# Introduction to MD Workflows and Tools

**Prof. Vijay S. Pande, PhD**

Departments of Chemistry, Computer Science, and Structural Biology
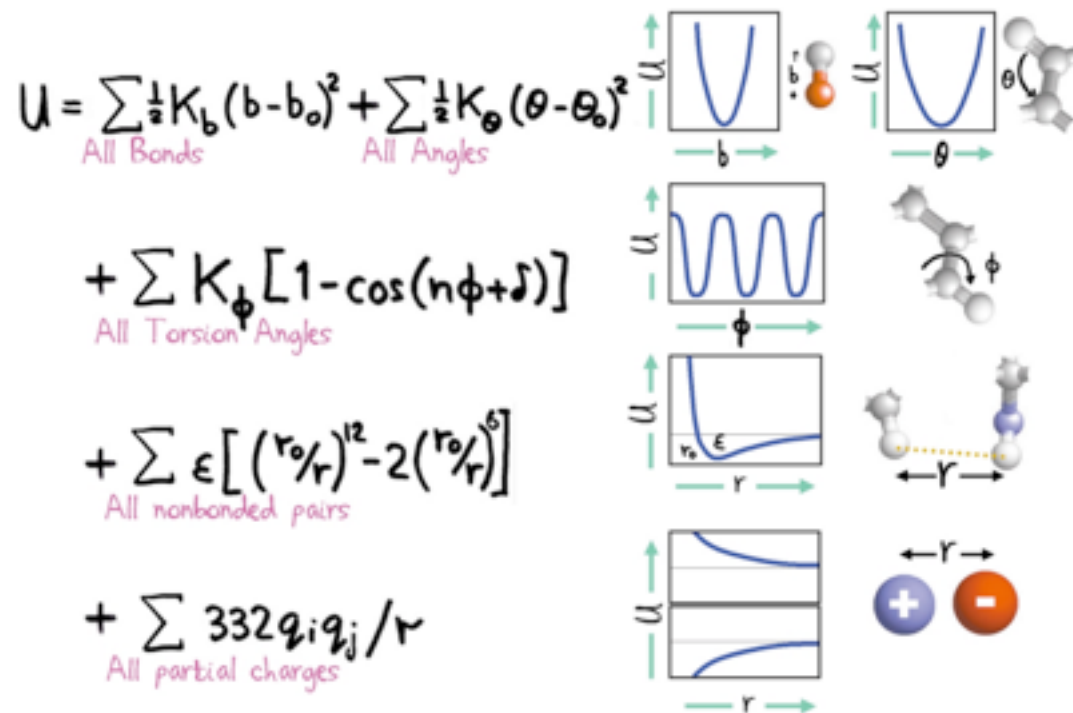Director, Biophysics Program
Director, Folding@home Distributed Computing Project
Stanford University
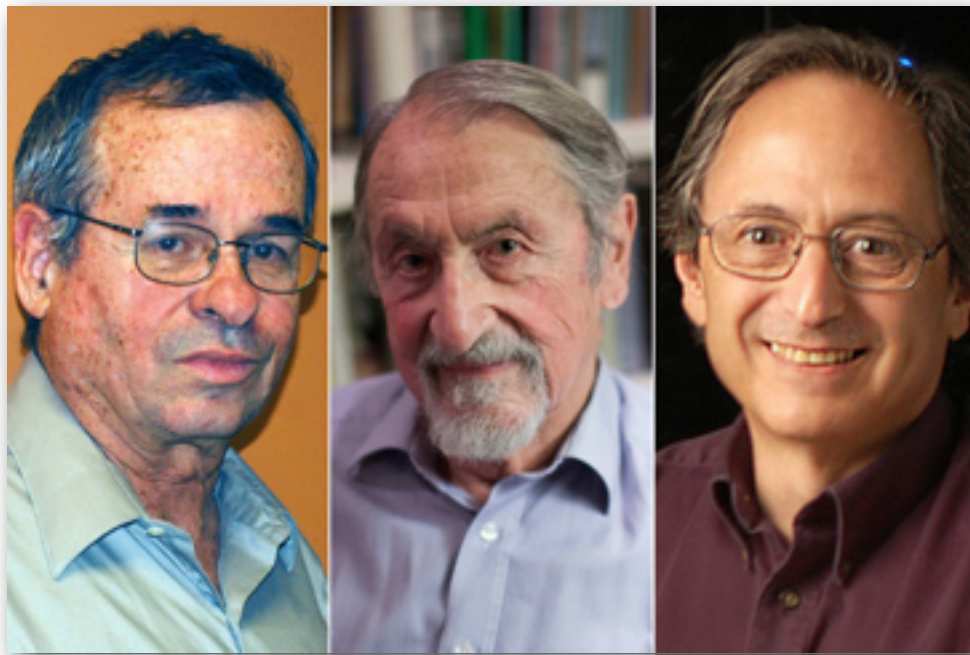
# The dream: simulating molecular dynamics

**Basic idea:** calculate forces between atoms, then numerically integrate Newton's Equations

$$U = \sum_{\text{All Bonds}} \tfrac{1}{2} K_b (b - b_0)^2 + \sum_{\text{All Angles}} \tfrac{1}{2} K_\theta (\theta - \theta_0)^2$$

$$+ \sum_{\text{All Torsion Angles}} K_\phi [1 - \cos(n\phi + \delta)]$$

$$+ \sum_{\text{All nonbonded pairs}} \epsilon \left[ \left(\tfrac{r_0}{r}\right)^{12} - 2\left(\tfrac{r_0}{r}\right)^{6} \right]$$

$$+ \sum_{\text{All partial charges}} 332 q_i q_j / r$$

M. Levitt, *Nature Structural Biology* **8** 392 (2001)
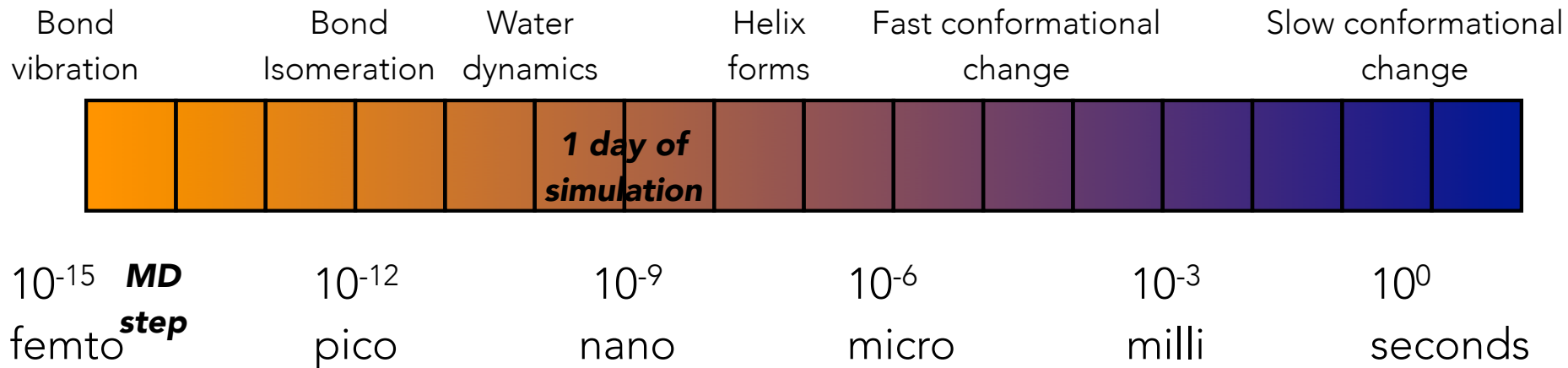
# The dream: simulating molecular dynamics

**Basic idea:** calculate forces between atoms, then numerically integrate Newton's Equations
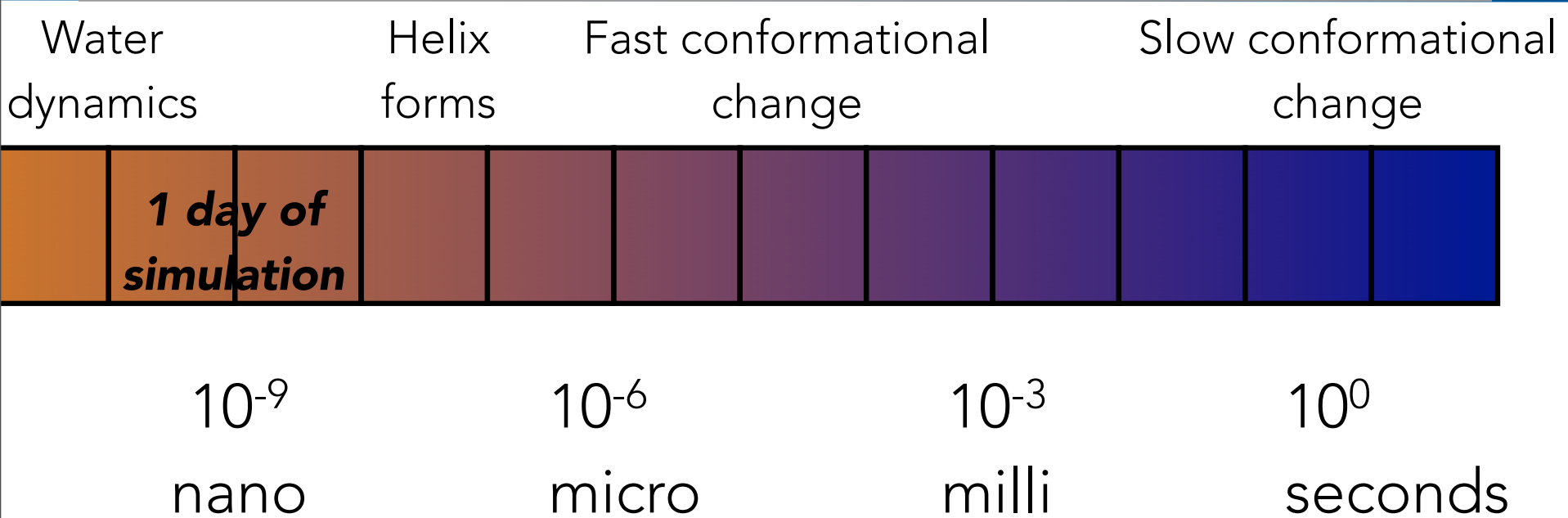


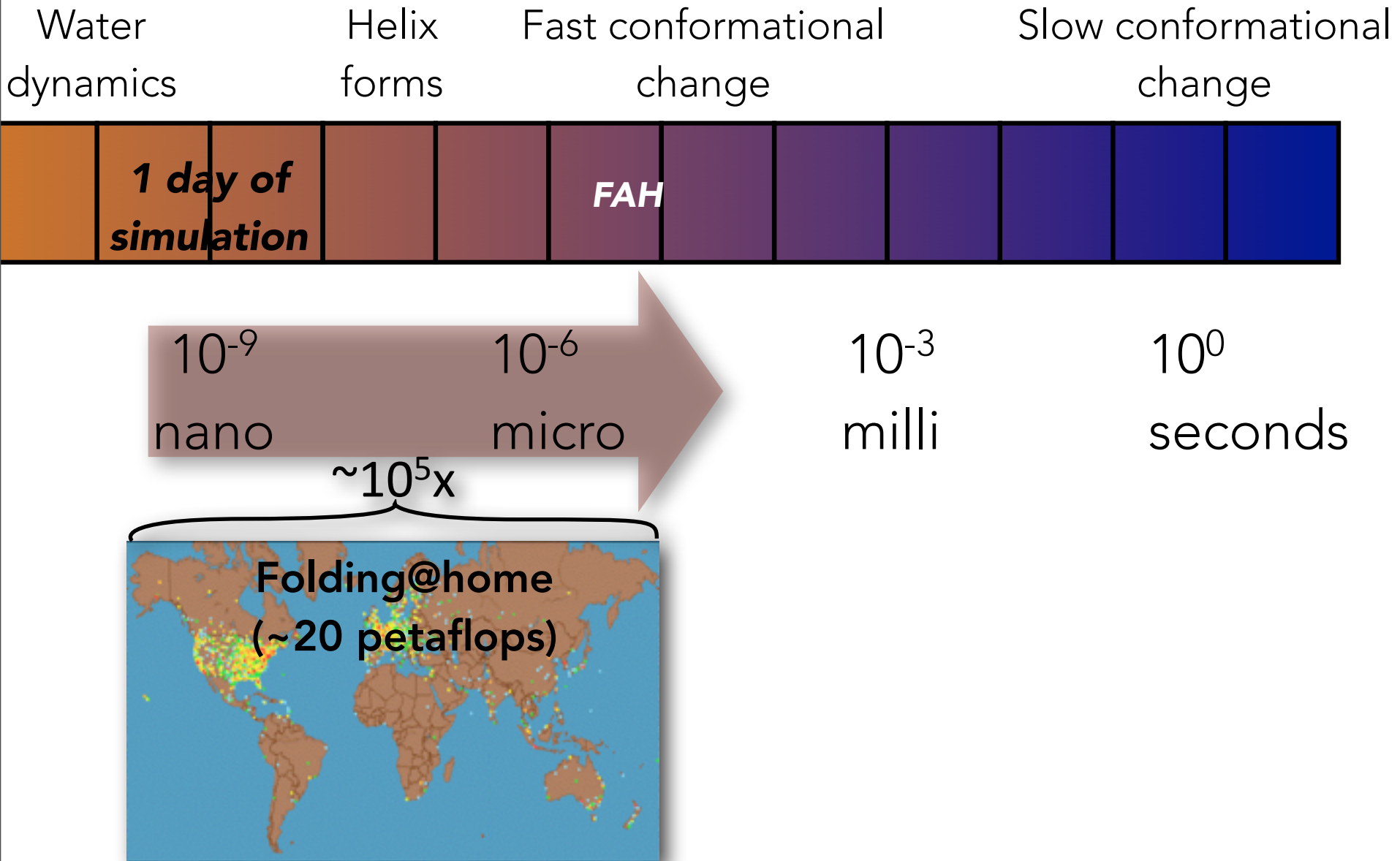*2013 Nobel Prize in Chemistry Awarded to Karplus, Levitt, and Warshel*
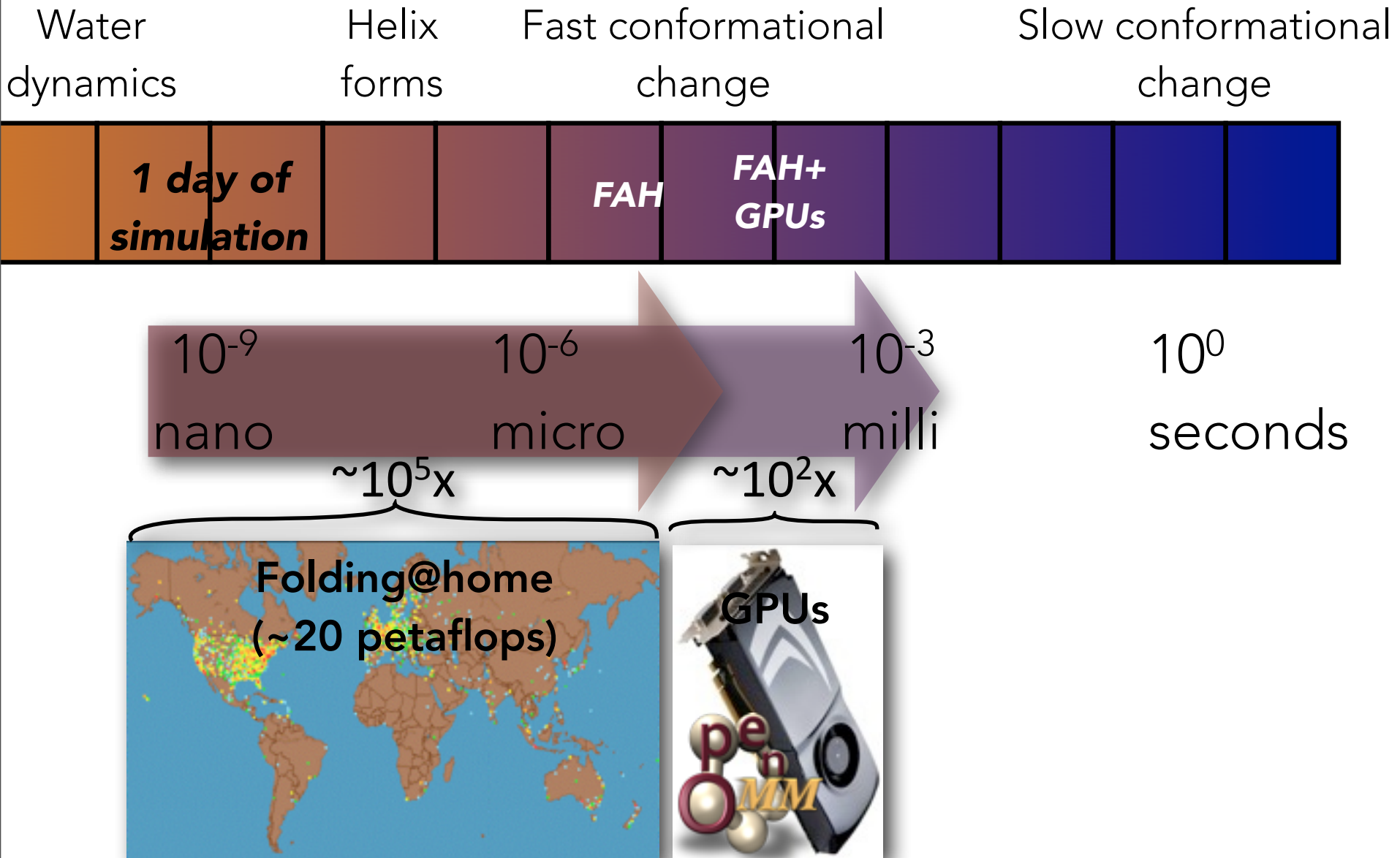
# The nightmare: long time scales

Bond vibration     Bond Isomeration     Water dynamics     Helix forms     Fast conformational change     Slow conformational change

**1 day of simulation**

$10^{-15}$ **MD step**    $10^{-12}$    $10^{-9}$    $10^{-6}$    $10^{-3}$    $10^{0}$

femto    pico    nano    micro    milli    seconds

# The nightmare: long time scales

Water dynamics      Helix forms      Fast conformational change      Slow conformational change

*1 day of simulation*

$10^{-9}$
nano

$10^{-6}$
micro

$10^{-3}$
milli

$10^{0}$
seconds

# The nightmare: long time scales

Water dynamics   Helix forms   Fast conformational change   Slow conformational change

**1 day of simulation**   *FAH*

$10^{-9}$
nano

$10^{-6}$
micro

$10^{-3}$
milli

$10^{0}$
seconds

~$10^5$x

**Folding@home
(~20 petaflops)**

# The nightmare: long time scales

Water dynamics

Helix forms

Fast conformational change

Slow conformational change

*1 day of simulation*

FAH

FAH+ GPUs

$10^{-9}$ nano

$10^{-6}$ micro

$10^{-3}$ milli

$10^{0}$ seconds

$\sim 10^{5}$x

$\sim 10^{2}$x



**Folding@home (~20 petaflops)**

**GPUs**

# The nightmare: long time scales

Water dynamics

Helix forms

Fast conformational change

Slow conformational change

**1 day of simulation**

*FAH*

*FAH+ GPUs*

*FAH+ GPUs + MSMs*

$10^{-9}$ nano

$10^{-6}$ micro

$10^{-3}$ milli

$10^{0}$ seconds

$\sim 10^5$x

$\sim 10^2$x

$\sim 10^2$x

**Folding@home (~20 petaflops)**

**GPUs**

**MSMBuilder technology**

# OpenMM suite of applications

Fast MD

ΔG calcs
(Chodera Lab)

ForceBalance
(Pande Lab)

Odin

ensemble
refinement

(Pande Lab)

MSM Accelerator: parallelize

MSM Builder: analyze

MSM Explorer: visualize

 = rapid development + rapid execution

OpenMM is an app, API, and library for rapid molecular dynamics.
Easy to modify and incorporate into any code.

# History of OpenMM

*2005*  Buck, Vishal
(Hanrahan, Darve, Pande)

*2006*  Elsen, Houston, Vishal
(Hanrahan, Darve, Pande)

**CUDA** (Buck, NVIDIA)
*2007*

*2007/8*  **FAH/ATI:** Houston, Friedrichs
(Pande, Simbios, ATI)

*Brook code*

**FAH/NVIDIA:** LeGrand, Friedrichs, Eastman (Pande, Simbios, NVIDIA)
*2008*

*2009*  **Open MM:** Friedrichs, et al
(Pande, Simbios, ATI)

*2012*  **OpenMM 4.0:** Eastman, Friedrichs et al (Simbios, Pande)

# OpenMM: JAC benchmark

| | CUDA (GTX Titan) | OpenCL (GTX Titan) | OpenCL (HD 7970) |
|---|---|---|---|
| **Implicit hbonds** | 284 | 183 | 120 |
| **Implicit hangles** | 524 | 324 | 104 |
| **RF 2fs** | 162 | 124 | 83.5 |
| **RF 5fs** | 330 | 233 | 90.2 |
| **PME 2fs** | 104 | 61 | 49.3 |
| **PME 5fs** | 226 | 132 | 63.0 |

Joint AMBER-CHARMM DHFR Benchmark in ns/day

# OpenMM roadmap

- **OpenMM 6**
  - Normal mode analysis script
  - AMOEBA OpenCL implementation
  - Constant pH implementation (JDC)
  - YANK release soon (JDC)
  - test/validate ABSINTH implicit solvent
  - More modeling tools within OpenMM app

- Further development Rosetta force field
- Triclinic boxes
- A more accurate SASA calculation for use with GB models
- Parameterize GB/VI at different temperatures
- CHARMM27 force field
- Thermodynamic ensemble validation tests
- PME for Lennard-Jones

http://wiki.simtk.org/openmm/RoadmapTimeline

# Licensing and distribution

- **API & reference BSD license, GPU kernels are LGPL**
  - free & open
  - we want LGPL to have a community owned set of GPU kernels
  - we're looking for collaborations for new features

- **But, please cite us**
  - P. Eastman, M. S. Friedrichs, J. D. Chodera, R. J. Radmer, C. M. Bruns, J. P. Ku, K. A. Beauchamp, T. J. Lane, L.-P. Wang, D. Shukla, T. Tye, M. Houston, T. Stich, C. Klein, M. R. Shirts, and V. S. Pande.  OpenMM 4.0: A Reusable, Extensible, Hardware Independent Library for High Performance Molecular Simulation.  *Journal of Computational and Theoretical Chemistry* **9** 461–469 (2013).

# How can we simulate experimentally relevant, long timescales?

*The power of
Markov State Models*

# Comparing two approaches

# Comparing two approaches



$15M **ANTON** Specialized
hardware from D.E. Shaw can
compute 14μs/day

# Comparing two approaches



$15M **ANTON** Specialized hardware from D.E. Shaw can compute 14μs/day



$0.3M GPU cluster + **OpenMM +MSMB** can also compute 14μs/day at ~1/50th the cost

# Comparing two approaches



*50x more powerful = 50x less expensive*

$15M **ANTON** Specialized hardware from D.E. Shaw can compute 14μs/day

$0.3M GPU cluster + **OpenMM +MSMB** can also compute 14μs/day at ~1/50th the cost

# Comparing two approaches

*50x more powerful = 50x less expensive*

$15M **ANTON** Specialized hardware from D.E. Shaw can compute 14µs/day

$0.3M GPU cluster + **OpenMM +MSMB** can also compute 14µs/day at ~1/50th the cost

# Comparing two approaches

*50x more powerful = 50x less expensive*

$15M **ANTON** Specialized hardware from D.E. Shaw can compute 14μs/day

$0.3M GPU cluster + **OpenMM +MSMB** can also compute 14μs/day at ~1/50th the cost

# Comparing two approaches

**OpenMM:** Over 100ns/day on 24,000 atom JAC

**MSM Builder:** http://msmbuilder.org
**OpenMM:** http://openmm.org

*50x more powerful = 50x less expensive*

$15M **ANTON** Specialized hardware from D.E. Shaw can compute 14µs/day

$0.3M GPU cluster + **OpenMM +MSMB** can also compute 14µs/day at ~1/50th the cost

# Our Goals

- **Build a model which can predict everything**
  - kinetics, thermodynamics, structure

- **Build a model which can yield powerful visualizations**
  - movies of key phenomena

- **Broad applicability**
  - works on many systems
  - easy to use, easily automated

# Comparison to other methods

# Comparison to other methods

## Popular methods

Accelerated MD      Anton      MSM      Metadynamics

Milestoning      Path-based methods      Replica Exchange

# Comparison to other methods

## Popular methods

Accelerated MD  Anton  MSM  Metadynamics

Milestoning  Path-based methods  Replica Exchange

## Our goals

# Comparison to other methods

## Popular methods

Accelerated MD          Anton          MSM          Metadynamics

Milestoning          Path-based methods          Replica Exchange

## Our goals

(1)    Fast sampling, including orthogonal degrees
         of freedom

## Popular methods

Accelerated MD        Anton        MSM        Metadynamics

Milestoning        Path-based methods        Replica Exchange

## Our goals

(1)    Fast sampling, including orthogonal degrees of freedom

## Popular methods

Accelerated MD      Anton      MSM      ~~Metadynamics~~

Milestoning      Path-based methods      Replica Exchange

## Our goals

(1) Fast sampling, including orthogonal degrees of freedom

(2) Can discover end points

# Comparison to other methods

## Popular methods

Accelerated MD          Anton          MSM          ~~Metadynamics~~

~~Milestoning~~          ~~Path-based methods~~          Replica Exchange

## Our goals

(1)    Fast sampling, including orthogonal degrees of freedom

(2)    Can discover end points

# Comparison to other methods

## Popular methods

Accelerated MD          Anton          MSM          ~~Metadynamics~~

~~Milestoning~~          ~~Path-based methods~~          Replica Exchange

## Our goals

(1)  Fast sampling, including orthogonal degrees of freedom

(2)  Can discover end points

(3)  Can predict kinetics (& thermodynamics, & structure)

# Comparison to other methods

## Popular methods

~~Accelerated MD~~     Anton          MSM          ~~Metadynamics~~

~~Milestoning~~     ~~Path-based methods~~     Replica Exchange

## Our goals

(1)   Fast sampling, including orthogonal degrees of freedom

(2)   Can discover end points

(3)   Can predict kinetics (& thermodynamics, & structure)

# Comparison to other methods

## Popular methods

~~Accelerated MD~~  Anton  MSM  ~~Metadynamics~~

~~Milestoning~~  ~~Path-based methods~~  ~~Replica Exchange~~

## Our goals

(1) Fast sampling, including orthogonal degrees of freedom

(2) Can discover end points

(3) Can predict kinetics (& thermodynamics, & structure)

# What are Markov State Models (MSMs)?

MSMs **automatically** build a **Master Equation** with MD simulation, typically with **many short (~μs) trajectories**

$$\frac{dp_i}{dt} = \sum_l \left[ k_{l,i} p_l - k_{i,l} p_i \right]$$

MSMs **automatically** build a **Master Equation** with MD simulation, typically with **many short (~μs) trajectories**

$$\frac{dp_i}{dt} = \sum_l \left[ k_{l,i} p_l - k_{i,l} p_i \right]$$

probability of being in state $i$

# What are Markov State Models (MSMs)?

MSMs **automatically** build a **Master Equation** with MD simulation, typically with **many short (~µs) trajectories**

$$\frac{dp_i}{dt} = \sum_l \left[ k_{l,i} p_l - k_{i,l} p_i \right]$$

probability of being in state
$i$

rate of change between states

# What are Markov State Models (MSMs)?

MSMs **automatically** build a **Master Equation** with MD simulation, typically with **many short (~μs) trajectories**

$$\frac{dp_i}{dt} = \sum_l [k_{l,i}p_l - k_{i,l}p_i]$$

probability of being in state $i$

rate of change between states

### *with the goals of:*

(1) aiding simulators **reach long timescales** and

(2) **gaining novel insight** from their simulations

MSMs **automatically** build a **Master Equation** with MD simulation, typically with **many short (~µs) trajectories**

$$\frac{dp_i}{dt} = \sum_l \left[ k_{l,i} p_l - k_{i,l} p_i \right]$$

probability of being in state $i$

rate of change between states

***with the goals of:***

(1) aiding simulators **reach long timescales** and

(2) **gaining novel insight** from their simulations

**see the work of:** Andersen, Best, Bowman, Caflisch, Chodera, Deuflhard, Dill, Grubmüller, Huang, Hummer, Levy, Noé, Pande, Pitera, Roux, Schütte, Swope, Weber

# Short trajectories vs long timescales?

## Two state (Single Barrier) Case

$$A \xrightarrow{k} B$$



Probability of crossing = 1 - exp(-kt) ≈ kt

# Short trajectories vs long timescales?

## Two state (Single Barrier) Case

$$A \xrightarrow{k} B$$



Probability of crossing = 1 - exp(-kt) ≈ kt

$p \approx kt$ for short time

(y-axis: probability, x-axis: time)

# Short trajectories vs long timescales?

## Two state (Single Barrier) Case

$$A \xrightarrow{k} B$$



**Probability of crossing = 1 - exp(-kt) ≈ kt**

$p \approx kt$ for short time

for $k = 1/\mu s$, $t = 0.01\mu s$, $p = 1\%$
i.e. 1 out of 100 will cross!

# Key stages in MSM construction

# Key stages in MSM construction



MSMAccelerator
Round 2. Beta = 0.00

run multiple trajectories

# Key stages in MSM construction



MSMAccelerator
Round 2. Beta = 0.00

run multiple trajectories

MSMAccelerator
Starting states for round 3

build MSM,
choose new starting points

# Key stages in MSM construction



MSMAccelerator
Round 2. Beta = 0.00

run multiple trajectories

MSMAccelerator
Starting states for round 3

build MSM,
choose new starting points

MSMAccelerator
Round 10. Beta = 0.05

run new trajectories

# Key stages in MSM construction



run multiple trajectories

build MSM,
choose new starting points

run new trajectories

repeat until convergence

# Key stages in MSM construction



MSMAccelerator
Round 2. Beta = 0.00

run multiple trajectories

MSMAccelerator
Starting states for round 3

build MSM,
choose new starting points

Adaptive sampling pushes in all degrees of freedom, not just pre-chosen coordinates. This is very important in high dim spaces.

run new trajectories

repeat until convergence

## MSM Adaptive Sampling

## Single long trajectory

## MSM Adaptive Sampling

## Single long trajectory

## MSM Adaptive Sampling

## Single long trajectory



MSMAccelerator
Round 1. Beta = 0.00



One trajectory
Round 1



**efficient ⊥ sampling, trivial to parallelize**

# MSM vs long trajectory

(McGibbon, Kiss, Harrigan, Lane, VSP)
(movie by Harrigan, McGibbon)

## MSM Adaptive Sampling

## Single long trajectory



**efficient ⊥ sampling, trivial to parallelize**

**cross barriers much slower, much worse statistics**

# Comparison to other methods

- **aMD**
  - removes kinetic information
  - speeds on certain degrees of freedom — must know which ones are slow



MSMAccelerator
Round 1. Beta = 0.00

- **Metadynamics**
  - removes kinetic information
  - drives on pre-chosen degrees of freedom, misses key challenge of how to sample orthogonal dofs

- **Replica Exchange**
  - removes kinetic information
  - works best for energy barriers, not ΔG barriers

# Comparison to other methods

- **Highly parallel MD**
  - still requires the kinetic analysis.
  - many short trajectories are MUCH more efficient
  - very expensive (50x) given throughput: GPU cluster better at many short trajectories

# Making sense of MSMs: lumping

Macrostate chain $(\mathbf{y}_n)$       Microstate chain $(\mathbf{z}_n)$



$y_1 \longrightarrow z_1$

$\Delta\tau_{\mathrm{lag}} \downarrow$

$y_2 \longrightarrow z_2$

$\Delta\tau_{\mathrm{lag}} \downarrow$

$y_3 \longrightarrow z_3$

# Formalization of lumping

- The model in this case is the lumping $M : Z \rightarrow Y$, a mapping from microstates to macrostates.
- The model is parametrized by the transition probability matrix $T$, and the local equilibrium distributions for the microstates $\Theta$.
- We can factorize the evidence into two factors:

$$
\begin{aligned}
P(\mathbf{z}_n|M) &= \int dT d\Theta\, P(\mathbf{z}_n|T, \Theta, M) P(T, \Theta|M) \\
&= \int dT d\Theta\, P(\mathbf{y}_n|T, M) P(T|M) P(\mathbf{z}_n|\mathbf{y}_n, \Theta, M) P(\Theta|M) \\
&= \underbrace{\int dT\, P(\mathbf{y}_n|T, M) P(T|M)}_{\text{Macrostate Markov chain}} \times \\
&\quad \underbrace{\int d\Theta\, P(\mathbf{z}_n|\mathbf{y}_n, \Theta, M) P(\Theta|M)}_{\text{Microstates from equilibrium within macrostates}}
\end{aligned}
$$

# What can MSMs do?

# MSMs reach long timescales



0.000 us

Copernicus: A new paradigm for parallel adaptive molecular dynamics. *Supercomputing 2011* (2011)

# MD simulation has come a long way

# MD simulation has come a long way

# MSMs make quantitative predictions

# MSMs for protein-ligand binding

(Lawrenz, VSP)

*surface view*

*flexible loop ensnares ligand*

*ribbon + surface view*

Movie made with VMD

Monday, March 24, 14

# The cloud looks a lot like Folding@home

**Large-scale, distributed, heterogeneous, loosely coupled, no common filesystem**

# The cloud looks a lot like Folding@home

**Large-scale, distributed, heterogeneous, loosely coupled, no common filesystem**

# Recent results using Google cloud



Monday, March 24, 14

# Recent results using Google cloud



**nature chemistry**

JANUARY 2014 VOL 6 NO 1
www.nature.com/naturechemistry

GPCRS in the Cloud

---

## Cloud-based simulations on Google Exacycle reveal ligand modulation of GPCR activation pathways

Kai J. Kohlhoff[1,4]*, Diwakar Shukla[1,2], Morgan Lawrenz[2], Gregory R. Bowman[3], David E. Konerding[4], Dan Belov[4], Russ B. Altman[1,3]* and Vijay S. Pande[2]*

Simulations can provide tremendous insight into atomistic details of biological mechanisms, but micro- to millisecond timescales are historically only accessible on dedicated supercomputers. We demonstrate that cloud computing is a viable alternative and brings long-timescale processes within reach of a broader community. We used Google's Exacycle cloud-...

"The unprecedented millisecond simulation timescales presented here for GPCR activation require computing architectures capable of such extensive sampling. Cloud computing provides a promising new avenue to tackle these types of questions … **Our work on Google's Exacycle platform demonstrates that large-scale exploratory analysis in the cloud can deliver new insight into biological problems.** "

[1]Department of Bioengineering, Stanford University, 450 Serra Mall, Stanford, California 94305, USA, [2]Department of Chemistry, Stanford University, 450 Serra Mall, Stanford, California 94305, USA, [3]Department of Genetics, Stanford University, 450 Serra Mall, Stanford, California 94305, USA, [4]Google Inc., 1600 Amphitheatre Parkway, Mountain View, CA 94043, USA. [†]These authors contributed equally. *e-mail: kohlhoff@google.com; russ.altman@stanford.edu; pande@stanford.edu