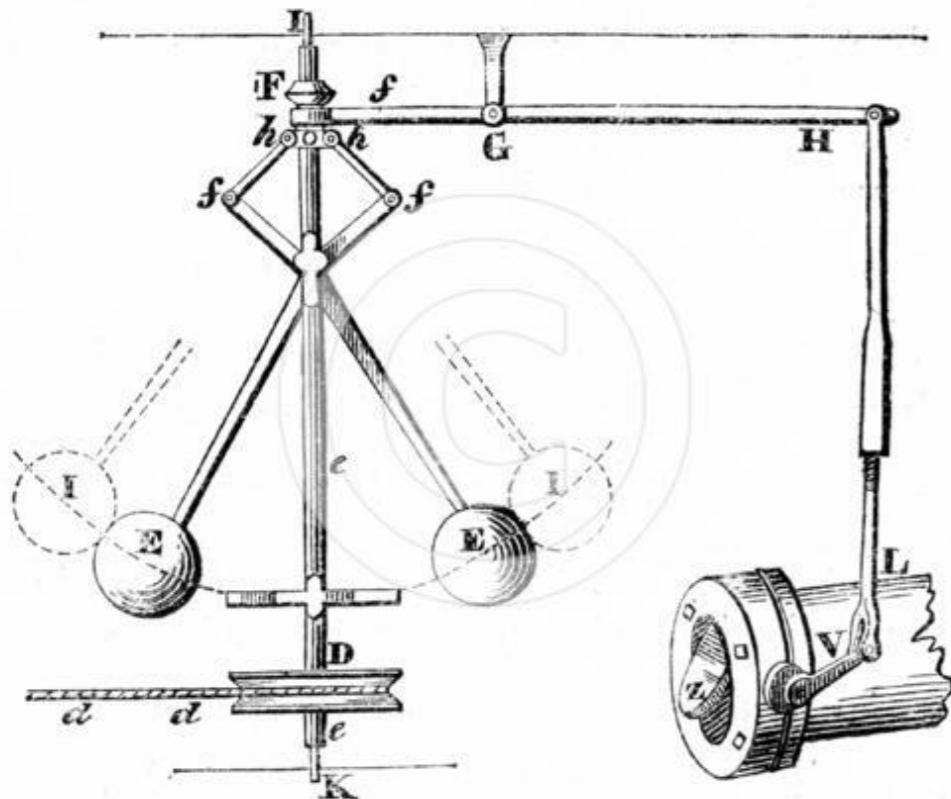
# Where I learned about Nosé Dynamics

Seminar, Cambridge University 1997:  
Nosé Dynamics



**Sir Henry Peter Francis Swinnerton-Dyer, 16th Baronet KBE FRS**  
**Number Theorist, Student of Littlewood, Polya and Sylvester Prizeholder**  
**Vice-Chancellor of Cambridge University 1973-83**

# 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

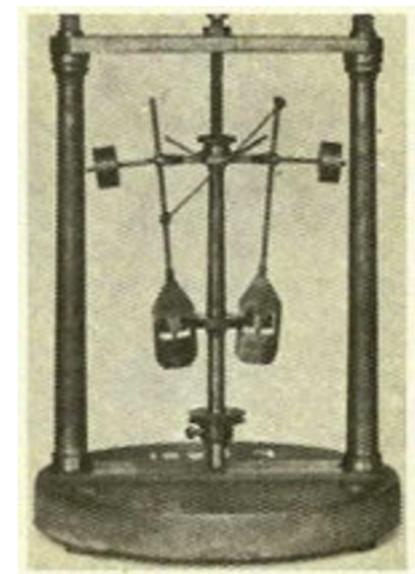
$$e^{-\beta[p^2/2+U(q)]} \times e^{-\beta\xi^2/2}$$

# Problems with the Gibbs Governor

**It doesn't actually work.**

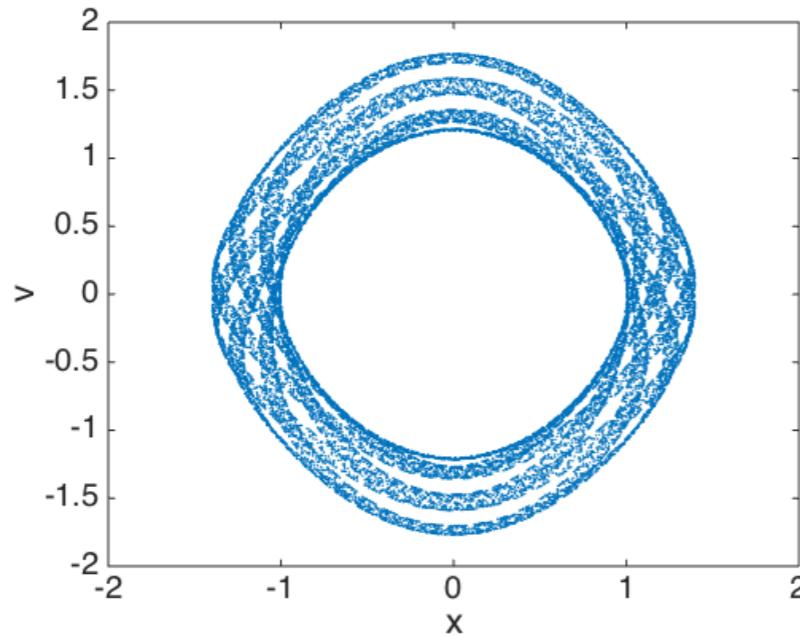
**It's not *the Gibbs Governor*. This is:**

*Undergraduate research  
project of Josiah Willard Gibbs*

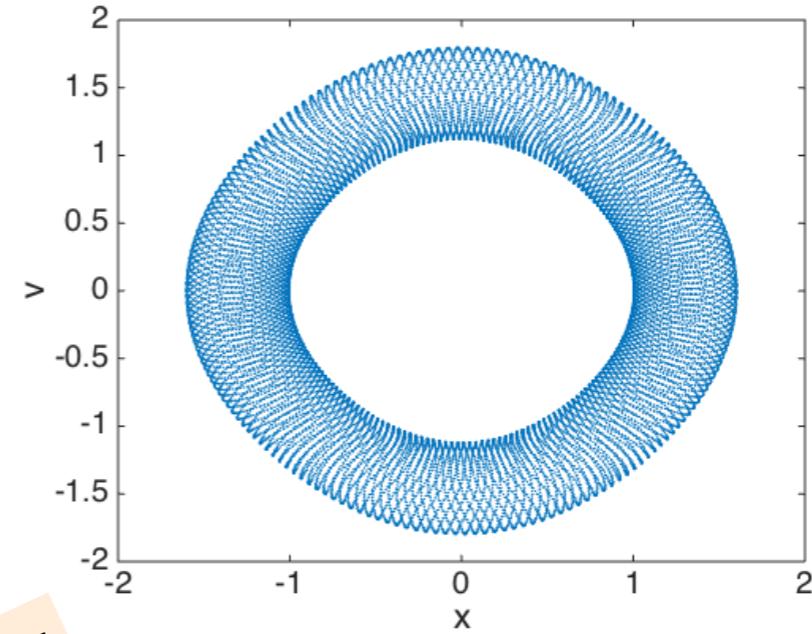


# Nosé-Hoover Dynamics for Harmonic Oscillator

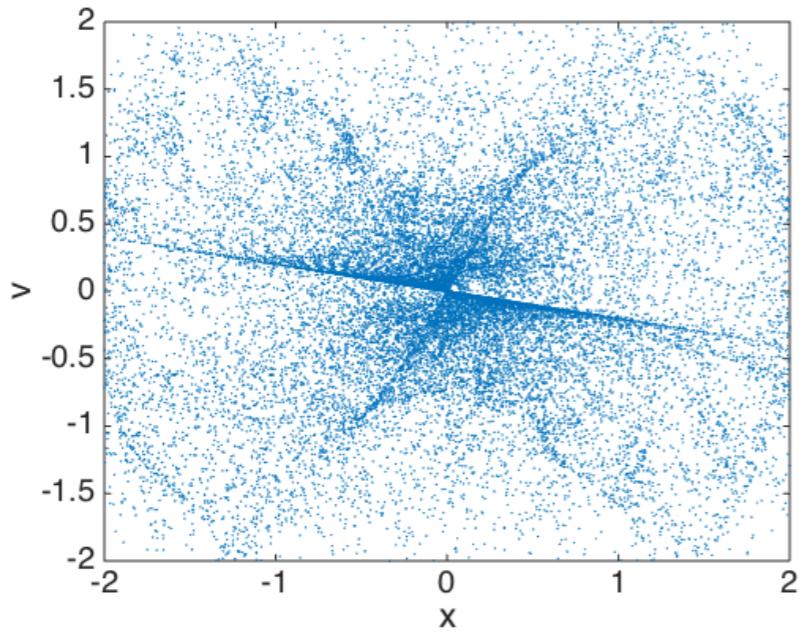
$\mu = 1$



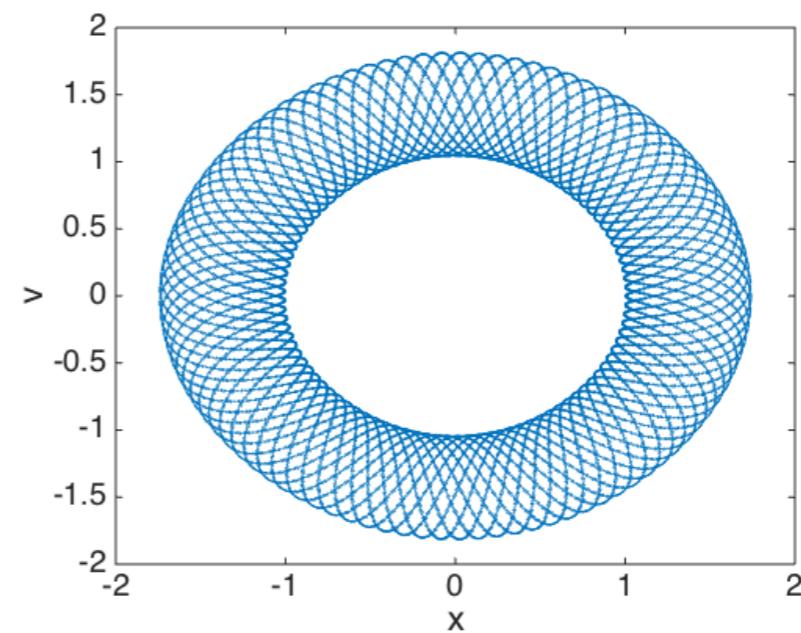
$\mu = 2$



$\mu = 1/2$

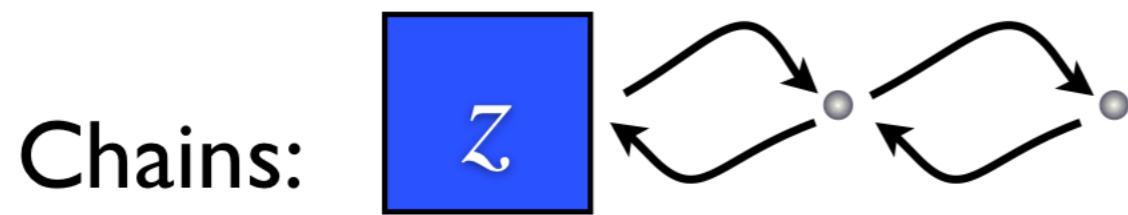


$\mu = 4$



All Wrong!

# Nosé-Hoover Chains also are not ergodic..



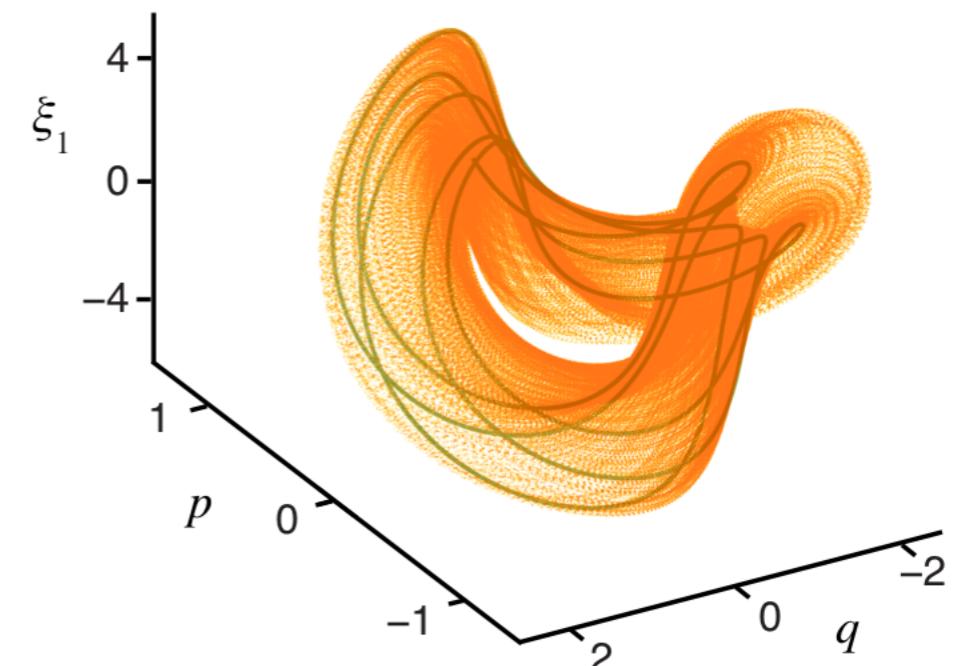
$$\dot{q} = p$$

$$\dot{p} = -q - \xi p$$

$$\dot{\xi}_1 = \mu_1^{-1}(p^2 - kT) - \xi_1 \xi_2$$

$$\dot{\xi}_2 = \mu_2^{-1}(\mu_1 \xi_1^2 - kT)$$

$$\mu_1 = 0.2, \quad \mu_2 = 1$$



# Stochastic version: **Nosé-Hoover-Langevin** dynamics

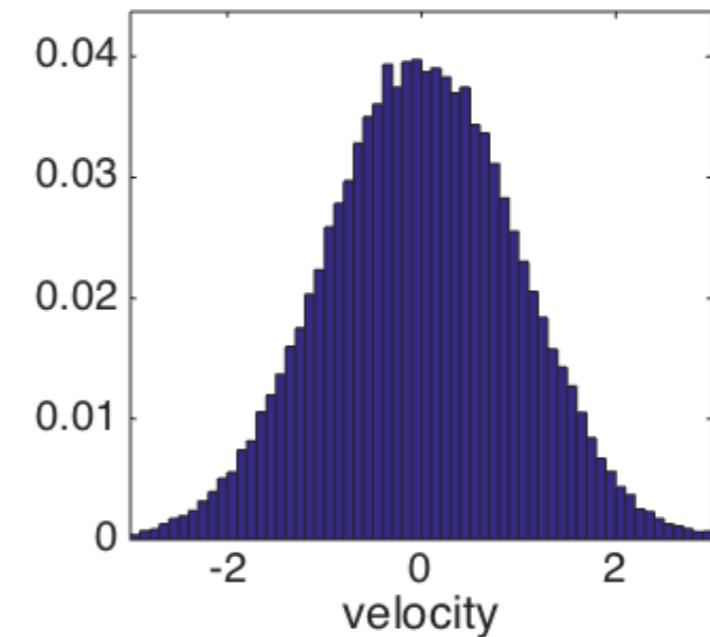
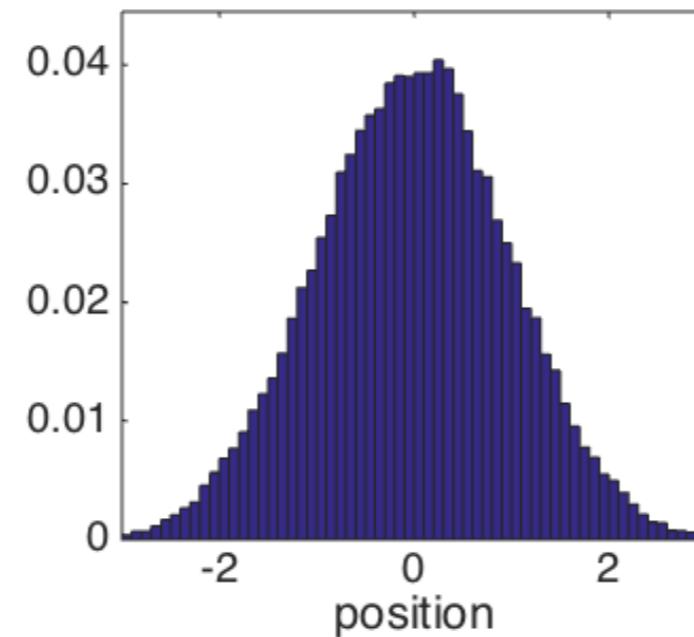
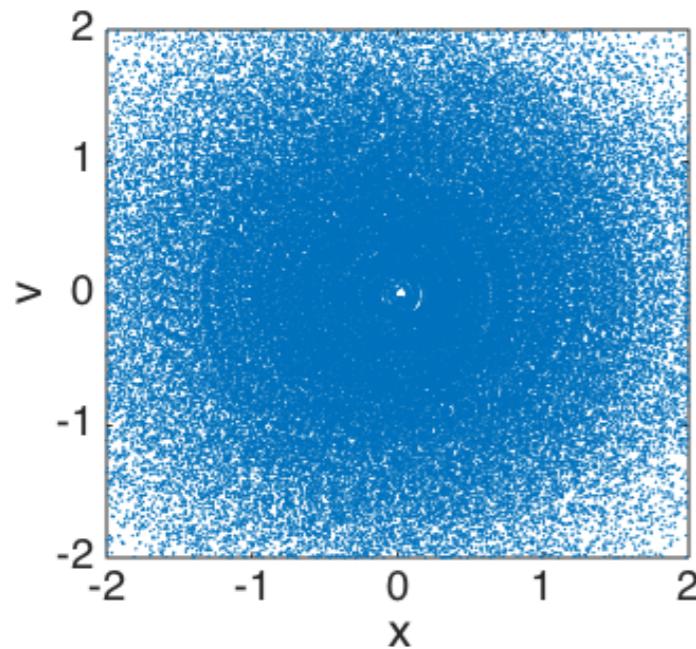
$$\dot{q} = p$$

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

$$\dot{\xi} = p^2 - kT - \gamma\xi + \sigma\eta(t)$$

scalar OU process

‘Histograms’



*matches theoretical behavior*

# Nosé-Hoover-Langevin

$$dq = M^{-1}p dt$$

$$dp = -\nabla U(q)dt - \xi p dt$$

$$d\xi = [p^T M^{-1} p - nk_B T] dt - \gamma \xi dt + \sqrt{2\beta^{-1}\gamma} dW_t$$

$\dot{q} = M^{-1}p$	<b>preserves</b>	$e^{-\beta p^T M^{-1} p / 2} e^{-\beta U(q)}$
$\dot{p} = -\nabla U$	<b>preserves</b>	$e^{-\beta p^T M^{-1} p / 2} e^{-\beta \xi^2 / 2}$
$d\xi = -\gamma \xi dt + \sqrt{2\beta^{-1}\gamma} dW_t$	<b>preserves ergodically</b>	$e^{-\beta \xi^2 / 2}$

# Ergodicity of NHL

NHL is clearly compatible with an extended Gibbs distribution meaning that

$$\mathcal{L}_{\text{NHL}}^\dagger [\rho_\beta e^{-\beta \xi^2/2}] = 0$$

We can also prove it is ergodic by using the theory developed for Langevin dynamics and explained in the previous lectures.

[L., Noorizadeh, Theil 2009]

# Harmonic system w/o resonance

---

$$H = \frac{p^T M^{-1} p}{2} + \frac{q^T B q}{2} \quad q, p \in \mathbb{R}^d$$

$$A = M^{-1} B, \quad A\varphi_k = \omega_k \varphi_k$$

$$\omega_k \neq \omega_l, \quad k \neq l$$

Theorem: **NHL is ergodic** on a large part of  $\mathbb{R}^{2d+1}$

$$f_0 = \begin{bmatrix} M^{-1}p \\ -Bq - \xi p \\ p^T M^{-1} p - d\beta^{-1} - \xi \end{bmatrix}, \quad f_1 = \mathbf{e}_{2d+1}$$

Example: clamped harmonic spring chain.

# Marginalization

If a system is ergodic for

$$\tilde{\rho}(q, p, \xi) = e^{-\beta(p^T M^{-1} p / 2 + U(q))} e^{-\beta \xi^2 / 2}$$

then the paths of the system provide Gibbs-weighted averages

$$\lim_{\tau \rightarrow \infty} \int_0^\tau \varphi(q(t), p(t), \xi(t)) dt = \int_{\Omega} \varphi(q, p, \xi) \tilde{\rho}(q, p, \xi) dq dp d\xi$$

and

$$\lim_{\tau \rightarrow \infty} \int_0^\tau \psi(q(t), p(t)) dt = \int_{\Omega} \psi(q, p) \rho_\beta(q, p) dq dp$$

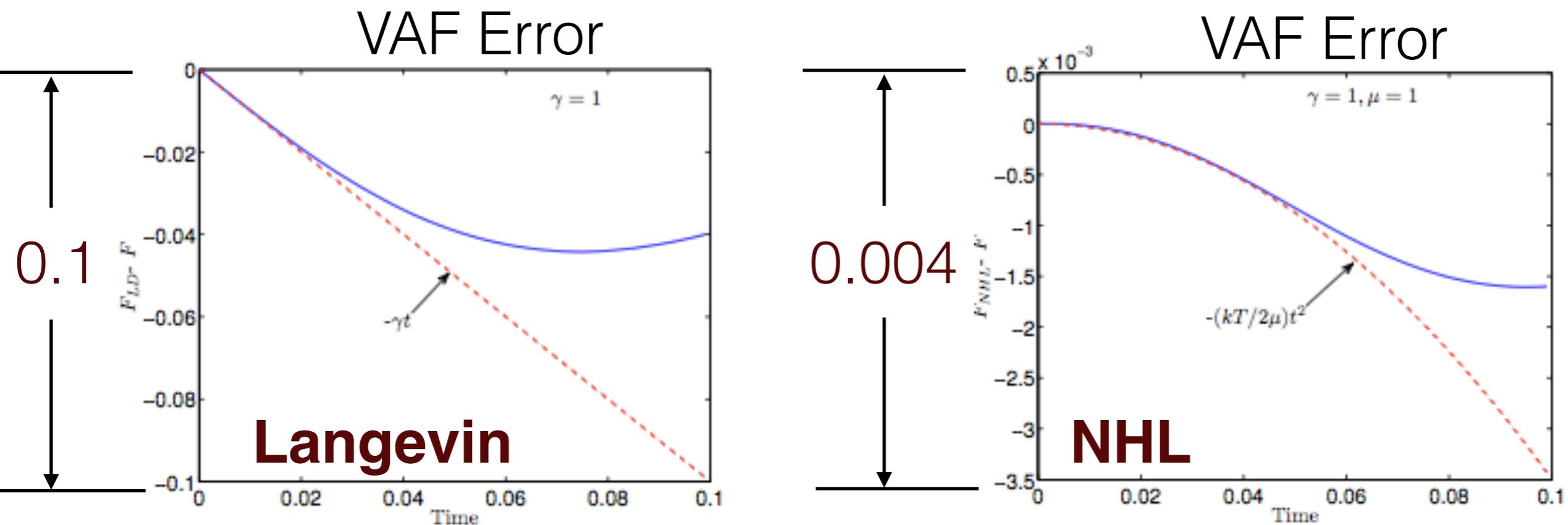
and

$$\lim_{\tau \rightarrow \infty} \int_0^\tau \eta(q(t)) dt = \int_{\Omega} \eta(q) e^{-\beta U(q)} dq$$

# “Gentle” property of NHL

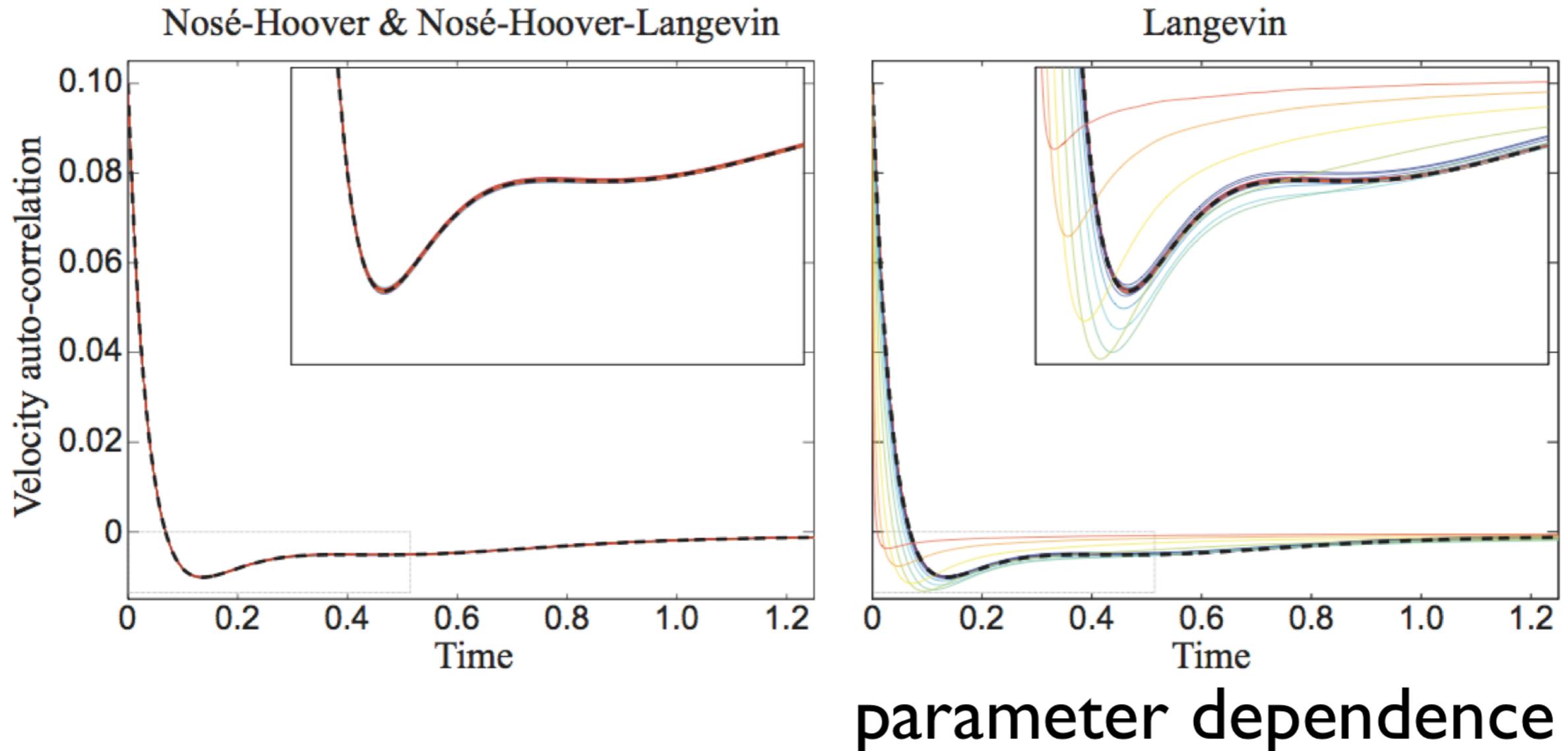
We can show that NHL is a “gentle” thermostat: dynamical properties are mildly perturbed for a given rate of convergence of kinetic energy.

[L., Noorizadeh and Penrose, J. Stat. Phys., 2011]



Similar (but less smooth) for “Stochastic Velocity Rescaling” of G. Bussi, D. Donadio and M. Parrinello

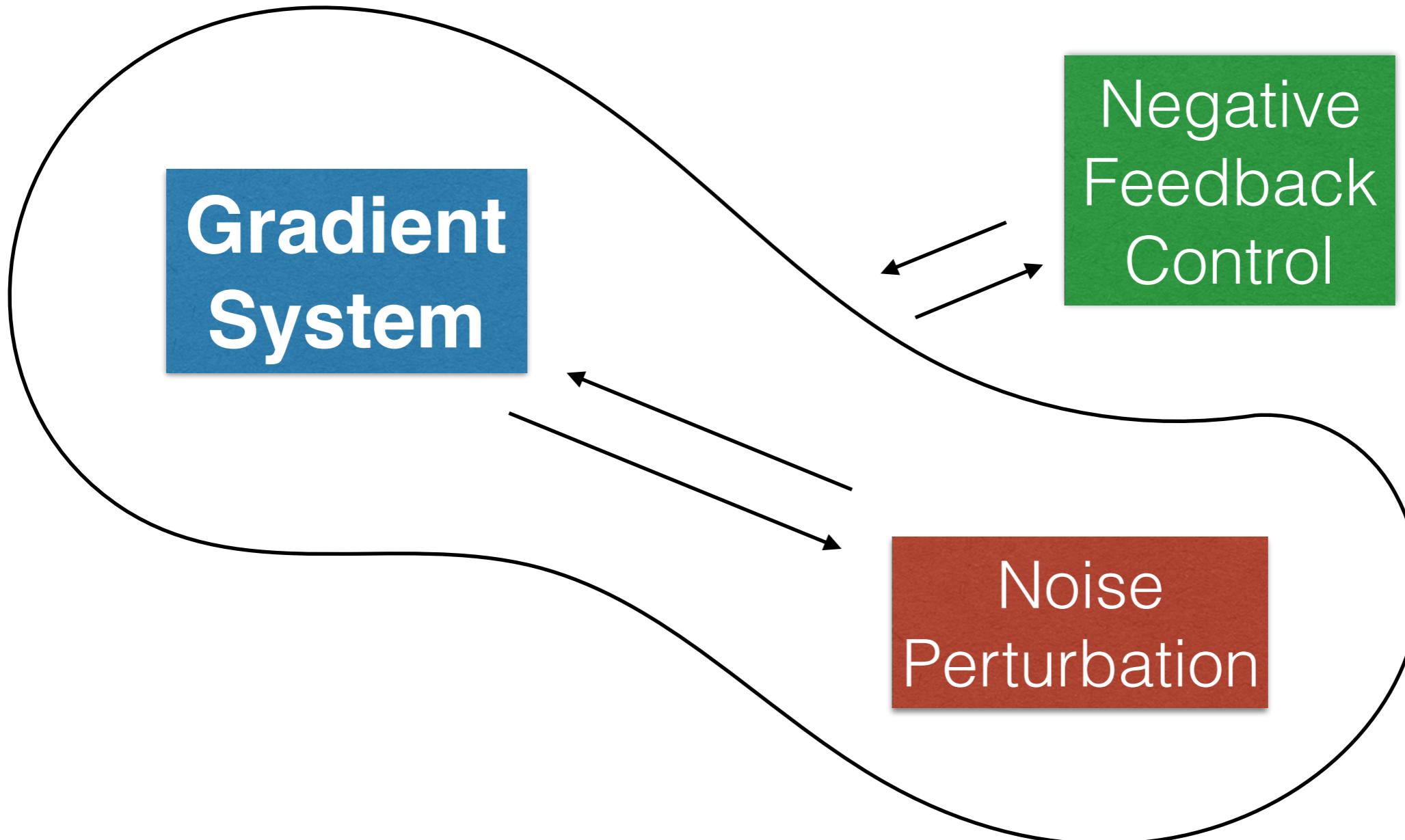
# Autocorrelation functions (LJ System)



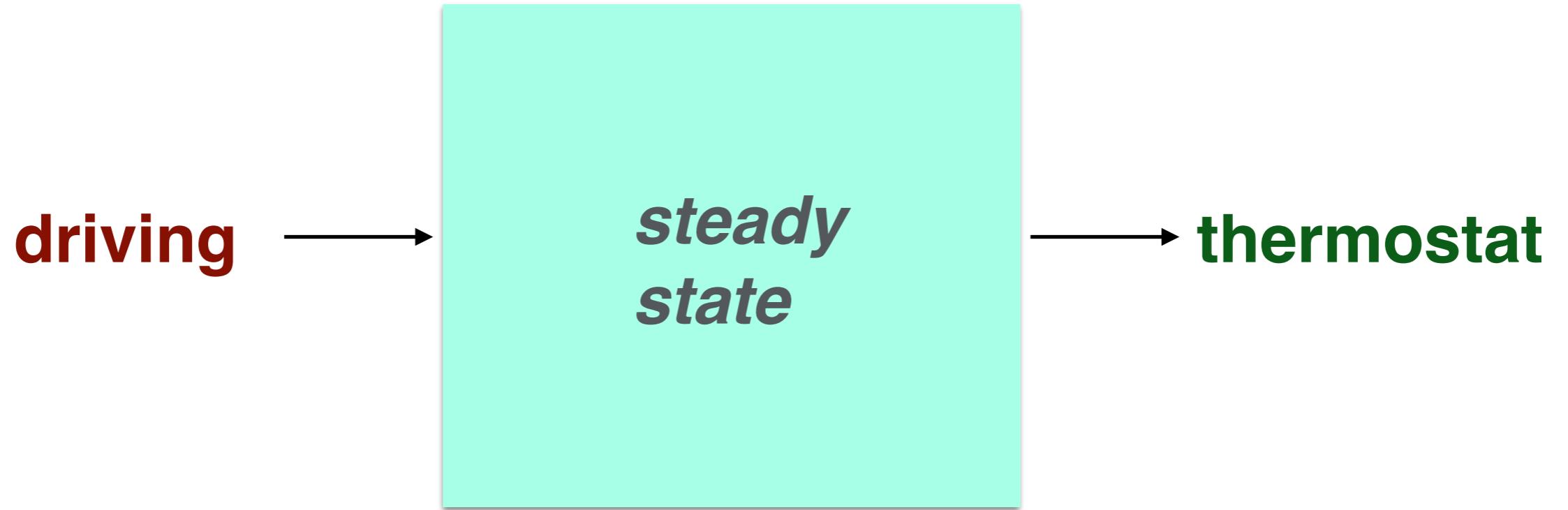
Even the deterministic method seems to work for complicated systems.

# Adaptive Thermostats

Jones & L., J. Chemical Physics, 2011



**Use negative feedback loop  
control to stabilize the system against  
force perturbation (even unknown)**



$$\dot{q} = \frac{p}{m}, \quad \dot{p} = F(q) - \xi p + \sigma_{\text{heat}} \dot{w}_{\text{heat}}, \quad \dot{\xi} = \left[ \frac{p^2}{m} - kT \right] / Q$$

$$\mathcal{L}^\dagger (\rho_\beta(q, p) \times \hat{\rho}(\xi)) = 0, \quad \hat{\rho}(\xi) = \exp(-\beta Q(\xi - \xi_{\text{heat}})^2/2)$$

& ergodic for this equilibrium distribution

# Bayesian Sampling

Understand choice of parameters  $q$  given observations  $X$

$$X = \{x_1, x_2, \dots, x_N\}$$

Posterior probability density (**from Bayes' Theorem**):

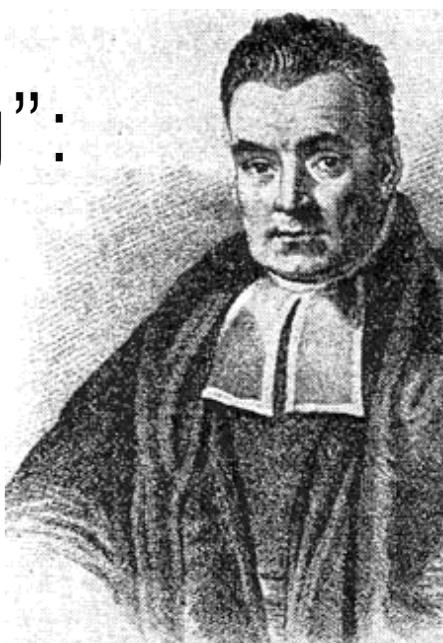
$$p(X|q)p(q) = p(q|X)$$

$$p(q|X) \propto \exp(-U(q)), \quad U(q) = -\log p(X|q) - \log p(q)$$

*model*                    *prior*

Use **Maximum Likelihood Estimate** / “Subsampling”:

$$\log p(X|q) \approx \frac{N}{\tilde{N}} \sum_{i=1}^{\tilde{N}} \log p(x_i|q) \quad \tilde{N} \ll N$$



Discipuli Domini Salini Druimmond qui vigesimo-septimo die  
 Februario MDCCXIX subscripserunt. 1419

Arch	Bennet	Alex Prokat
Geo.	Carruthers	David Lindsay
GEO:	Gordon	Geo. Doug. 2
Geo.	McLennan	Gul. Taylor
W.	Preston	Ja: Barclay
Gul.	Horsburgh	Jo: Gilmurst
Hen:	Ker	Jo: Horsley
Ioan:	Boston	John Patoun
Jo: Carruthers		John Rusker son
Ioan Morison		Jo: Smith
John Paxton		Jo: Thomson
Jo Faell		Isa. Maddox
Mich: Robertson		Rob: Cleland
Pat Mardock		Rob: Douglass
Simon Elliot		Rob: Richardson
Thomas Carmichael		Th. Bayes
		Tho: Morrison

1719

Th. Bayes

# Methods for Gibbs sampling with a noisy gradient

- Ignore the perturbation
- Estimate the perturbation/correct for it
- SGLD (Langevin with a diminishing stepsize sequence)
- Adaptive thermostats Ad-L, Ad-NH, ...

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

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

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

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

1. Stepsize-dependent dynamics (like in B.E.A.)
2. Distorts temperature
3. Easy to correct - if we know  $\sigma$
4. Computing/estimating  $\sigma$  can be difficult in practice

Options:

Monte-Carlo based approach [[Ceperley et al, ‘Quantum Monte Carlo’ 1999](#)]

Stochastic Gradient Langevin Dynamics [[Welling, Teh, 2011](#)]

Adaptive Thermostat [[Jones and L., 2011](#)]

# The Adaptive Property

*Jones & L. 2011*

Applying Nosé-Hoover Dynamics to a system which is driven by white noise restores the canonical distribution.

## Adaptive (Automatic) Langevin

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

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

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

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

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

# Discretization

[With X. Shang, 2015]

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**

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

typically anticipate 2nd order (IM)

# 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 also implies large friction so restricted phase space exploration** (even if better accuracy). So caution is needed...

$$\mathcal{L}^\dagger(f_2\tilde{\rho})=\mathcal{L}_2^\dagger\tilde{\rho}$$

$$\mathcal{L}^\dagger = -p \partial_q + U'(x) \partial_p + \xi \partial_p(p\cdot) + \frac{\gamma}{\beta} \partial_{pp} - \frac{1}{\mu}(p^2-\beta^{-1}) \partial_\xi$$

$$\varepsilon=1/\hat\gamma=1/\mu$$

$$\left(\bar{\mathcal{L}}_{\rm O}^\dagger+\varepsilon\mathcal{L}_{\rm H}^\dagger\right)\left(\hat{f}_{2,0}+\varepsilon\hat{f}_{2,1}+O(\varepsilon^2)\right)\rho_\beta=-\varepsilon\mathcal{P}\mathcal{L}_2^\dagger\tilde{\rho}_\beta$$

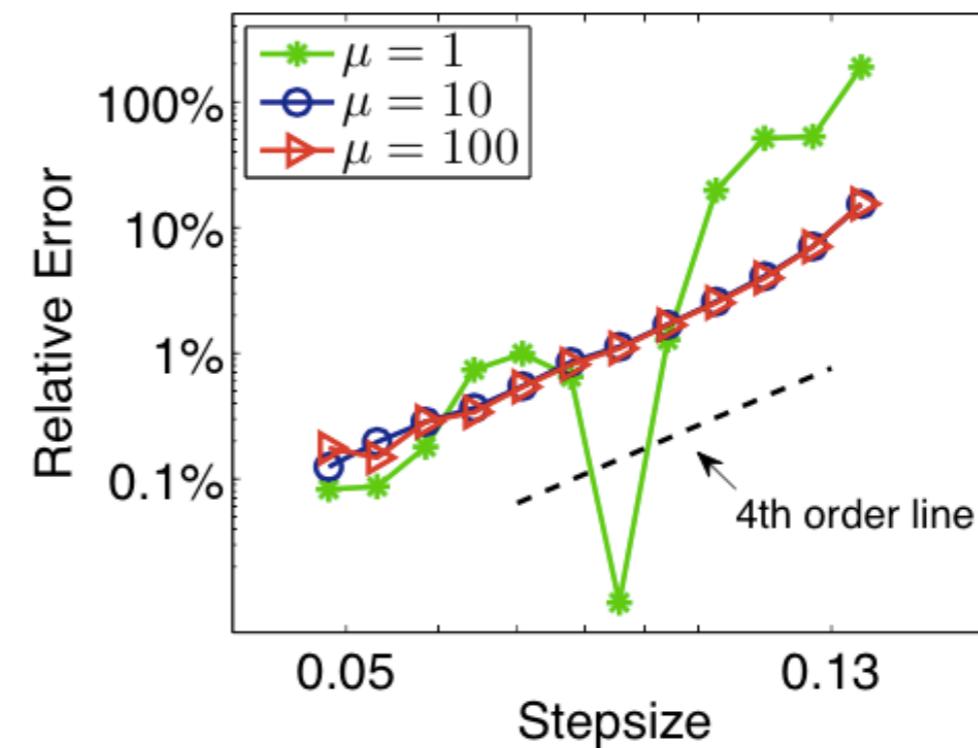
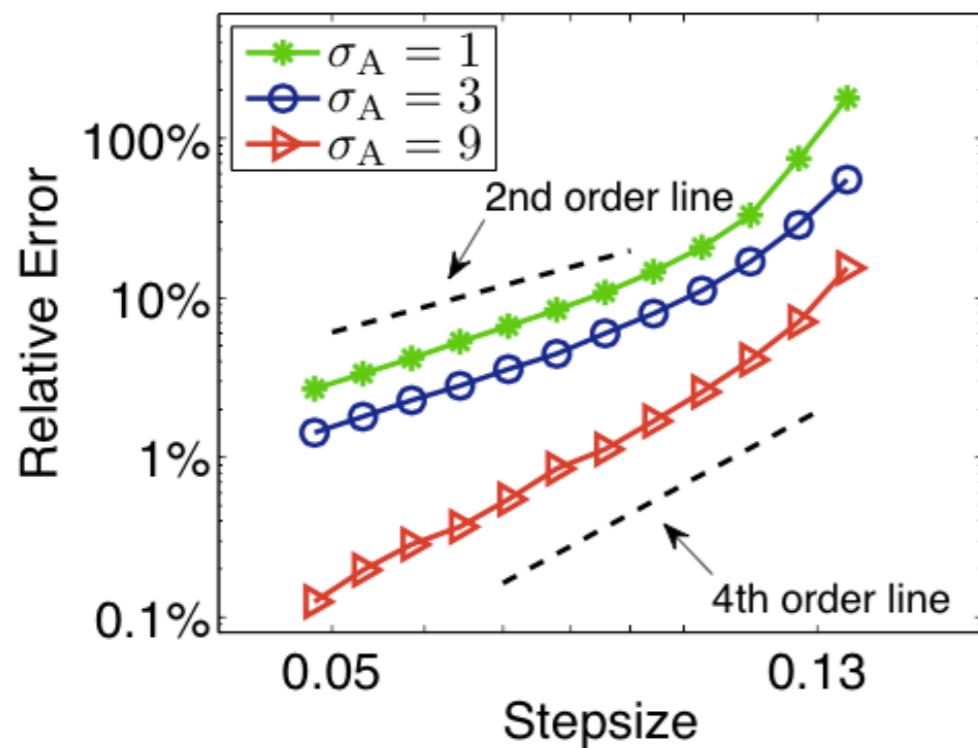
$$\mathcal{P}\mathcal{L}_2^\dagger\tilde{\rho}_\beta=\left(\frac{\beta}{12}\left[3pU'(x)U''(x)-p^3U'''(x)\right]+\frac{\hat{\gamma}}{12}\left[3U''(x)-3\beta p^2U''(x)+\frac{1}{\mu}\left(6\beta p^4-28p^2+10\beta^{-1}\right)\right]\right)\rho_\beta$$

$$\hat{f}_{2,0}\equiv\hat{f}_{2,0}^{\text{BADODAB}}=\frac{1}{8}\left(U''(x)-\beta p^2U''(x)\right)$$

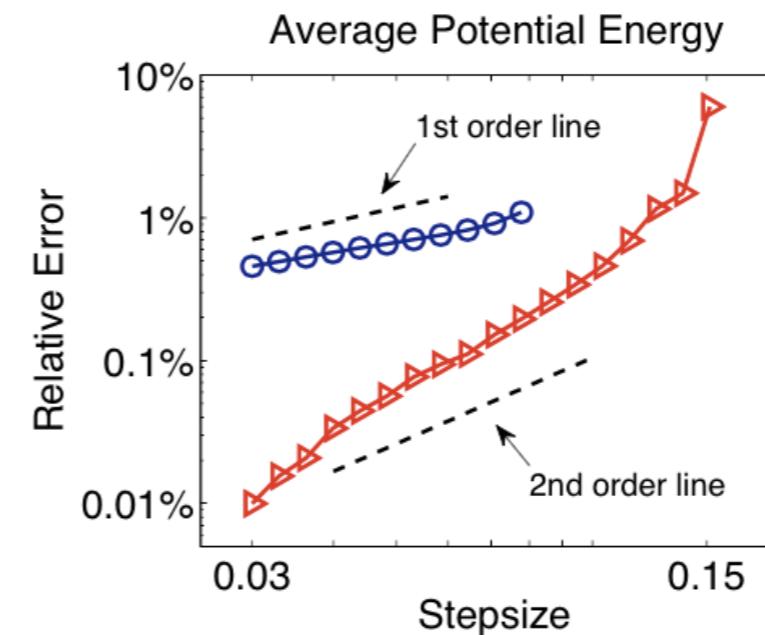
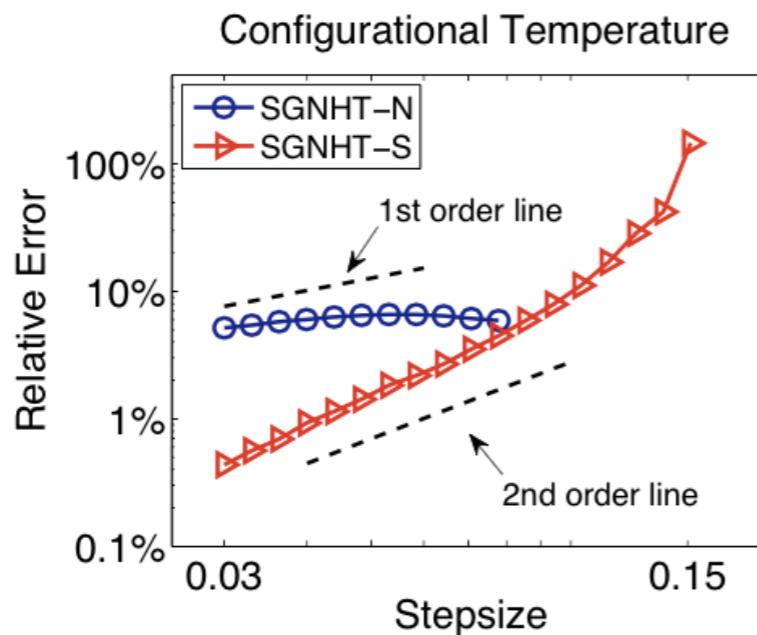
$$\langle \phi(x) \rangle_{\text{BADODAB}} = \langle \phi(x) \rangle + h^2 \langle \phi(x) \hat{f}_{2,0}^{\text{BADODAB}} \rangle + O(\varepsilon h^2 + h^4)$$

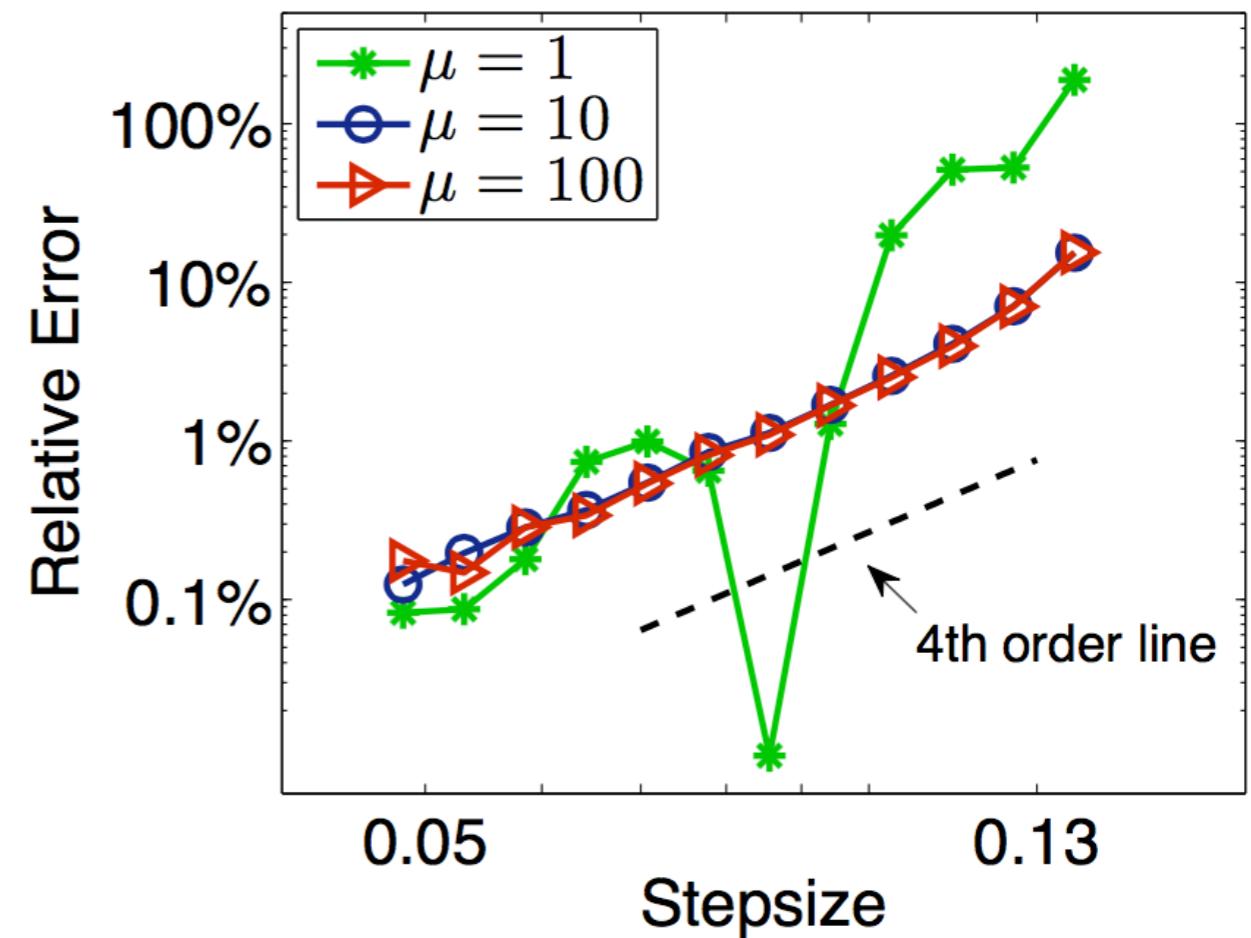
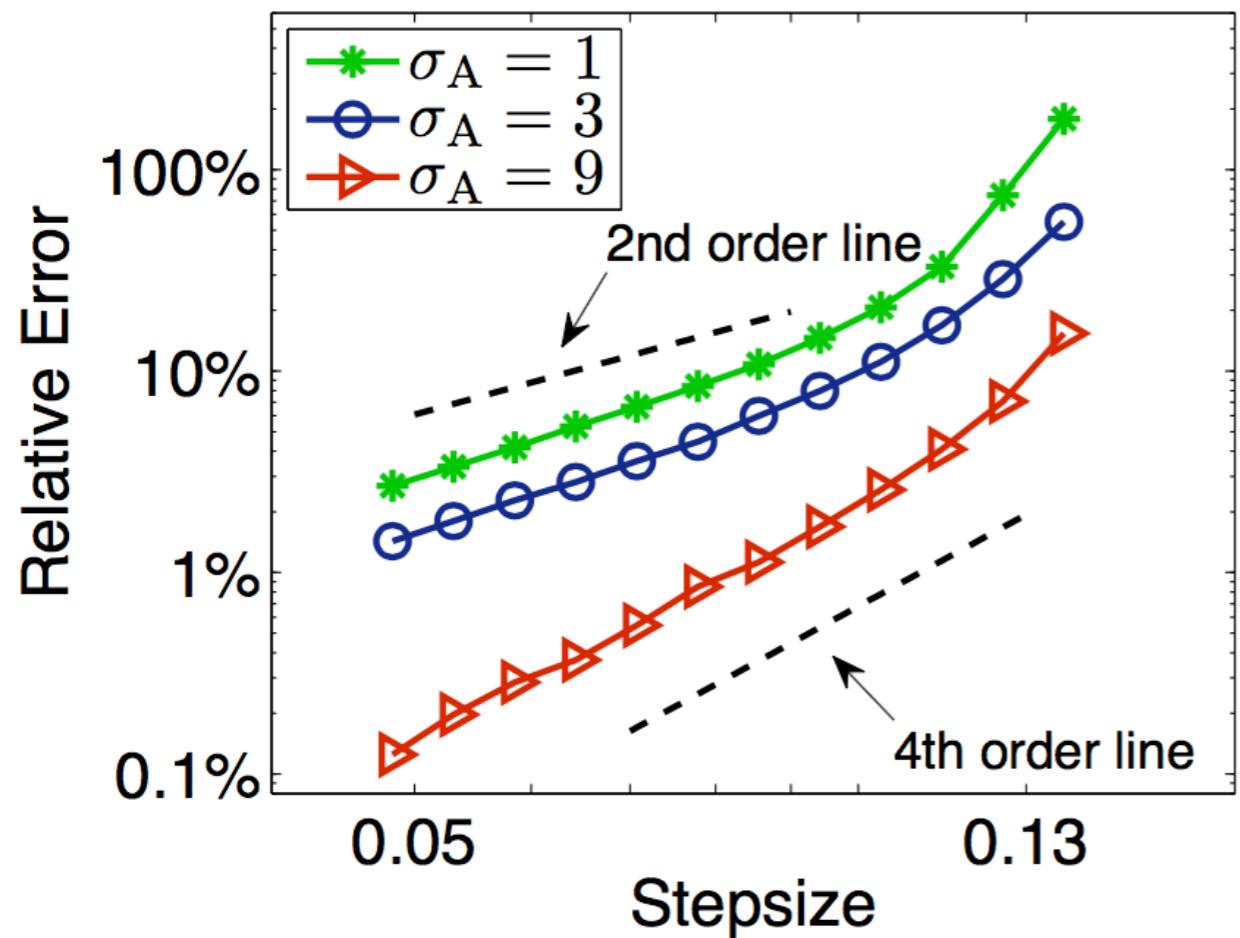
# 500 LJ particles, clean gradient

configurational temperature



Comparison with Chen et al.





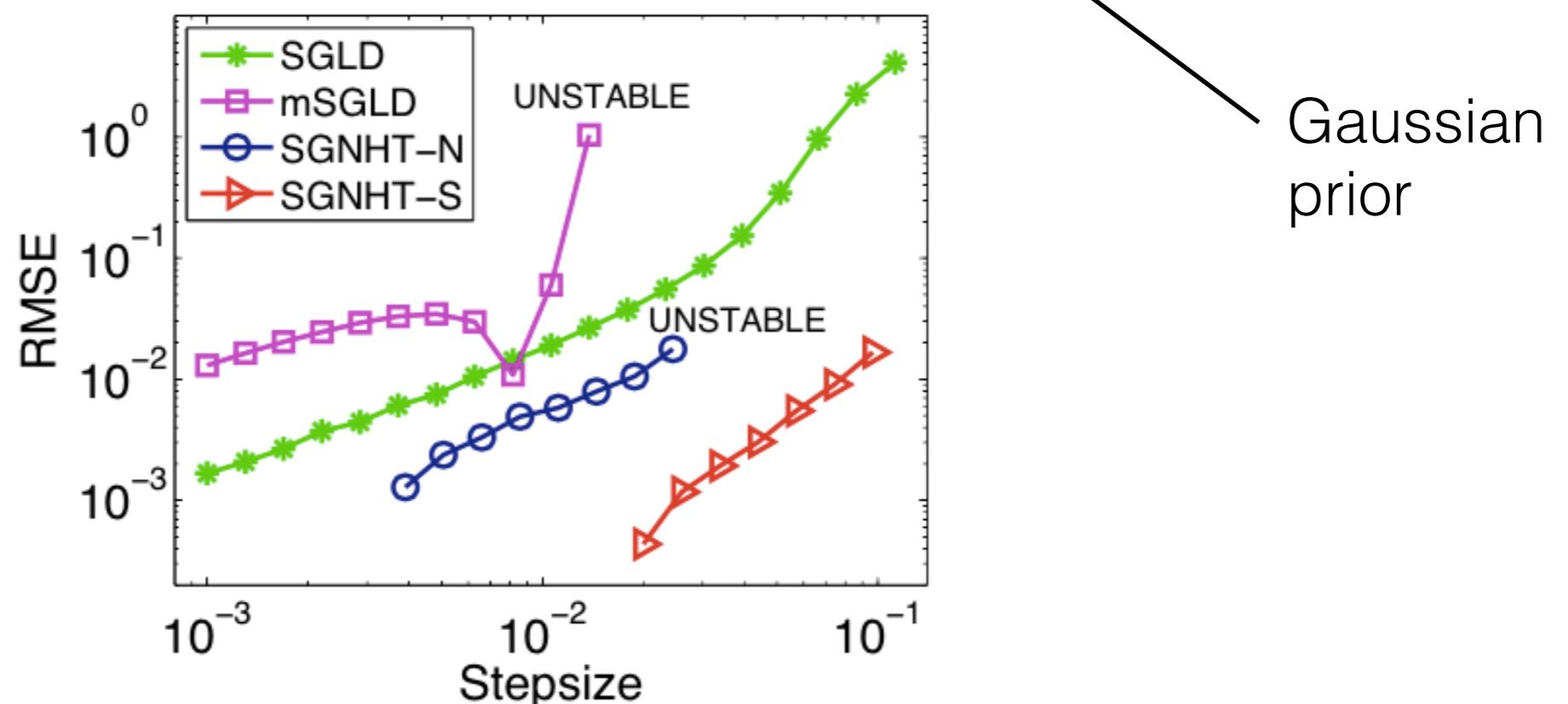
- Fourth order convergence to the invariant measure
- Large friction ( $\hat{\gamma} \propto \sigma_A^2$ ) and thermal mass ( $\mu$ ) limits
- Only one force calculation required at each step

# Bayesian Logistic Regression

$\pi(y_i|\mathbf{x}_i, \boldsymbol{\beta}) = f(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$      $f$ : logistic function

↑  
covariates e.g. age, income, ...  
  
data e.g. voting intention  
  
posterior parameter distribution

$$\pi(\boldsymbol{\beta}) \propto \exp\left(-\frac{1}{2}\|\boldsymbol{\beta}\|^2\right) \prod_{i=1}^N f(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$$



# 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?**

# Methods for Gibbs sampling with a noisy gradient

- Ignore the perturbation
- Estimate the perturbation/correct for it
- SGLD (Langevin with a diminishing stepsize sequence)
- Adaptive thermostats Ad-L, Ad-NH, ...

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

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

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

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

1. Stepsize-dependent dynamics (like in B.E.A.)
2. Distorts temperature
3. Easy to correct - if we know  $\sigma$
4. Computing/estimating  $\sigma$  can be difficult in practice

Options:

Monte-Carlo based approach [Ceperley et al, ‘Quantum Monte Carlo’ 1999]

Stochastic Gradient Langevin Dynamics [Welling, Teh, 2011]

Adaptive Thermostat [Jones and L., 2011]

# The Adaptive Property

*Jones & L. 2011*

Applying Nosé-Hoover Dynamics to a system which is driven by white noise restores the canonical distribution.

## Adaptive (Automatic) Langevin

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

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

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

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

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

# Discretization

[With X. Shang, 2015]

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**

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

typically anticipate 2nd order (IM)

# 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 also implies large friction so restricted phase space exploration** (even if better accuracy). So caution is needed...

$$\mathcal{L}^\dagger(f_2\tilde{\rho})=\mathcal{L}_2^\dagger\tilde{\rho}$$

$$\mathcal{L}^\dagger = -p \partial_q + U'(x) \partial_p + \xi \partial_p(p\cdot) + \frac{\gamma}{\beta} \partial_{pp} - \frac{1}{\mu}(p^2-\beta^{-1}) \partial_\xi$$

$$\varepsilon=1/\hat\gamma=1/\mu$$

$$\left(\bar{\mathcal{L}}_{\rm O}^\dagger+\varepsilon\mathcal{L}_{\rm H}^\dagger\right)\left(\hat{f}_{2,0}+\varepsilon\hat{f}_{2,1}+O(\varepsilon^2)\right)\rho_\beta=-\varepsilon\mathcal{P}\mathcal{L}_2^\dagger\tilde{\rho}_\beta$$

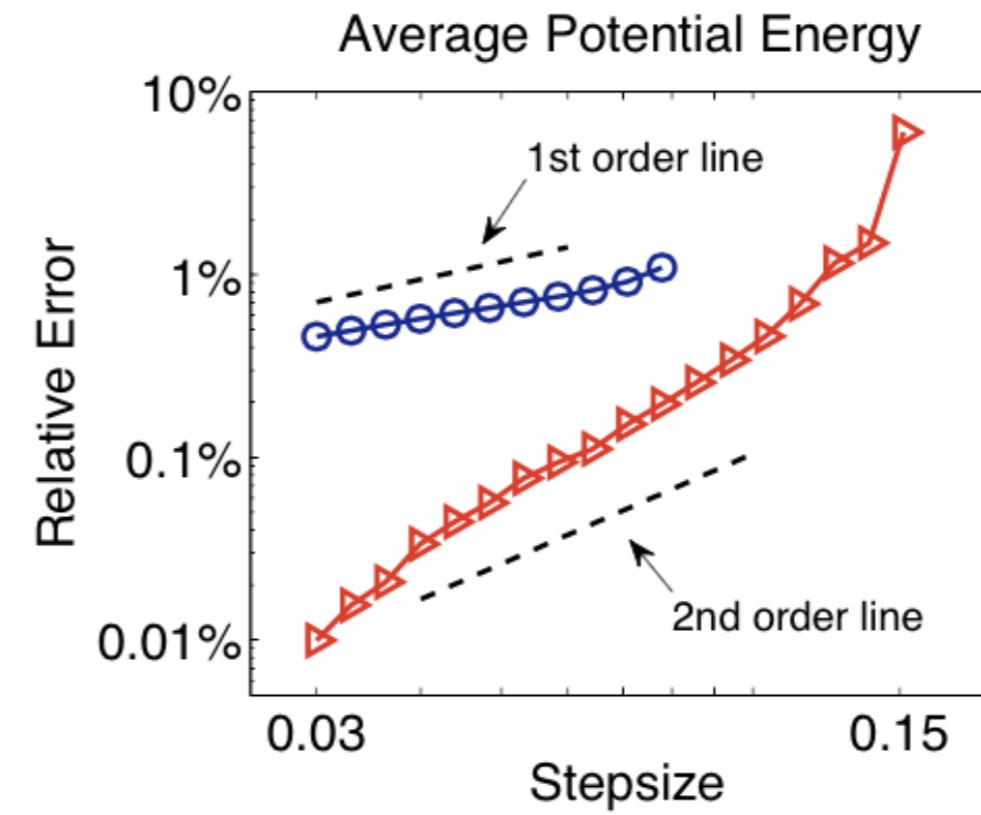
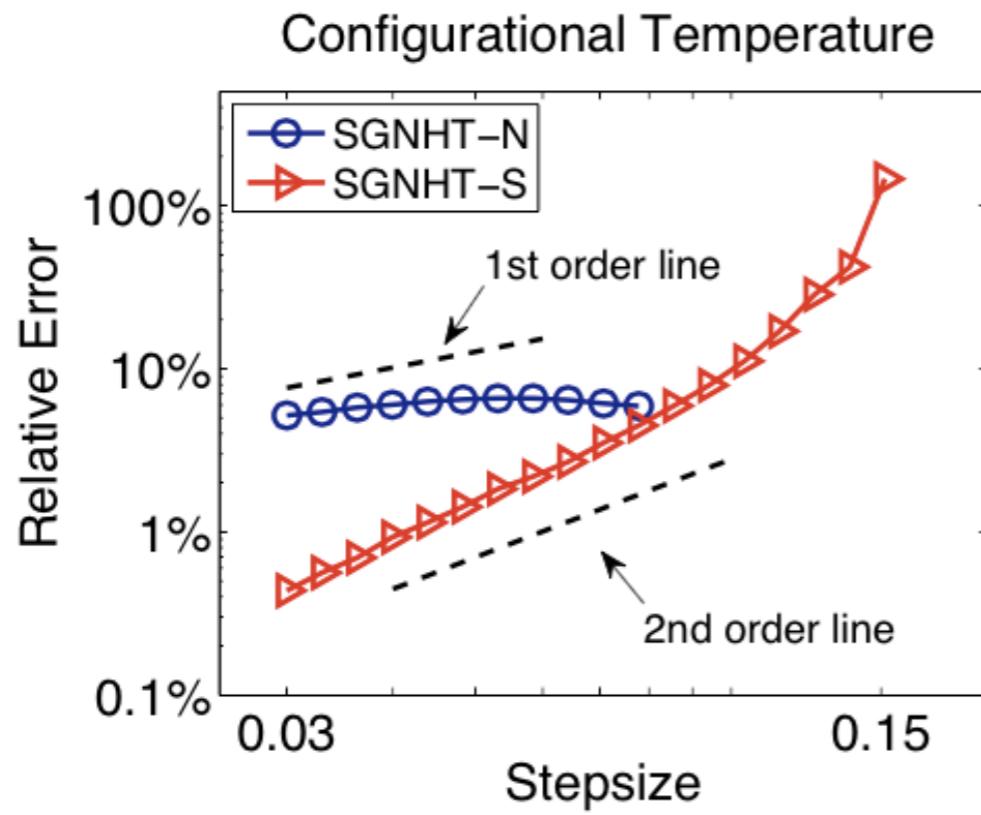
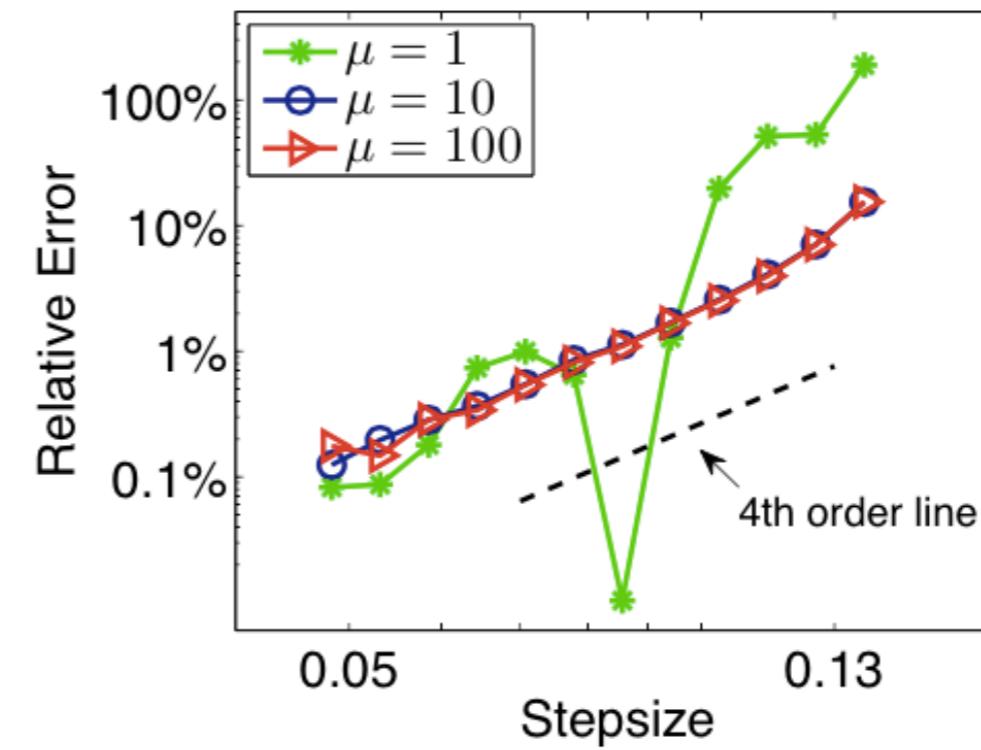
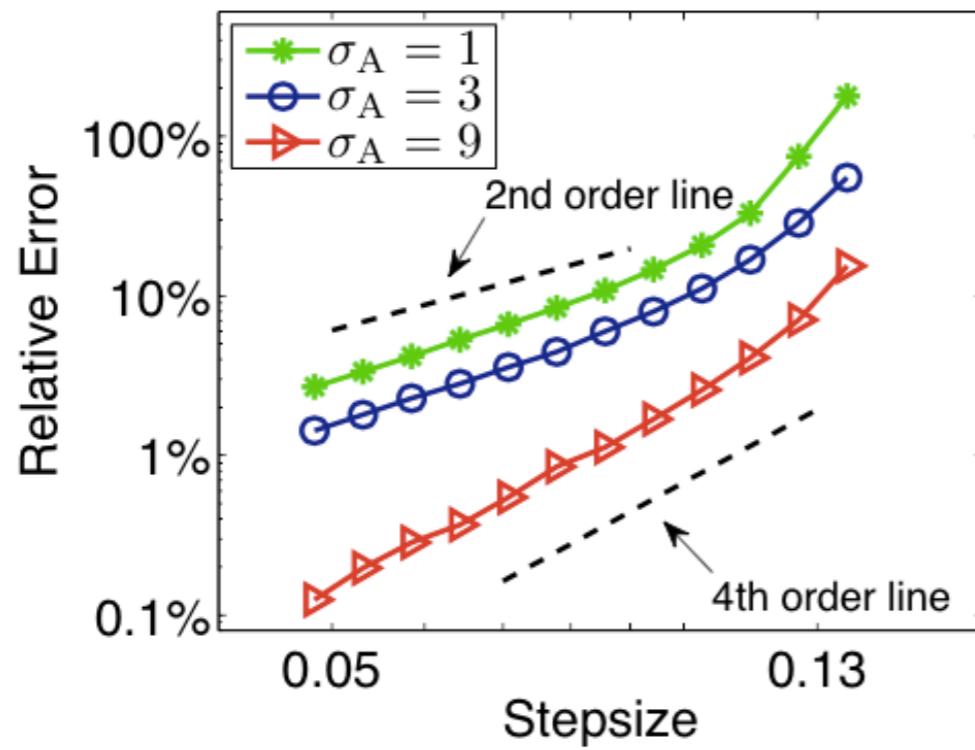
$$\mathcal{P}\mathcal{L}_2^\dagger\tilde{\rho}_\beta=\left(\frac{\beta}{12}\left[3pU'(x)U''(x)-p^3U'''(x)\right]+\frac{\hat{\gamma}}{12}\left[3U''(x)-3\beta p^2U''(x)+\frac{1}{\mu}\left(6\beta p^4-28p^2+10\beta^{-1}\right)\right]\right)\rho_\beta$$

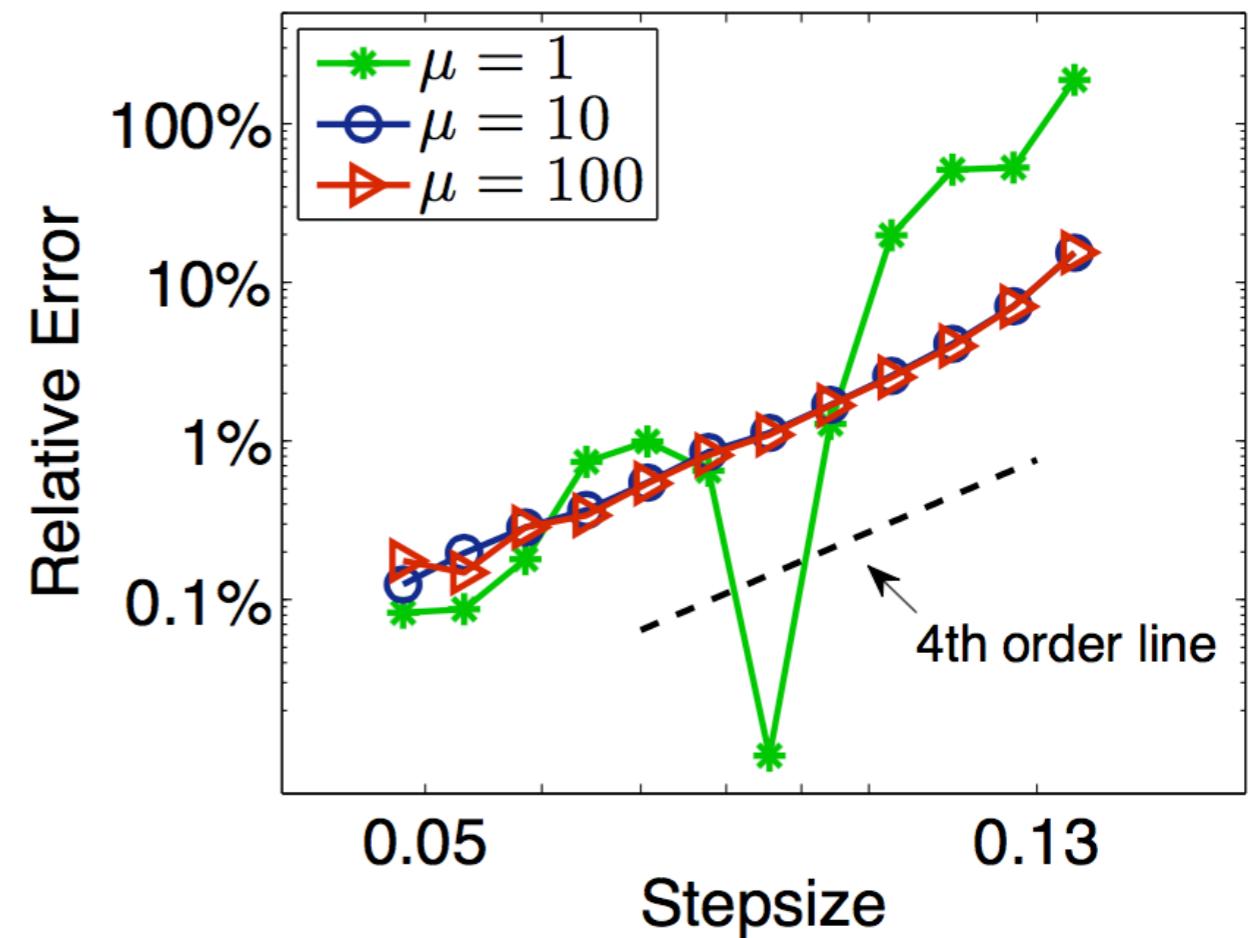
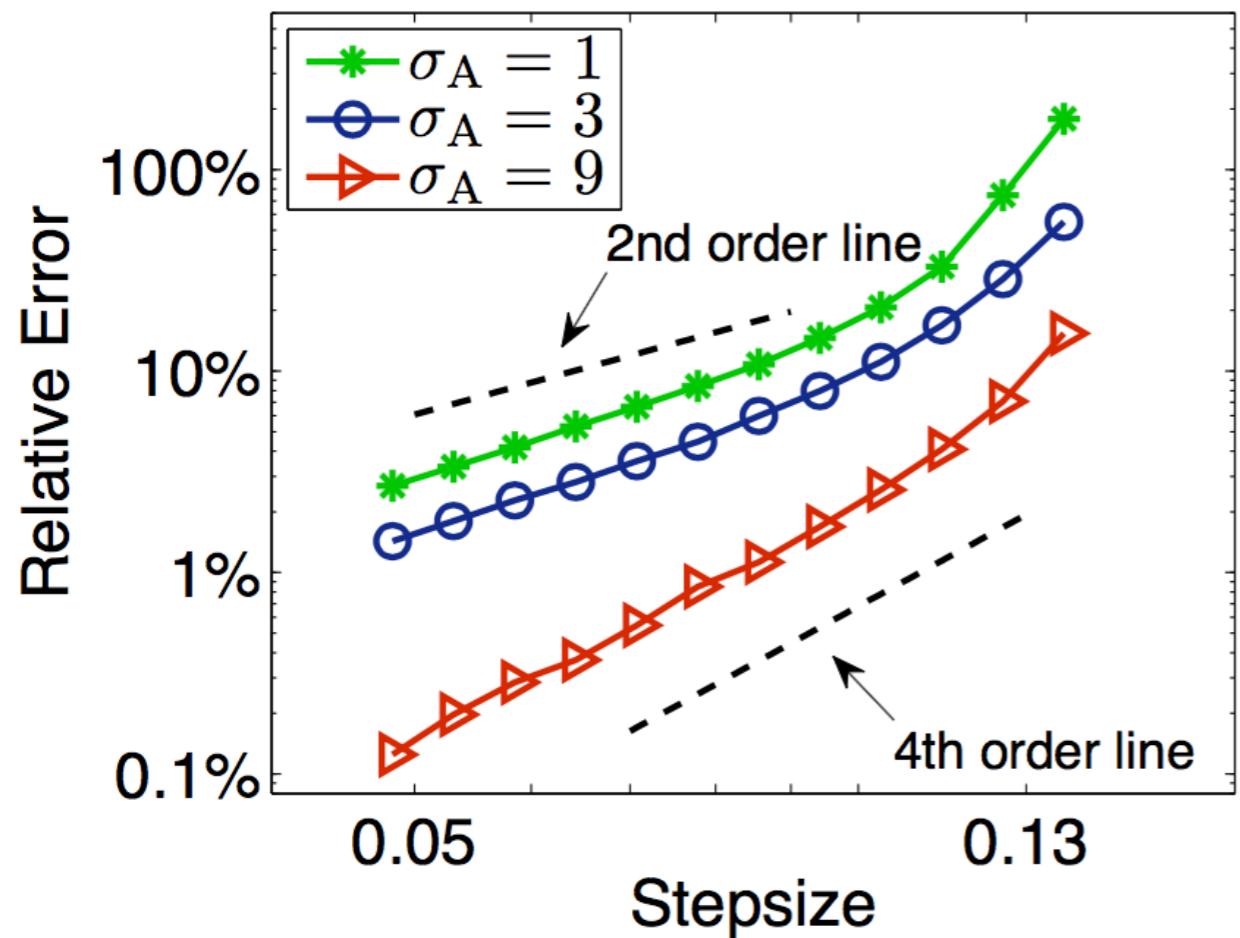
$$\hat{f}_{2,0}\equiv\hat{f}_{2,0}^{\text{BADODAB}}=\frac{1}{8}\left(U''(x)-\beta p^2U''(x)\right)$$

$$\langle \phi(x) \rangle_{\text{BADODAB}} = \langle \phi(x) \rangle + h^2 \langle \phi(x) \hat{f}_{2,0}^{\text{BADODAB}} \rangle + O(\varepsilon h^2 + h^4)$$

# 500 LJ particles, clean gradient

## configurational temperature





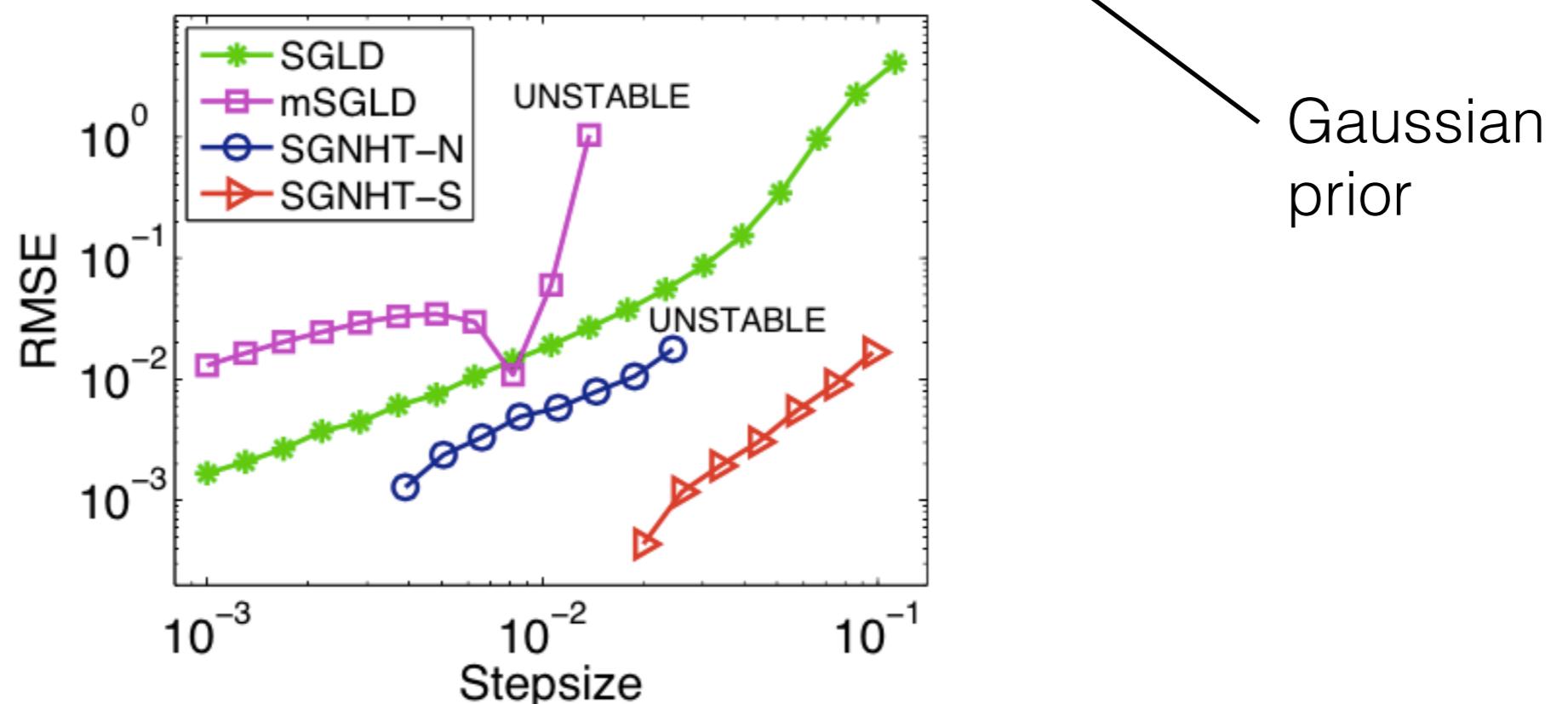
- Fourth order convergence to the invariant measure
- Large friction ( $\hat{\gamma} \propto \sigma_A^2$ ) and thermal mass ( $\mu$ ) limits
- Only one force calculation required at each step

# Bayesian Logistic Regression

$\pi(y_i|\mathbf{x}_i, \boldsymbol{\beta}) = f(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$      $f$ : logistic function

↑  
covariates e.g. age, income, ...  
  
data e.g. voting intention  
  
posterior parameter distribution

$$\pi(\boldsymbol{\beta}) \propto \exp\left(-\frac{1}{2}\|\boldsymbol{\beta}\|^2\right) \prod_{i=1}^N f(y_i \boldsymbol{\beta}^T \mathbf{x}_i)$$



# Covariance-Controlled Adaptive Langevin Dynamics

**Shang, Zhu, Leimkuhler and Storkey, NIPS 2015**

In the typical case, the noise may have a multivariate Gaussian distribution but with unknown (and evolving) covariance.

If we assume that we can obtain a covariance estimator then we can use this to enhance the accuracy of the SDEs.

**CCAdL=**

“Covariance Controlled Adaptive Langevin Dynamics” incorporates such a correction term together with an adaptive Langevin thermostat...

- Formulation

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

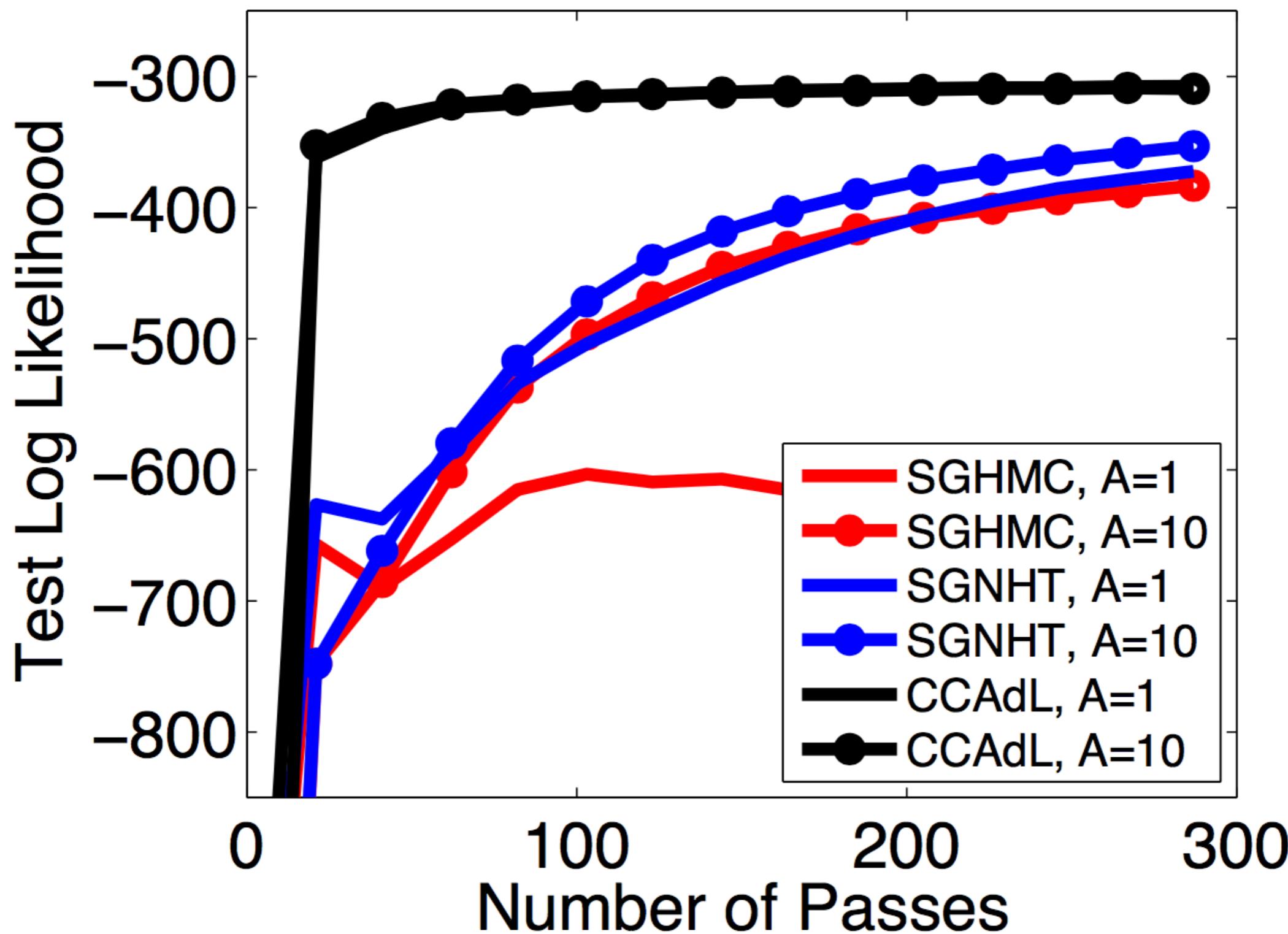
$$\begin{aligned} d\mathbf{p} = & -\nabla U(\mathbf{q})dt + \sqrt{h\Sigma(\mathbf{q})}\mathbf{M}^{1/2}d\mathbf{W} - (h/2)\beta\Sigma(\mathbf{q})\mathbf{p}dt \\ & - \xi\mathbf{p}dt + \sqrt{2\hat{\gamma}\beta^{-1}}\mathbf{M}^{1/2}d\mathbf{W}_A \end{aligned}$$

$$d\xi = \mu^{-1} [\mathbf{p}^T \mathbf{M}^{-1} \mathbf{p} - N_d k_B T] dt$$

with invariant distribution

$$\tilde{\rho}_\beta(\mathbf{q}, \mathbf{p}, \xi) \propto \exp(-\beta H(\mathbf{q}, \mathbf{p})) \exp\left(-\frac{\beta\mu}{2}(\xi - \hat{\gamma})^2\right)$$

- Parameter-dependent noise effectively dissipated by the additional covariance control term.



Binary classification of handwritten digits 7 and 9.

# Ensemble preconditioning MCMC simulation

**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)**

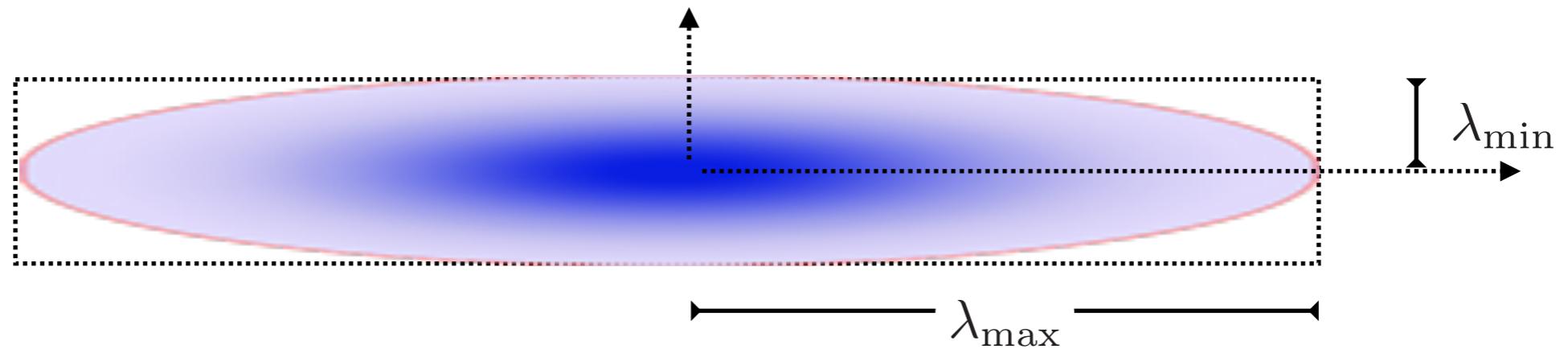
$$\tau_f = 1 + 2 \sum_{t=1}^{\infty} \text{cor}(f(x_t), f(x_0))$$

We would like to have  $\tau_f$  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:

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

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

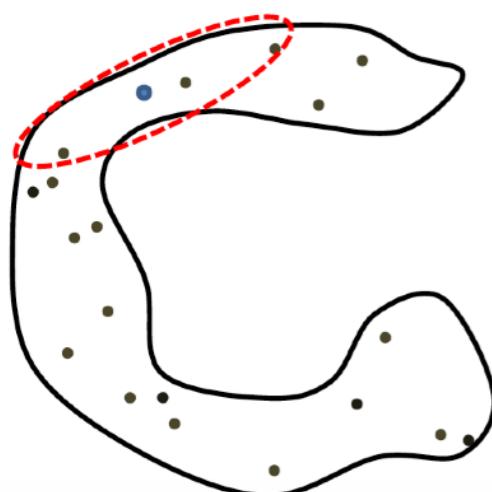
*each walker samples the  
same target distribution  $\pi(q)$*

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**

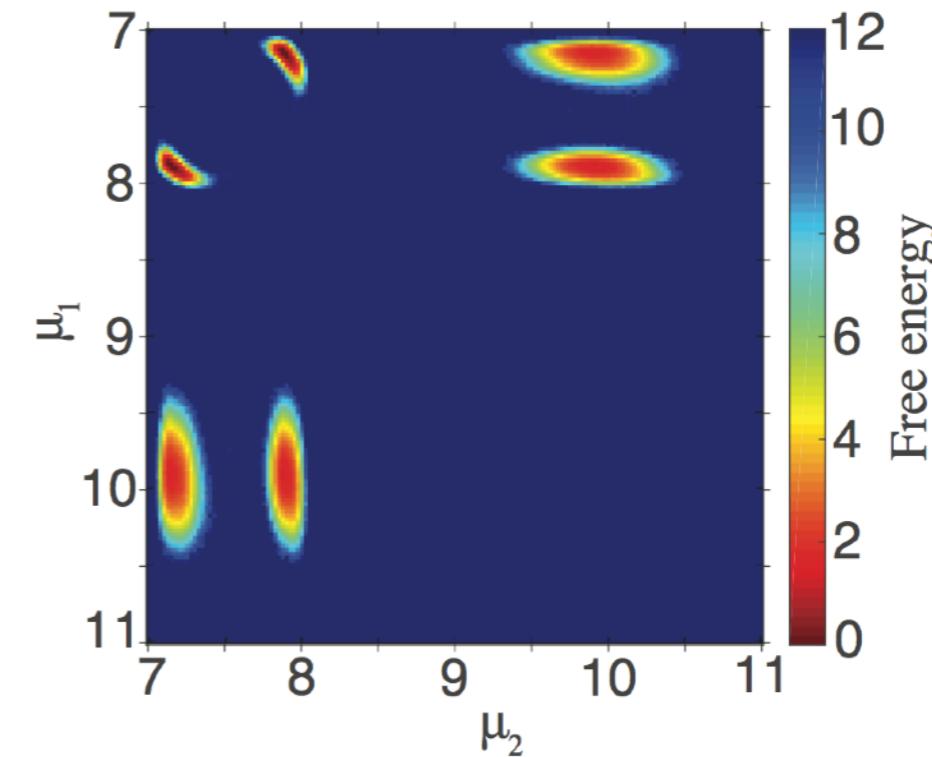
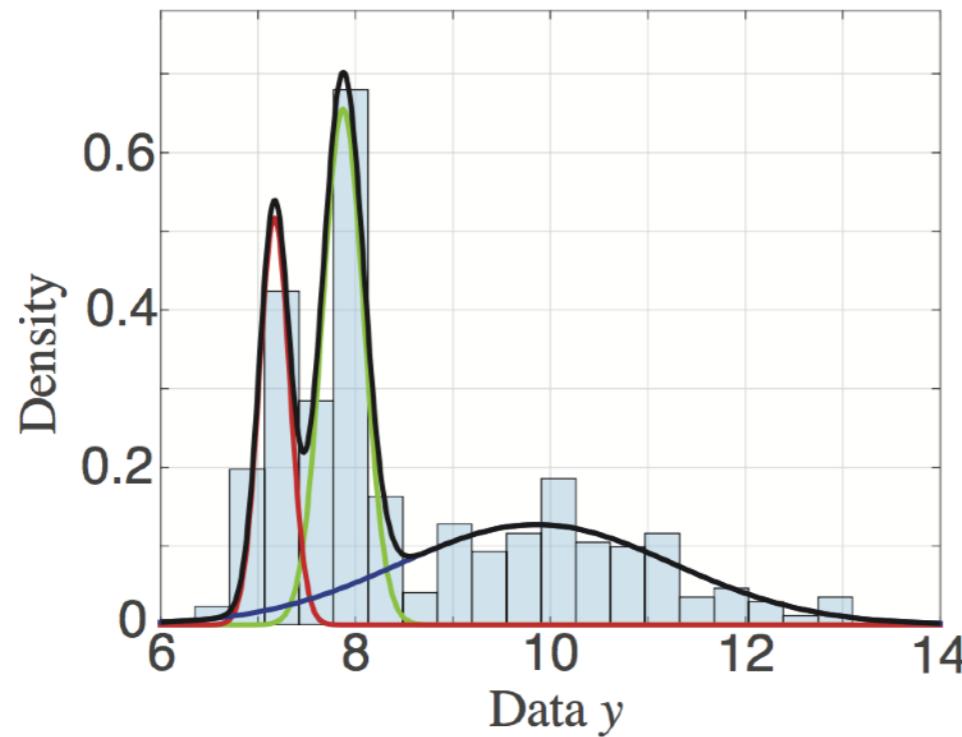
Discretization of SDEs: similar to **BAOAB**

# Gaussian Mixture Model: Hidalgo Stamps

## “Adventures in Stamp Collecting”

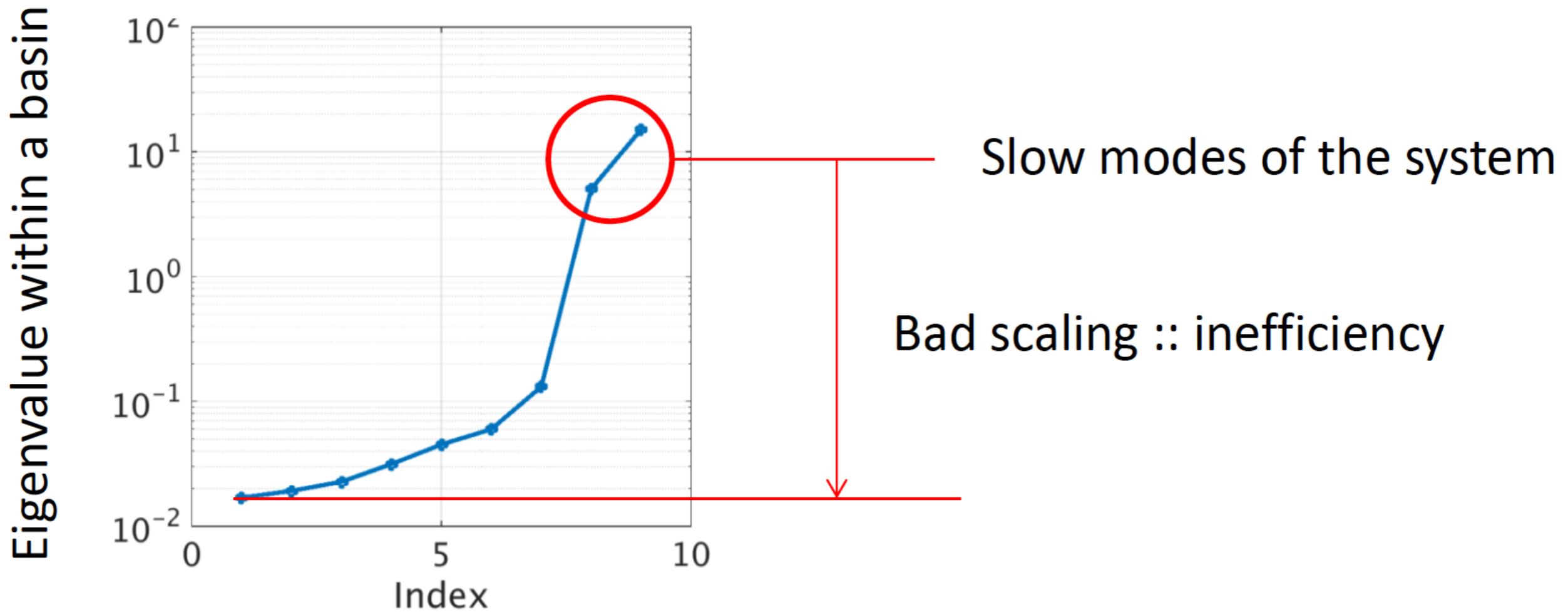
$$\sum_{n=1}^3 z_n \mathcal{N}(\mu_n, \lambda_n^{-2})$$

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:



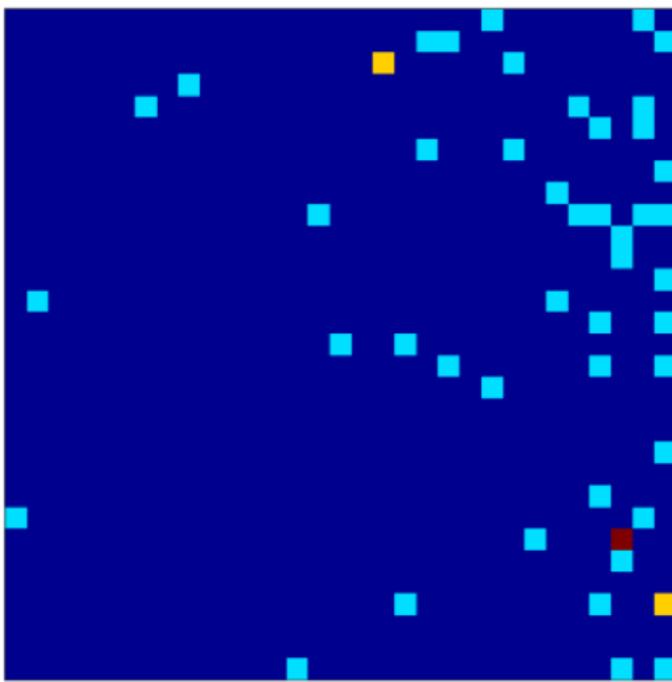
But the eigensystems are **different in different basins**, so the localized covariance is needed...

# Gaussian Mixture Model: Hidalgo Stamps

## Integrated Autocorrelation Times of Different Schemes

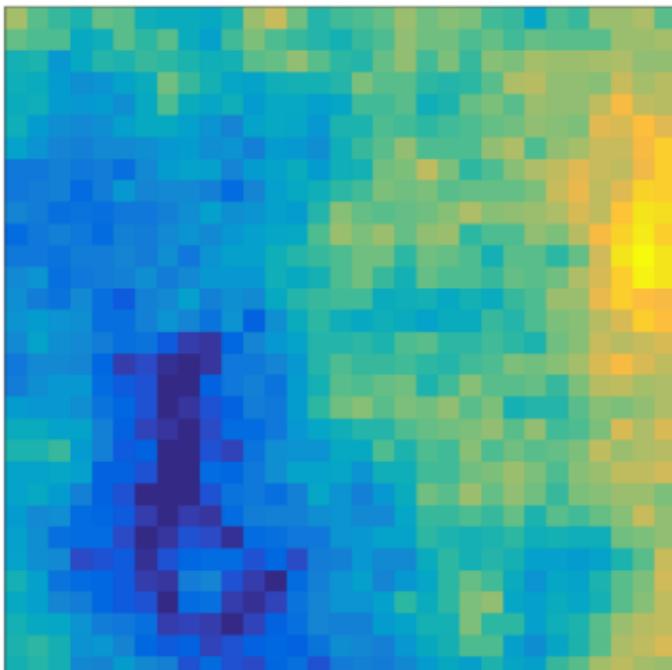
Scheme	$\min(z)$	$\max(\lambda)$	$\min(\mu)$	$\beta$
HMC	21495	<b>42935</b>	27452	7148
Langevin Dynamics	6825	<b>13279</b>	8384	4641
Ensemble Q-N	69	83	98	<b>115</b>

## Numerical test: Log-Gaussian Cox model



Observations  $X$

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



Break  $[0,1]^2$  into a  $32 \times 32$  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  $\mu$  and aim to infer likely  $Y$ , using  
hyperparameters  $\sigma^2, \beta$ , with prior  $\text{gamma}(2, \frac{1}{2})$

# Log Gaussian Cox Model

Scheme	$x$	$\sigma^2$	$\beta$	Efficiency
HMC	800.7	1041.6	<b>1318.7</b>	1.0
RMHMC	2158.9	34.0	<b>1502.0</b>	0.15
LD	405.1	140.6	<b>435.3</b>	3.5
... (no Metropolis)	81.6	20.5	<b>136.5</b>	11.2
EQN	71.9	49.2	<b>239.5</b>	5.4
... (no Metropolis)	64.4	8.8	<b>47.8</b>	26.8

**Ensemble Quasi-Newton python package**

[http://bitbucket.org/c\\_matthews/ensembleqn](http://bitbucket.org/c_matthews/ensembleqn)