

From Physical Modelling to Big Data Analytics: Examples and Challenges

Ben Leimkuhler
University of Edinburgh



RULE

EPSRC/NSF

ERC

ATI Summit on Big Data in the Physical Sciences

Outline

Part I: **Challenges**

Ideas from the Alan Turing Institute Scoping Workshop
“The Challenges of Data Intensive and Extreme Scale
Numerical Simulation in Physics, Materials Science and
Chemistry” (Jan 5-6, 2016, British Library)

Part II: **Example**

Simulations of physical states using noisy forces



Inference and machine learning algorithms

I. Challenges

ATI Scoping Workshops

researchers



sticky notes



facilitator



~30 white papers



ATI Scoping Workshops

Examples (particularly Physical Sciences relevant)

Theoretical and Computational Approaches to Large Scale Inverse Problems

Partial Differential Equations for Modelling, Analysing and Simulating Data Rich Phenomena

Big Data in Geoscience

The Challenges of Data Intensive and Extreme Scale Numerical Simulation in Physics, Materials Science and Chemistry

The Challenges of Data Intensive and Extreme Scale Numerical Simulation in Physics, Materials Science and Chemistry

Main Organizers: [Gabor Csanyi \(Cambridge, Engineering\)](#), [Detlef Hohl \(Shell\)](#), [Stephen Jarvis \(Warwick, Computer Science\)](#), [Ben Leimkuhler \(Edinburgh, Mathematics\)](#), [Mark Parsons \(Edinburgh, EPCC\)](#)

Themes

Extreme scale data-computing (exascale)
Numerics for data science
Data-centric materials and chemistry modelling

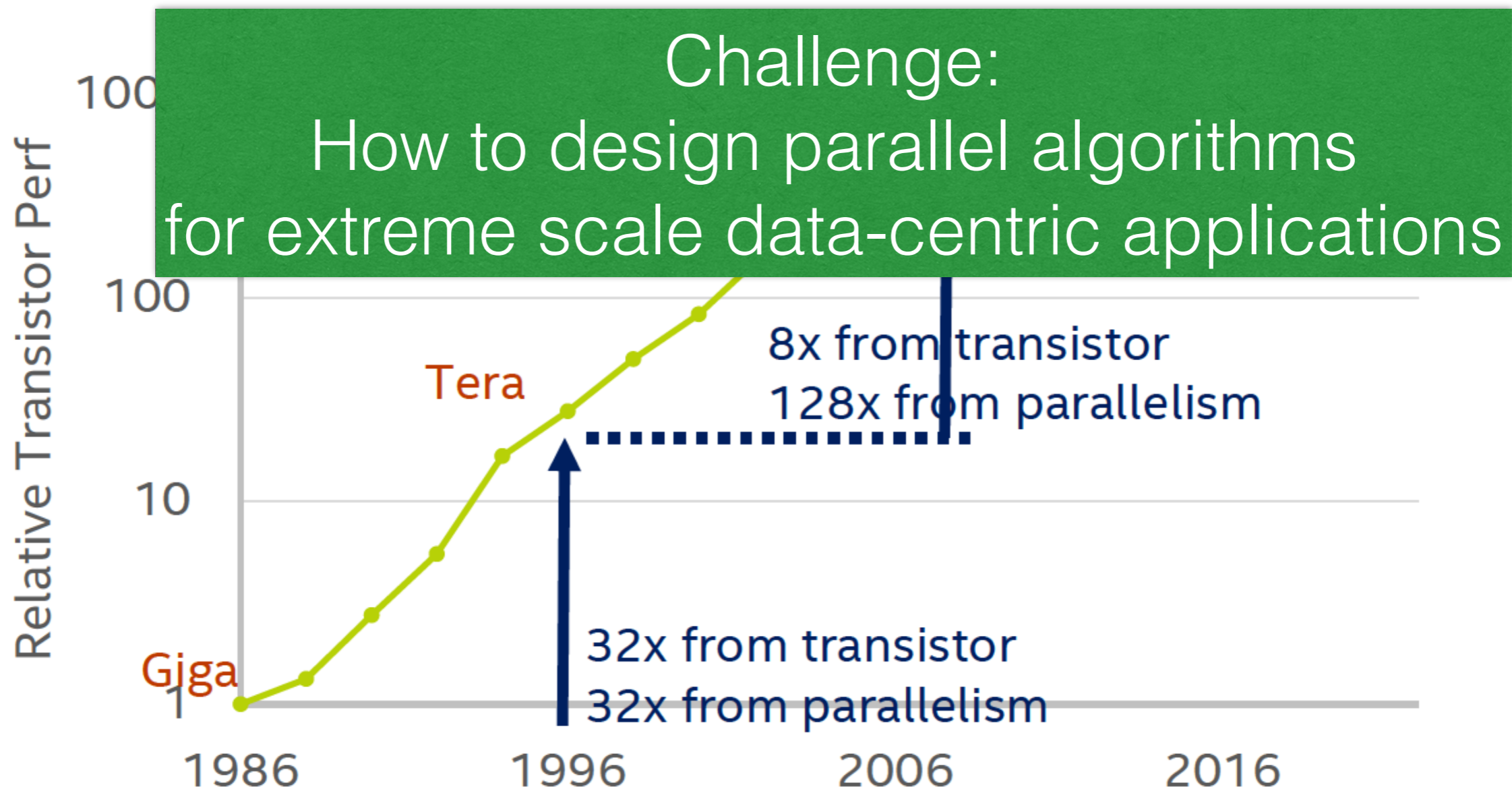
Participants

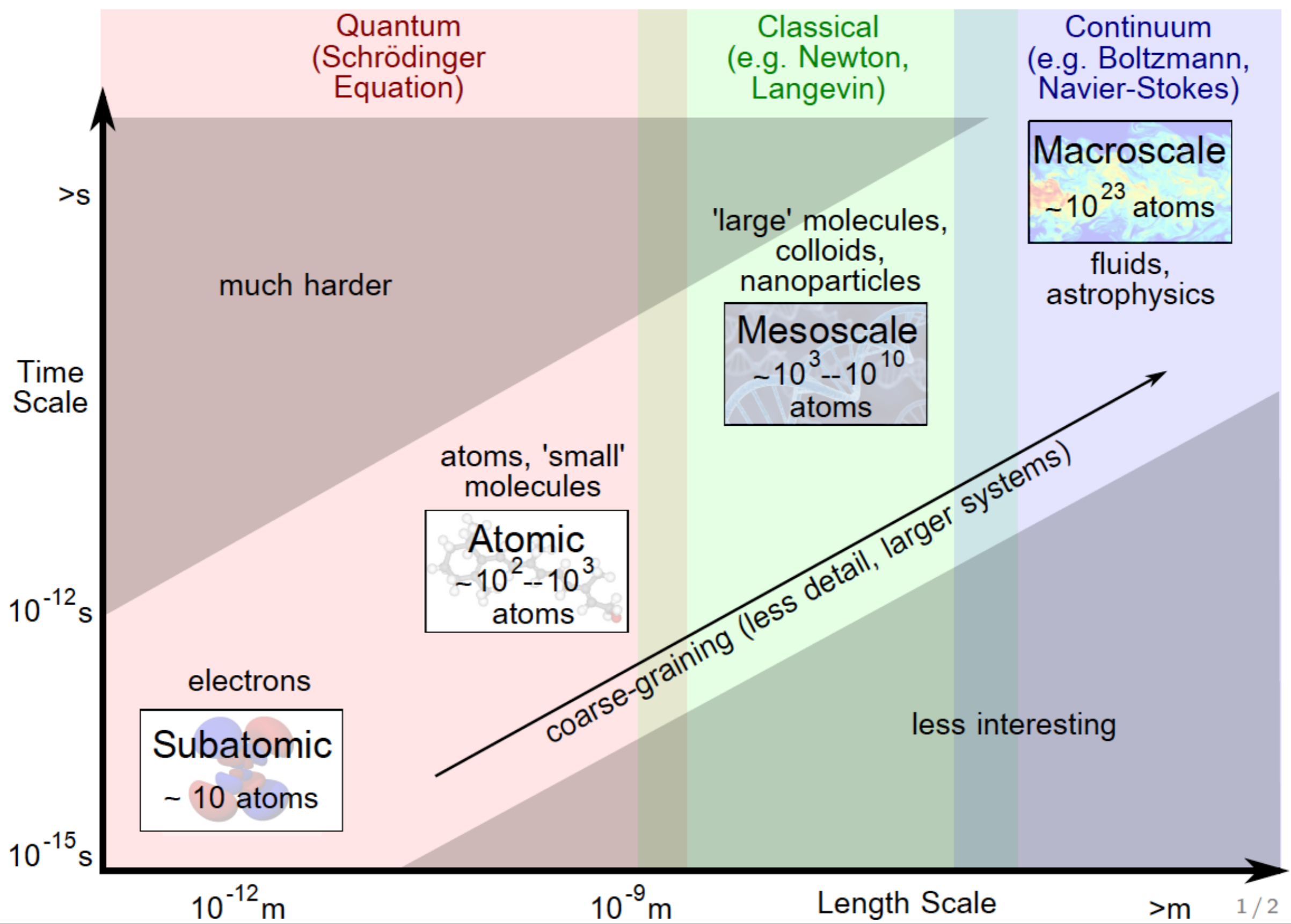
HPC scientific simulation community
Parallel computing experts
Mathematicians
Data scientists
Physical modellers (esp. materials/fluids)
Industry: Intel, Shell, Dassault (Biovia), Rolls Royce, ..

Increasing reliance on parallelism for HPC gains (as opposed to improvements from transistor) & Importance of energy considerations

Intel Exascale Labs

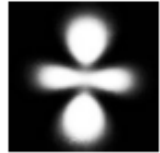
Implications to HPC Roadmap





Machine Learning Quantum Mechanics

Gabor Csanyi (Cambridge)



First principles simulation is extremely successful in materials science and chemistry

Traditionally $O(N^3)$ or worse
100 atoms \sim 100 CPU hours

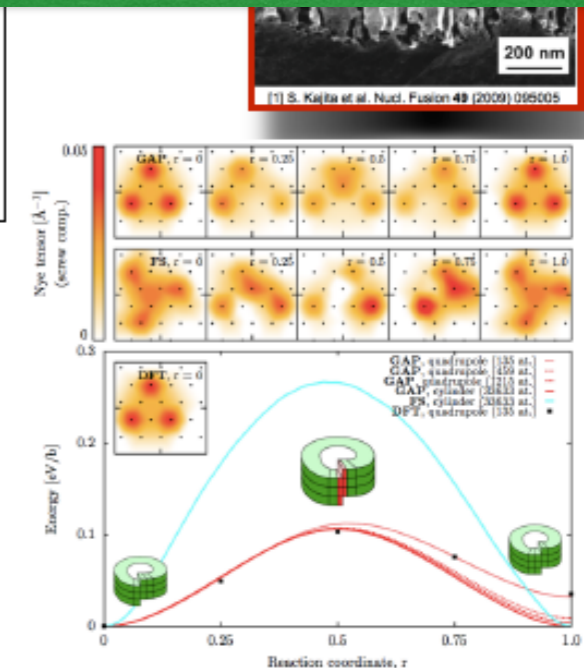
$$-i\hbar \frac{\partial \Psi}{\partial t} = \mathcal{H}\Psi$$

Challenge:

How to build efficient molecular algorithms that can learn forcefields on-the-fly with prescribed accuracy

Force Fields - $O(N)$
(interatomic potentials)
10 ms / atom / core

- Data is plentiful and “cheap” to generate
- Need hard accuracy guarantees
- Need error prediction
- Multiple scales of interactions
- Interpolation to 10^{-3} - 10^{-4} accuracy



Data-Centric Multiscale Modelling

Matthew Borg/Jason Reese (Edinburgh)

'Enhanced'
CFD

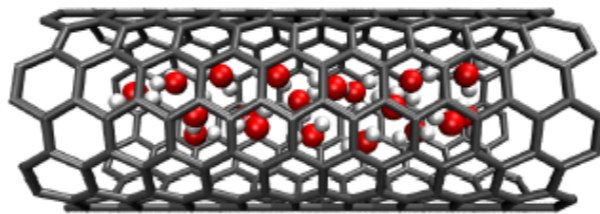
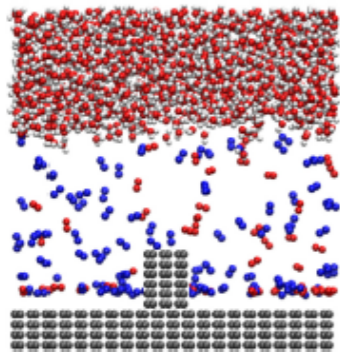
Challenge:

Automation of the modelling hierarchy across many orders of magnitude in spatial and temporal scales

CFD

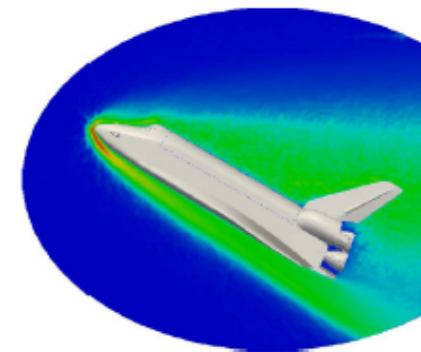
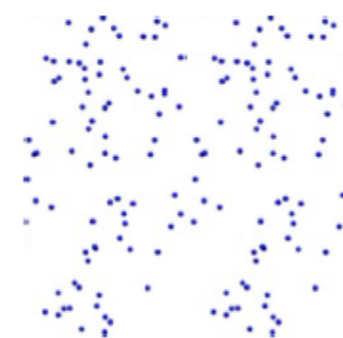
Molecular
Dynamics (MD)

(liquids, gases, solids; deterministic)



Direct Simulation
Monte Carlo (DSMC)

(rarefied gases; reactions; stochastic)



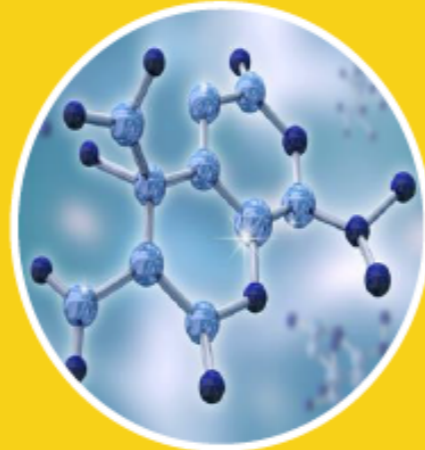
Modelling of flows, rock properties, seismic imaging

Shell Technology Centres



Geo Labs

Parallel Seismic Imaging
Micro-seismic
Induced Seismic



Sim Labs

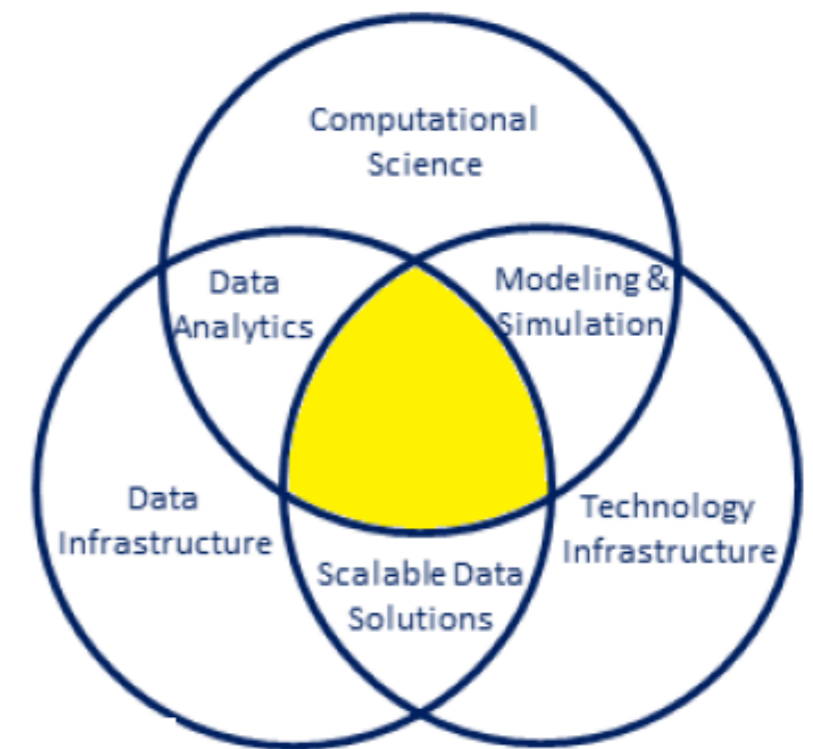
Computational Chemistry
Comp. Material Science
Flow Compute
Digital Rock



Data Labs

Adv. Analytics & Machine Learning
Sparse Modeling
Extreme statistics

High Performance Computing



Highlights

Spirited discussions about critical themes and **grand challenge problems**, e.g.

Complex flows (unsteady, non-Newtonian fluids)

Predictive molecular biology (e.g. rational drug design)

Virtual materials laboratory

Focus on enabling technologies, e.g.

Learning strategies for large scale sampling

Data-intensive scale-bridging/coupling strategies

Data-centric, hardware aware, extreme computing

Data Fusion workbench

II. Example

Molecular Dynamics

$$H(x, p) = \sum_{i=1}^N \frac{p_i^2}{2m_i} + U(x_1, x_2, \dots, x_N)$$

e.g empirical potentials

- Newton's equations + stochastic perturbations
- Computationally intensive
- Typical computations: averages in a statistical ensemble
- “Fast” force components, hence very small timesteps
- Lengths of simulations extremely limited
- Substantial share of worldwide supercomputing

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**

learning-based **QM/MM** methods [w./ G. Csanyi and others]

...Many applications in **Bayesian Inference &**

Big Data Analytics

What to do about the force error?

Langevin Dynamics

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

$$dp = -\nabla U dt - \gamma M^{-1}p dt + \sqrt{2\beta^{-1}\gamma} dW$$

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

$$\rho \propto e^{-\beta [p^T M^{-1} p / 2 + U(x)]}$$

Splitting Methods for Langevin Dynamics

$$\mathcal{L} = \mathcal{L}_A + \mathcal{L}_B + \mathcal{L}_O$$

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

$$\mathcal{L}_B = -\nabla U(x) \cdot \nabla_p$$

$$\mathcal{L}_O = -\gamma(M^{-1}p) \cdot \nabla_p + \gamma\beta^{-1}\Delta_p$$

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

Expansion of the invariant distribution

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

Leading order:

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

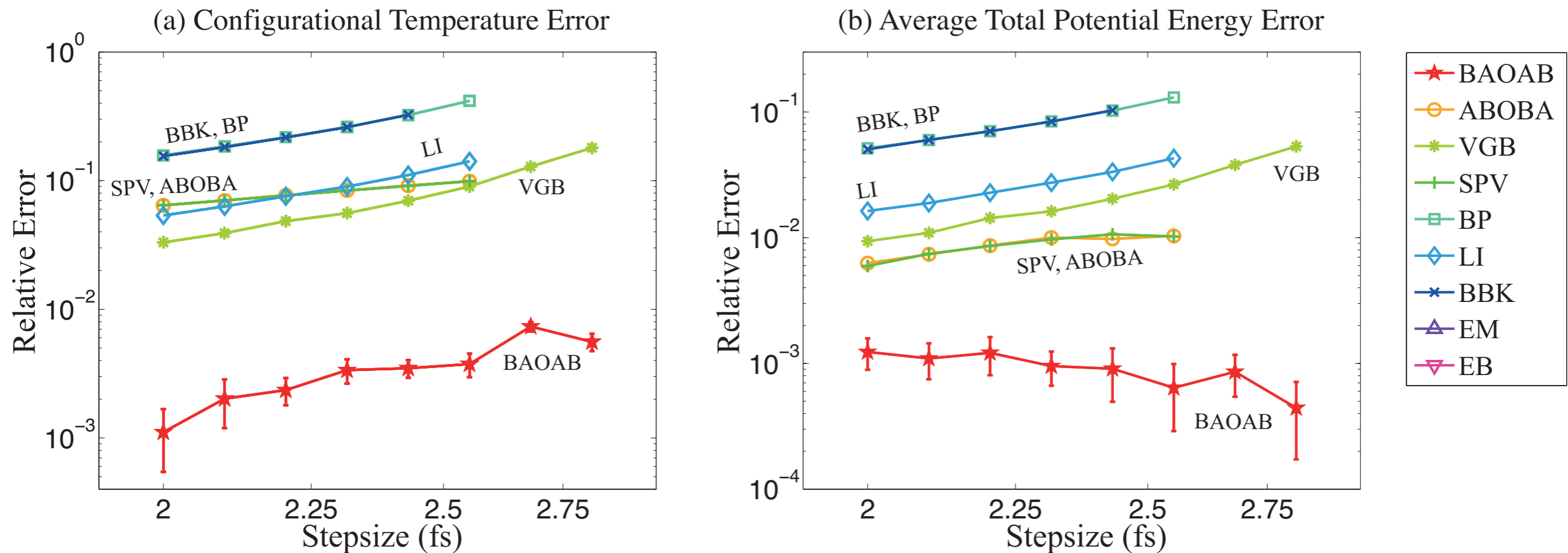
L. & Matthews, AMRX, 2013

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

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

Improvements for “real” molecular systems

for alanine dipeptide in flexible TIP3P



**BAOAB much better than alternatives for relevant
configurational quantities...**

but....

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 $\sigma(x) \approx \sigma$

$$\begin{aligned} h\tilde{F}(x) &= -h\nabla U(x) + h\eta \\ &= -h\nabla U(x) + \sqrt{h}(\sqrt{h}\eta) \end{aligned}$$

Like discretizing a stochastic differential equation with $O(h)$ variance!

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

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

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

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

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

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)

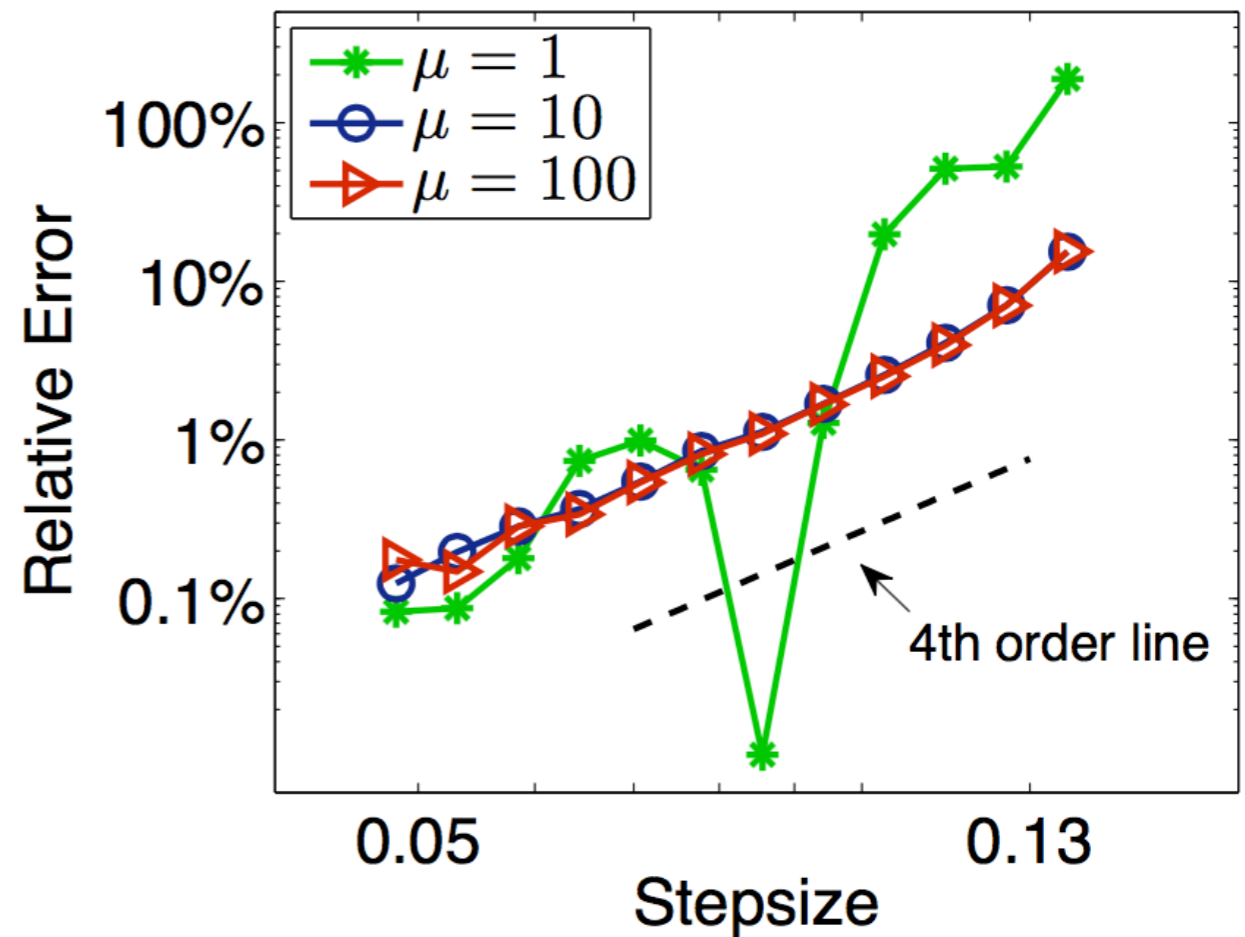
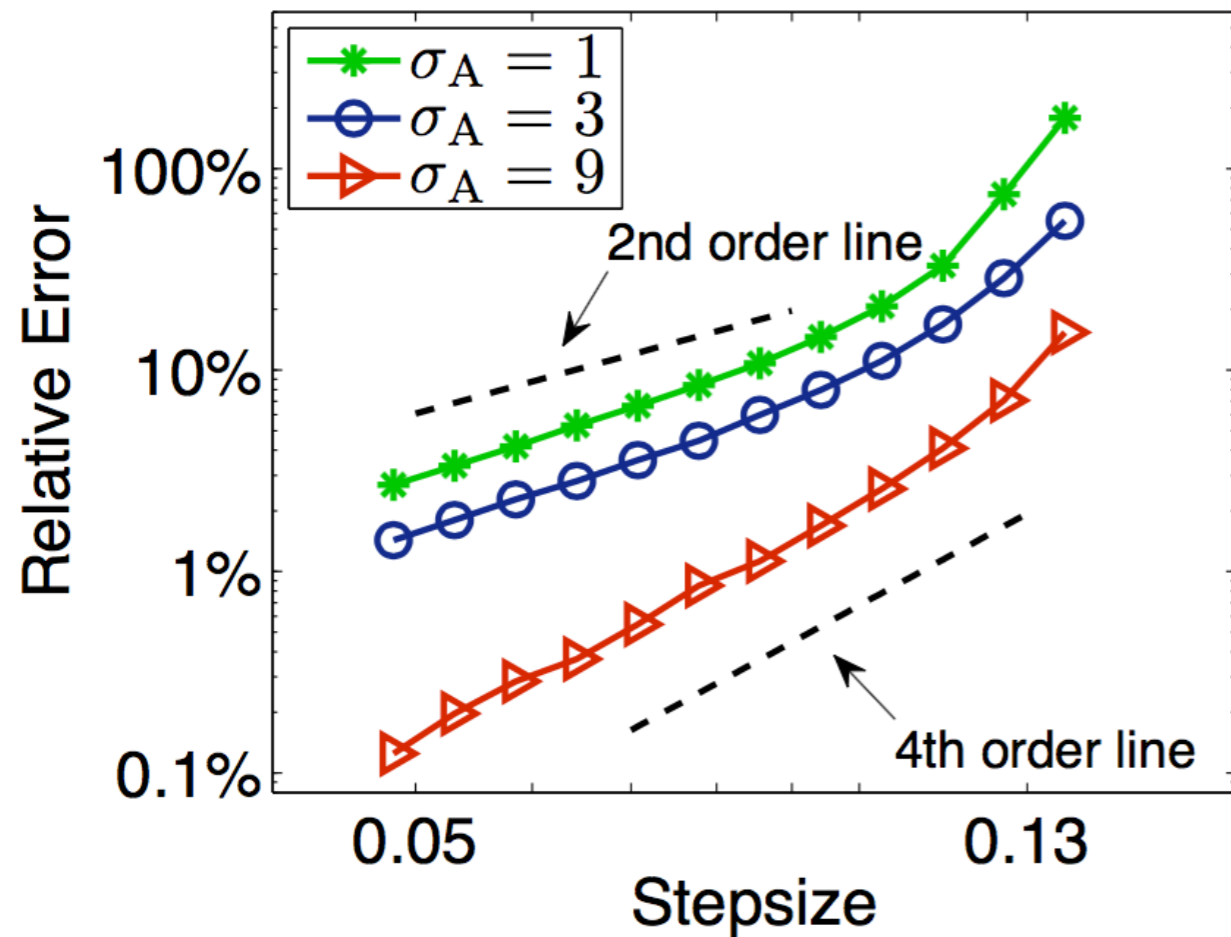
Superconvergence

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...

500 Lennard-Jones Particle MD



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

Bayesian Learning Application

Find best choice of parameters q given observations X

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

Challenges: data set very large

Ex: Netflix: 480000 users, 17000 ratings \Rightarrow 100M ratings!

Posterior probability density (from Bayes' Theorem):

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

Data Scientist Thomas Bayes, U of Edinburgh, Class of 1721

Use Maximum Likelihood Estimate/"Subsampling":

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



Bayesian Logistic Regression

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

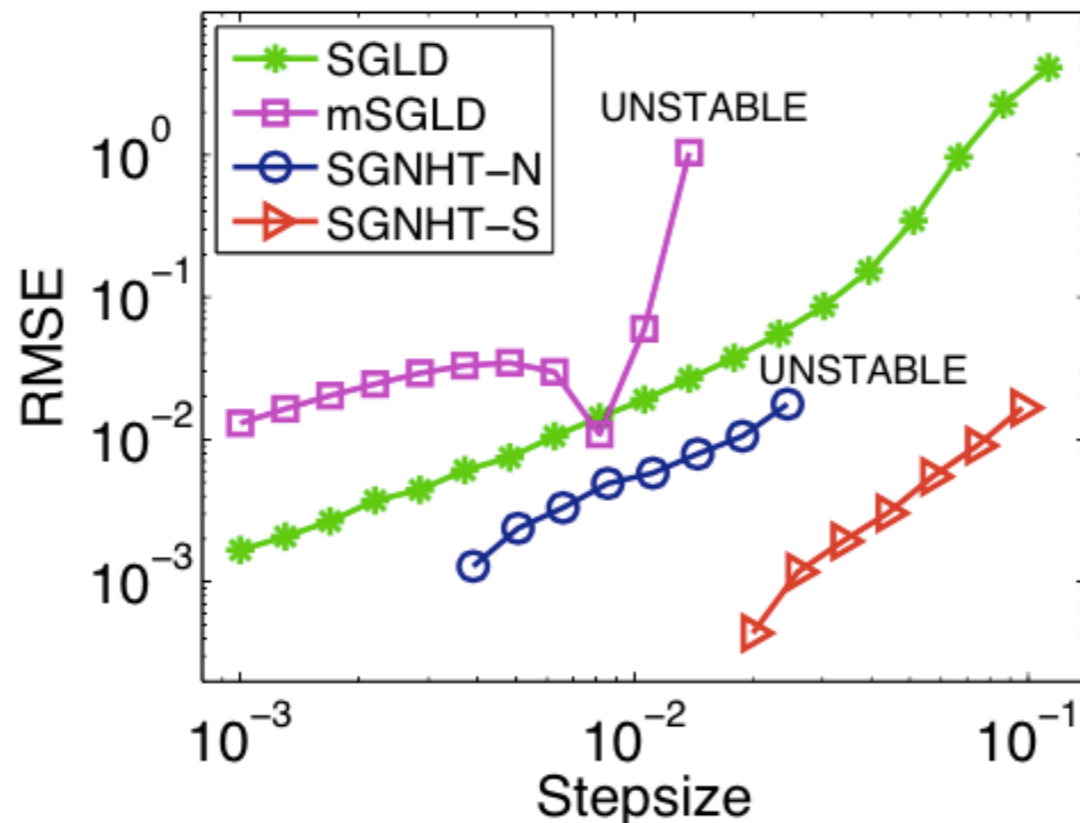
covariates e.g. age, income, ...

data e.g. voting intention

posterior parameter distribution

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

Gaussian prior

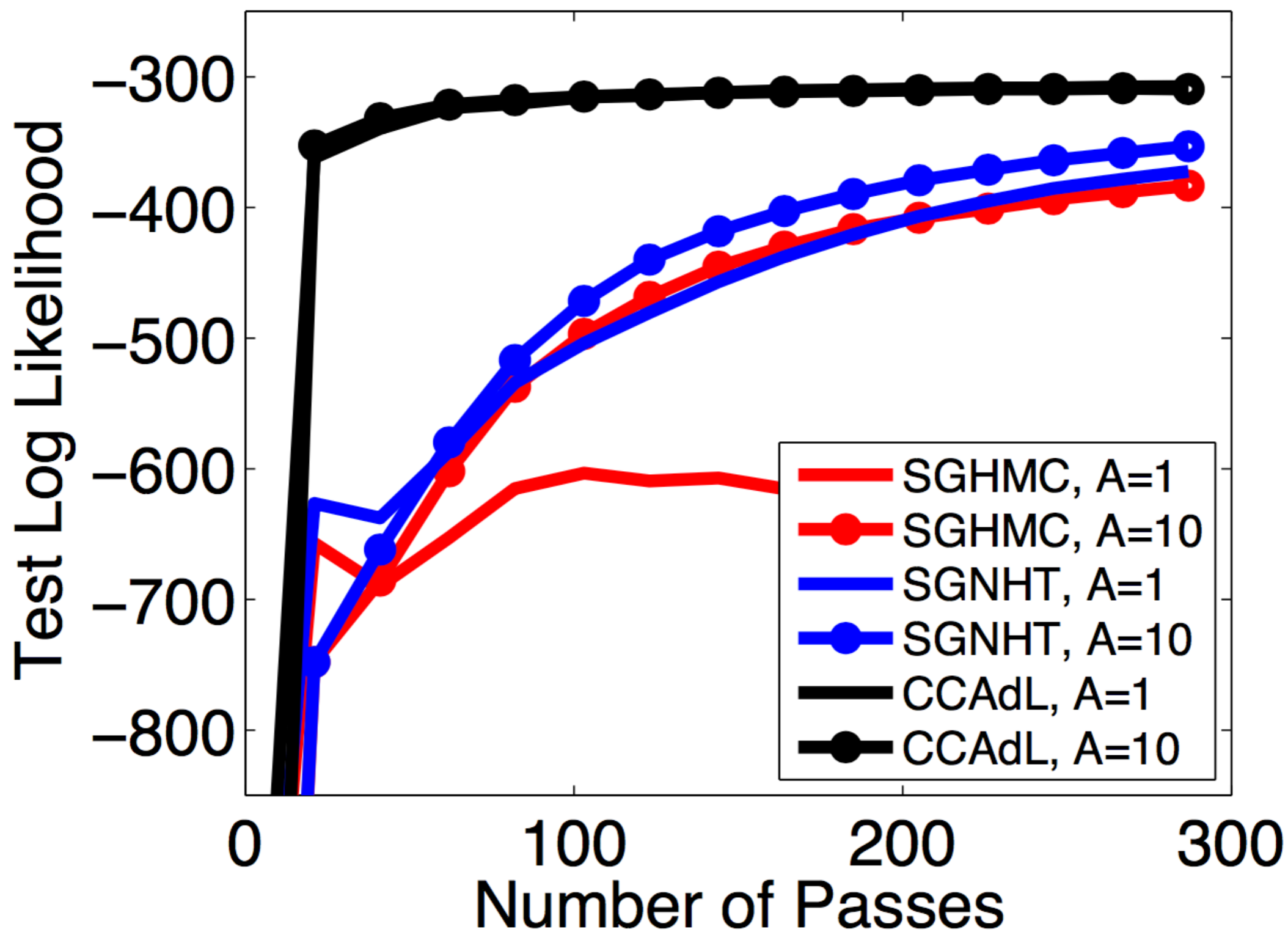


Covariance-Controlled Adaptive Langevin Dynamics

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...



Binary classification of handwritten digits 7 and 9.

Conclusions

There are many open and interesting challenges in data-centric scientific simulation.

Much of the interest lies in the interfaces between scale regimes and in incorporating data from experiment and observation

Sometimes, methods developed for solving large scale physical applications can find new uses in the world of big data analytics.