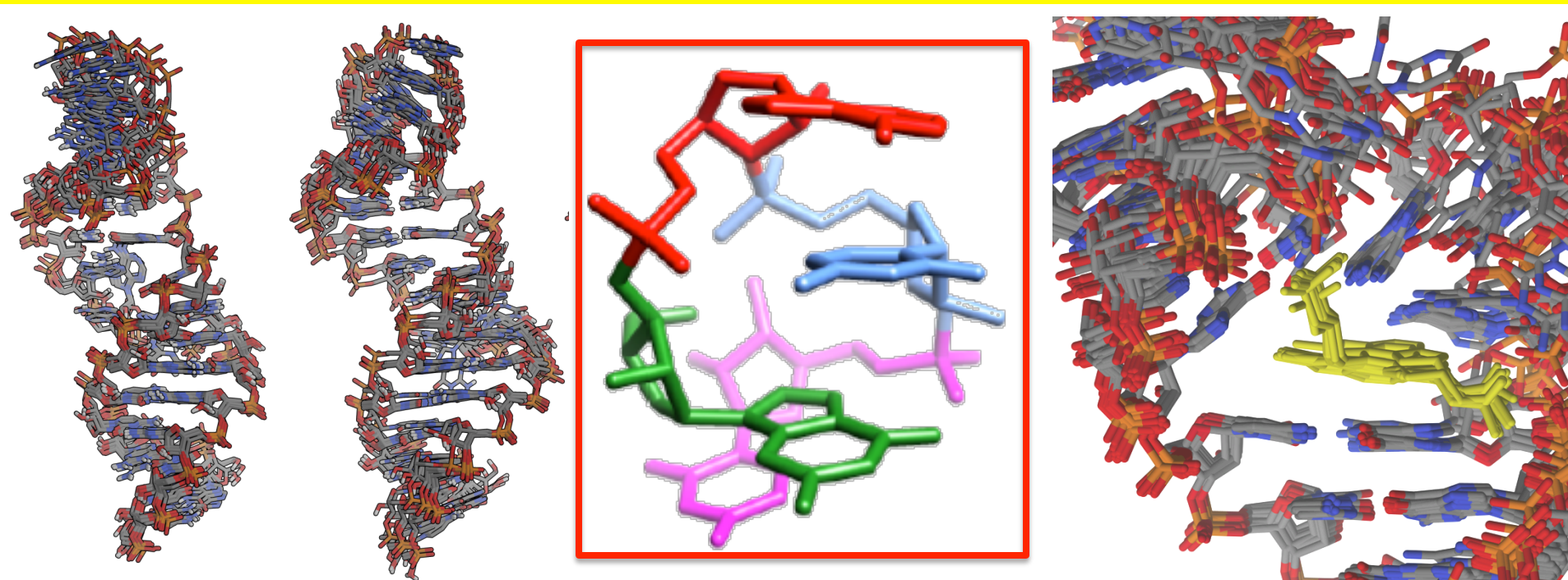


Using large-scale simulation to reproducibly converge nucleic acid structure and dynamics



Thomas E. Cheatham III
Associate Professor
Department of Medicinal Chemistry
College of Pharmacy, University of Utah

People:

Niel Henriksen, Hamed Hayatshahi, Dan Roe,
Julien Thibault, Kiu Shahrokh, Rodrigo Galindo,
Christina Bergonzo, Sean Cornillie

\$\$\$:



National Science Foundation
WHERE DISCOVERIES BEGIN

- R01-GM098102: “RNA-ligand interactions: sim. & experiment” ~2015
- R01-GM072049: “P450 dehydrogenation mechanisms” ~2014
- R01-GM081411: “...simulation ... refinement of nucleic acid” ~2013
- NSF CHE-1266307 “CDS&E: Tools to facilitate deeper data analysis, ...” ~2015
- NSF “Blue Waters” PetaScale Resource Allocation for AMBER RNA

Computer time:



D E Shaw Research

“Anton”
(3 past awards)



XRAC MCA01S027
~15M core hours



~14M GPU hours

!!!

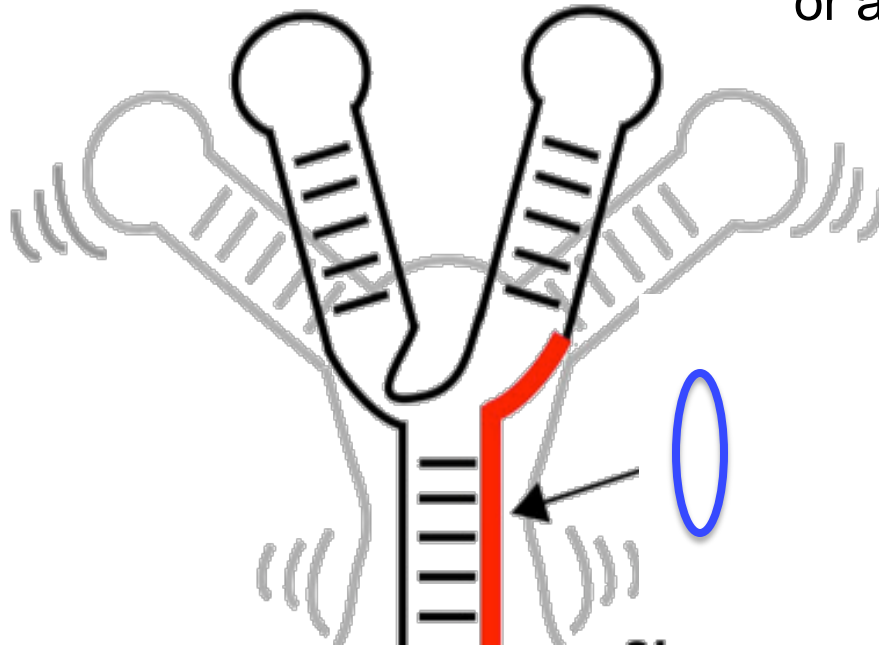


~3M hours

Accurate modeling of RNA and other biomolecules requires:
accurate and fast simulation methods
validated RNA, protein, water, ion, and ligand “force fields”
“good” experiments to assess results
dynamics and complete sampling: (convergence, reproducibility)

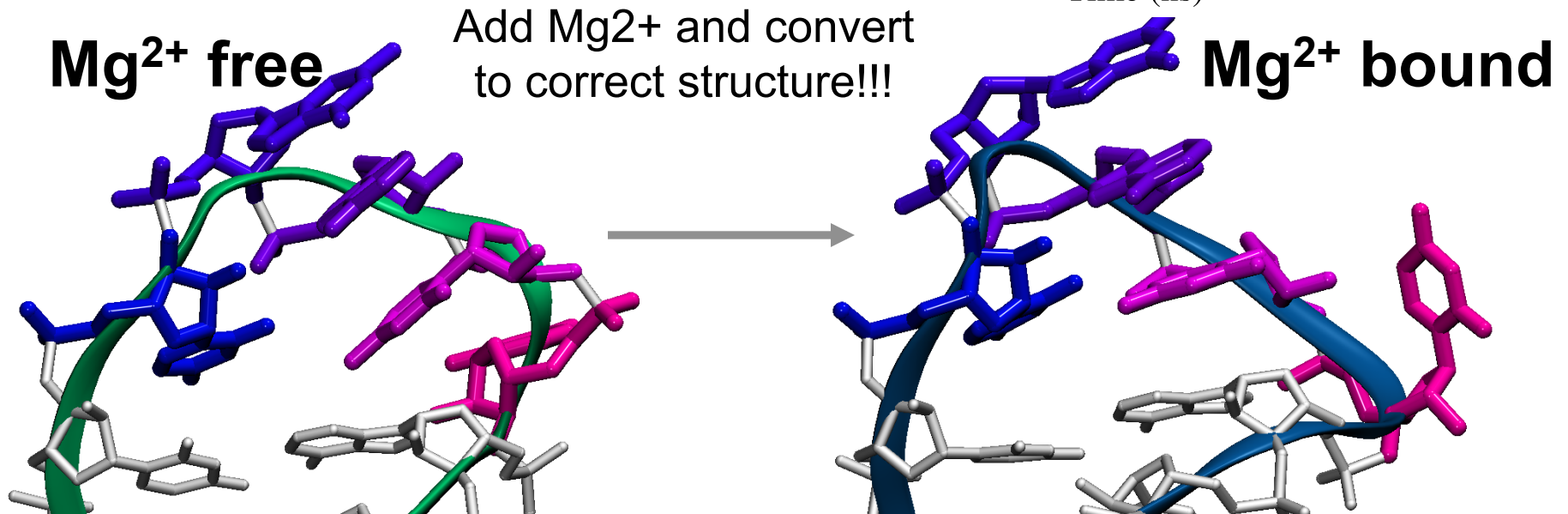
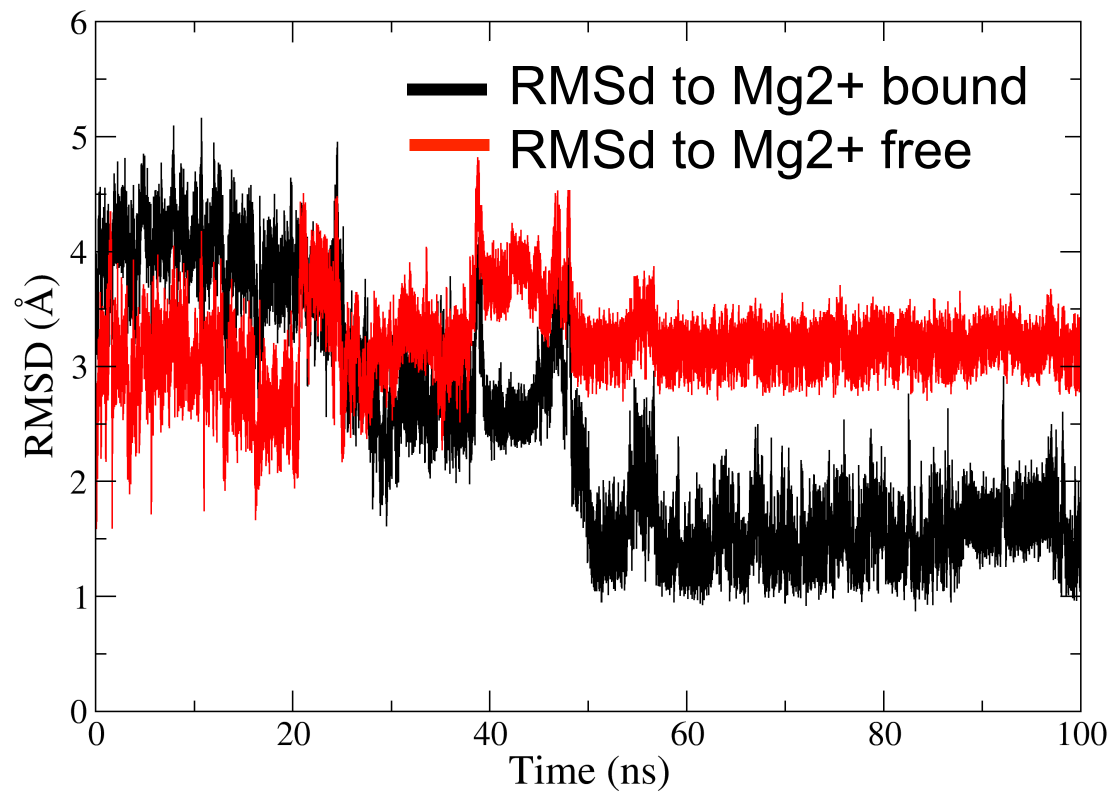
Question: Is the movement real or artifact?

True RNA dynamics
or artifact of force field?



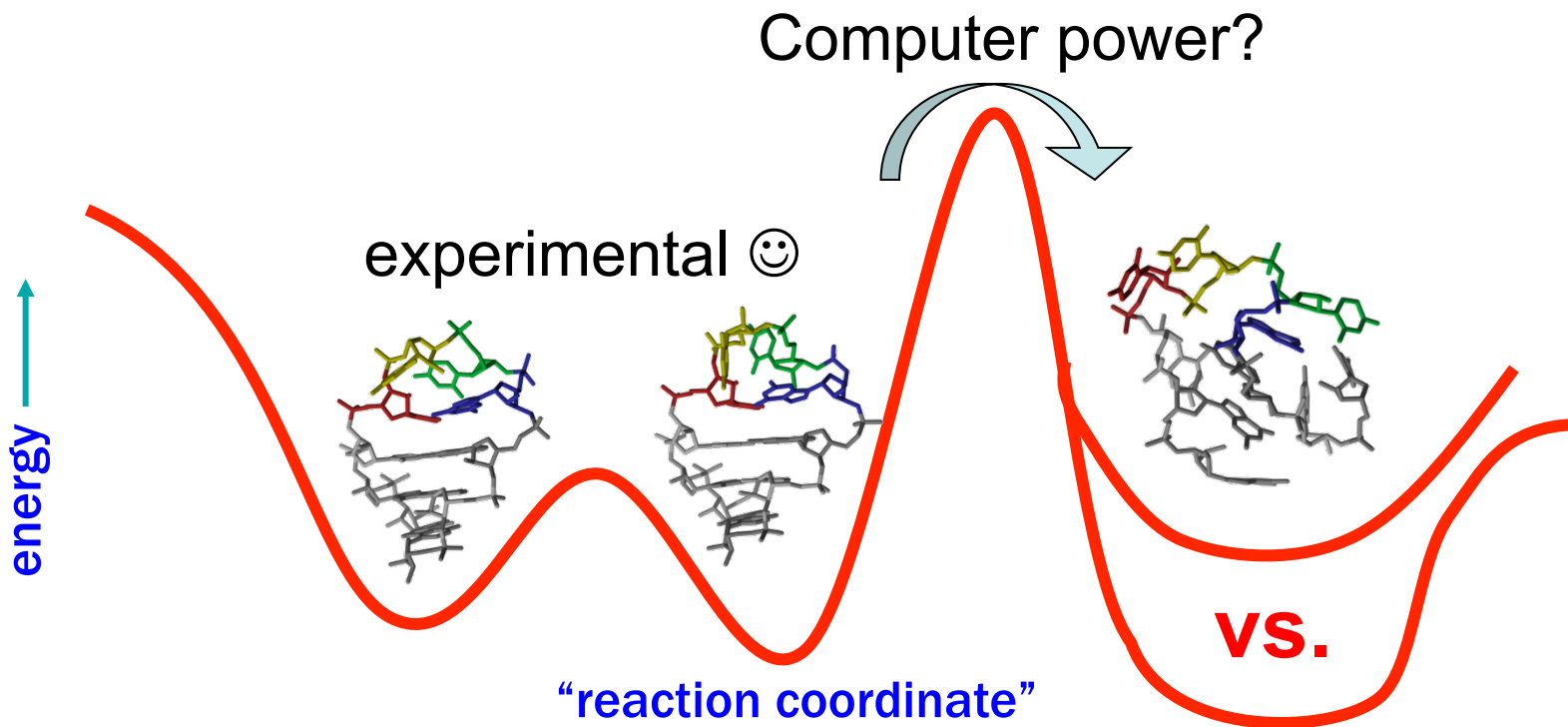
We're seeing
some progress!!!

(vsrSL5)



are the force fields reliable? (free energetics, sampling, dynamics)

Short simulations stay near experimental structure; longer simulations invariably move away and often to unrealistic lower energy structures...



How to fully sample conformational ensemble?



↑
16 μ s/day!

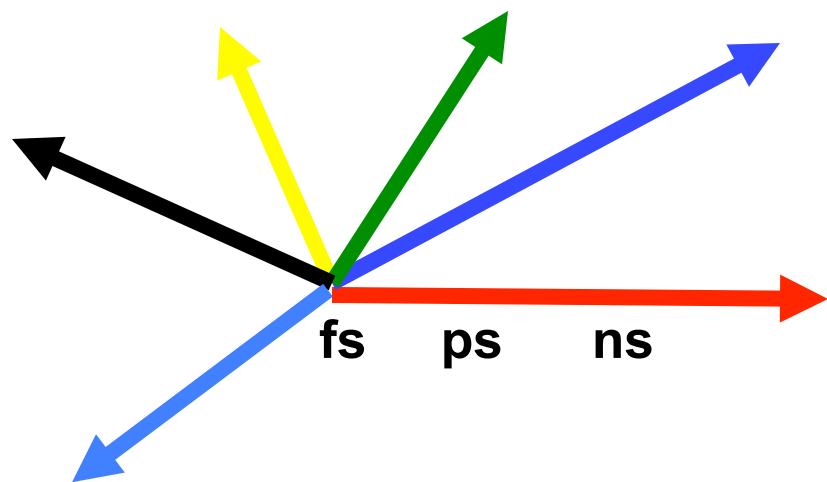
brute force – long contiguous in time MD
requires: special purpose / unique hardware

D.E. Shaw's Anton machine



Simulating protein movements using Anton could aid drug design.

SCIENCE/AAAS



ensembles of
independent
simulations

AMBER on GPUs



110 ns/day!

~1978 - present

amber

Assisted Model Building with Energy Refinement

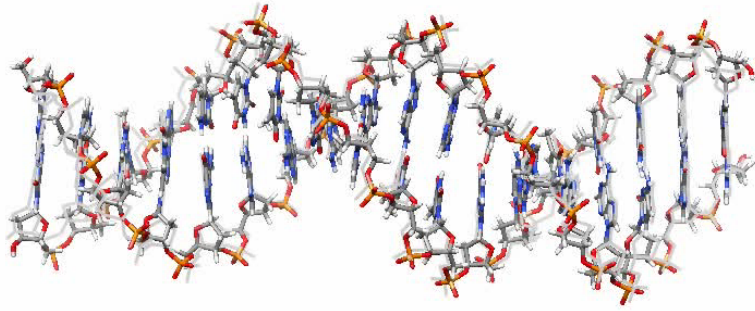
code vs. force field

Amber 14 released April, 2014

- 1.23x increase in GPU performance; peer-to-peer [fully deterministic, mixed SP/fixed precision, ||-ized]
- **support for M-REMD simulation and analysis**
- constant pH
- new TI methods
- more methods ported to GPU

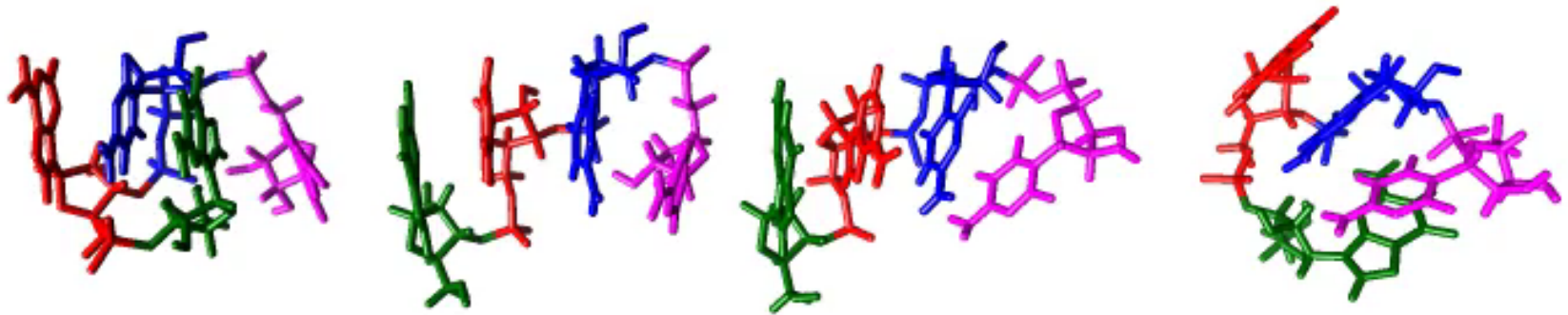
Today: two “long-time-to-develop” short stories...

- ✓ can we converge DNA duplex structure/dynamics?



Anonymous NIH R-01 reviewer in 2005:
“One has to wonder how many relatively short MD simulations have to be performed on short DNA fragments before what can be learned will have been learned...”

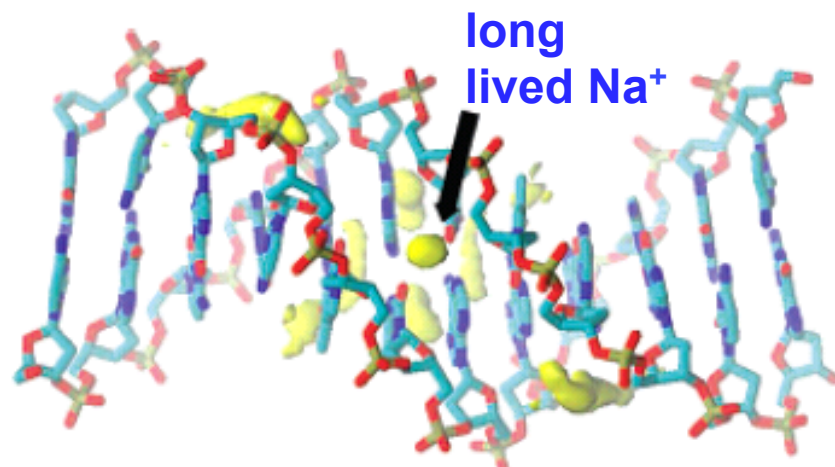
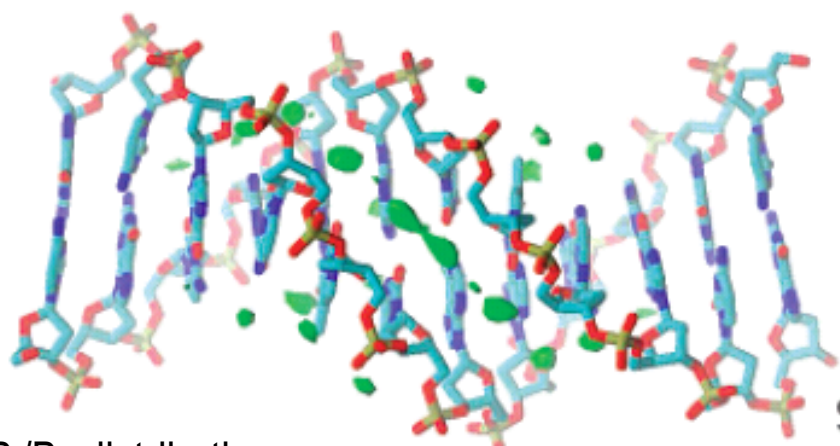
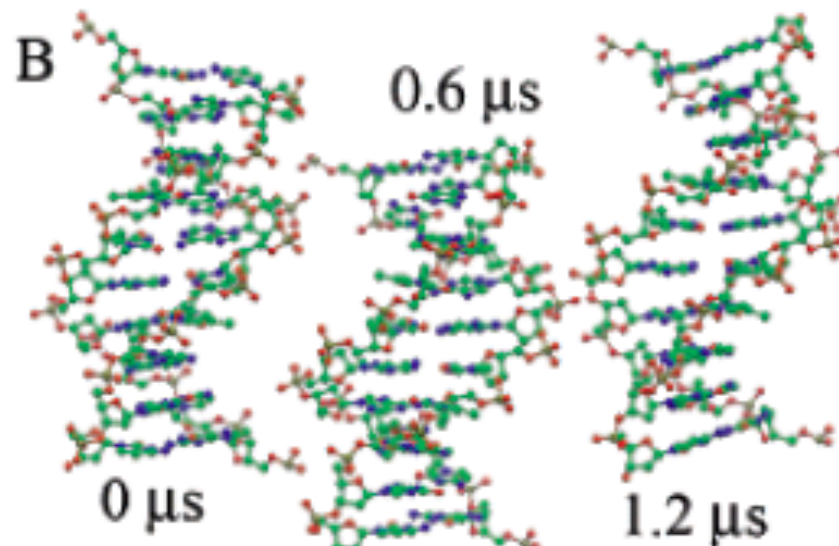
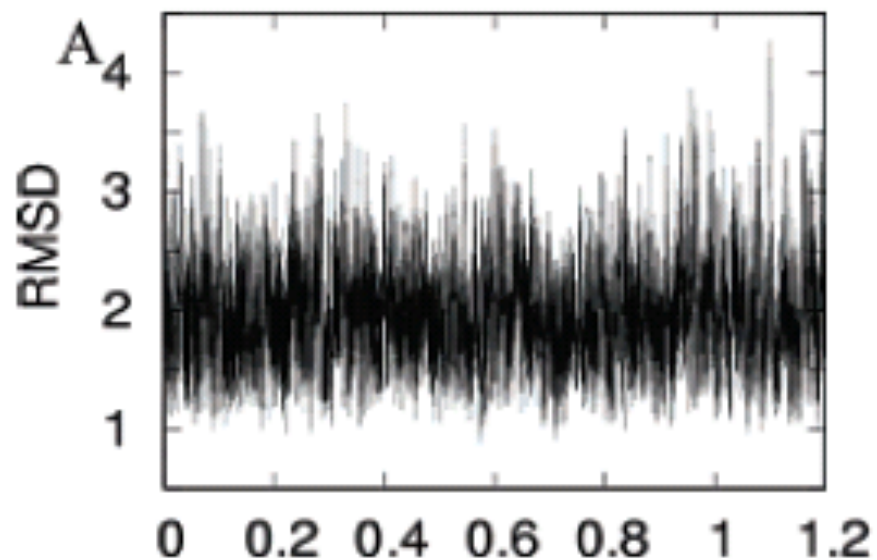
- ✓ sampling RNA structure *accurately* is difficult



Dynamics of B-DNA on the Microsecond Time Scale

J. AM. CHEM. SOC. 2007,
129, 14739–14745

Alberto Pérez,^{†,‡} F. Javier Luque,[§] and Modesto Orozco^{*,†,‡,||}



Convergence? Not yet...

Anton “testing” for ABC

ABC benchmark (50 ns, SPC/E + KCI)

GAAC: GCACGAA**CG**AACGAACGC

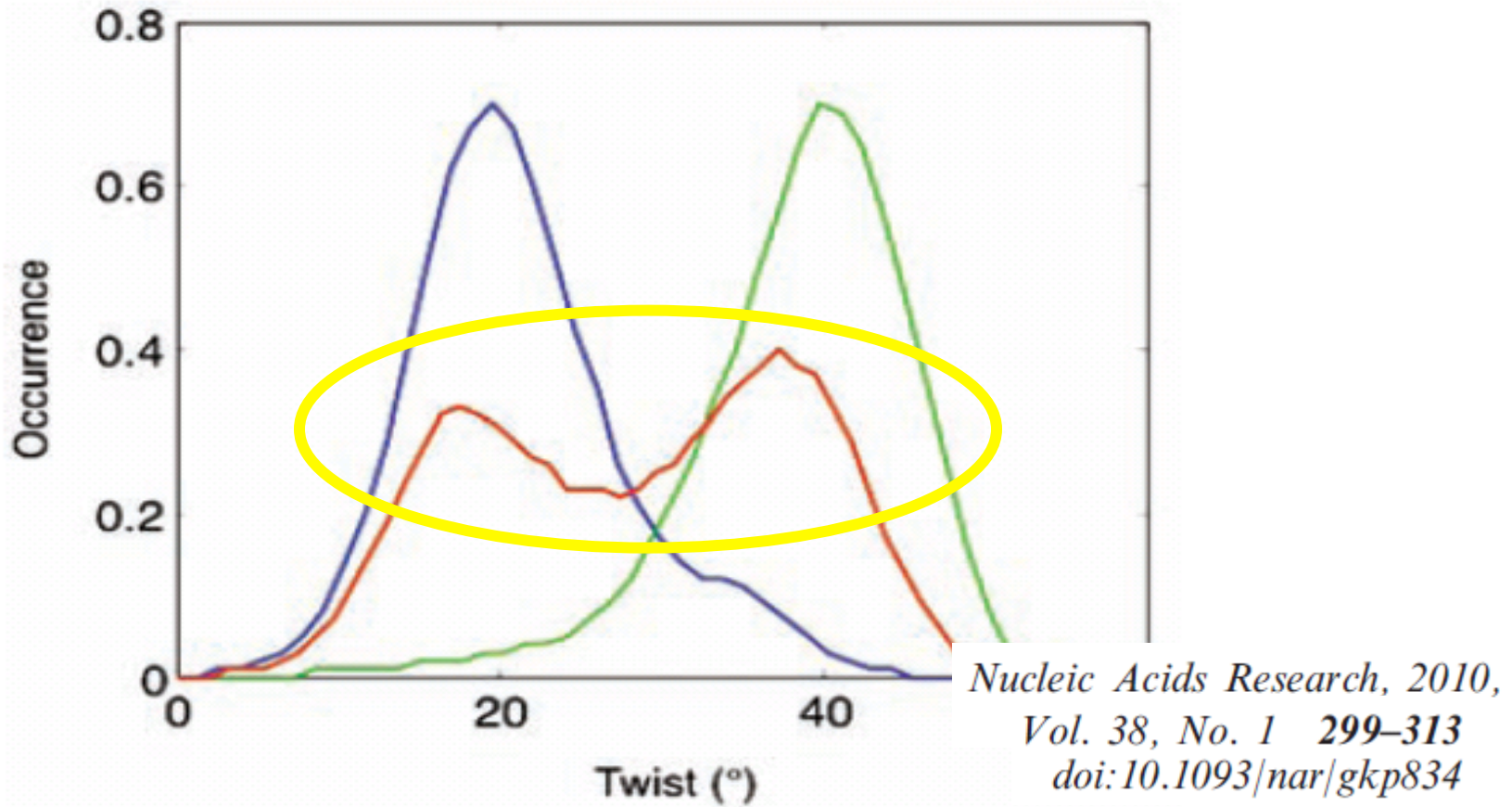


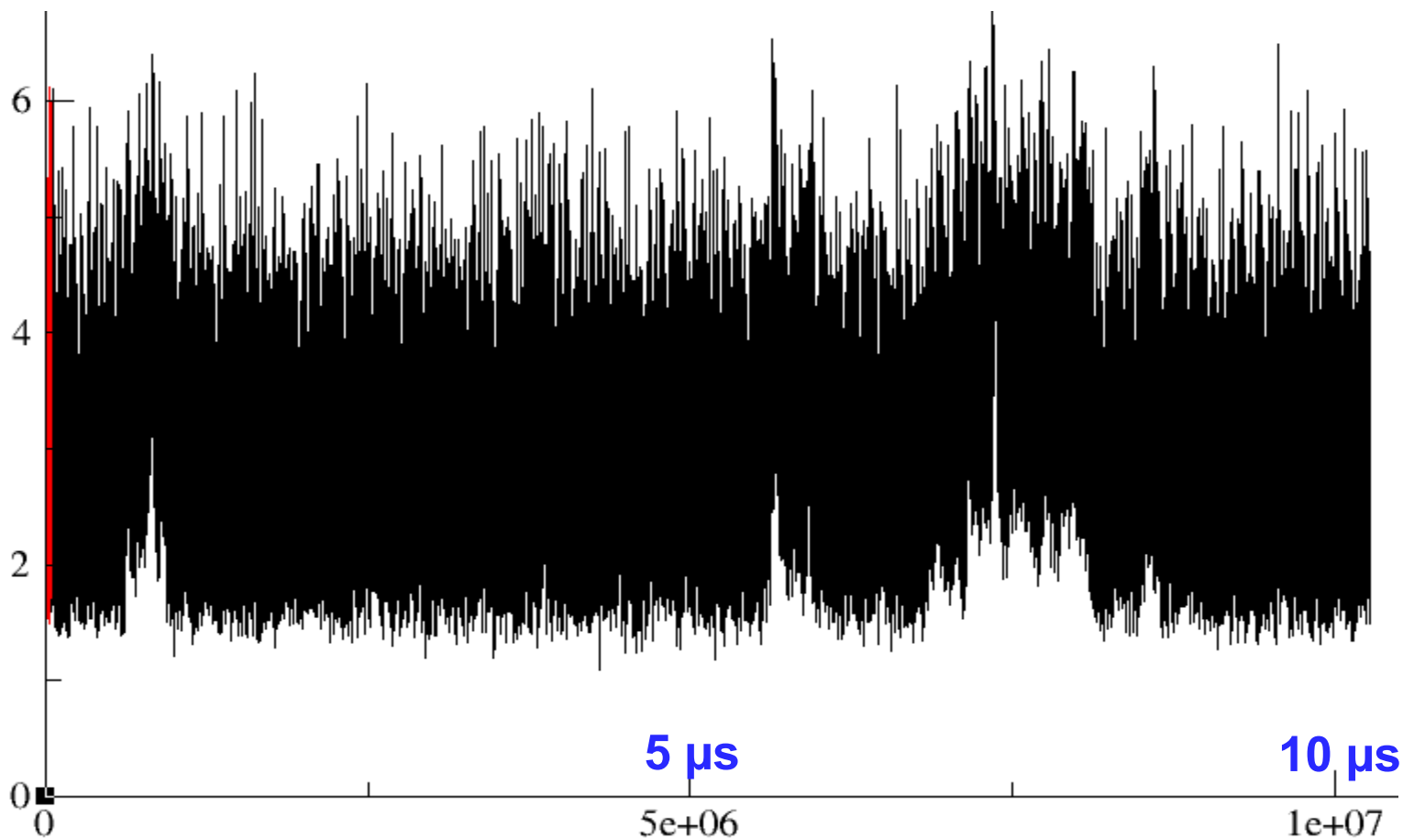
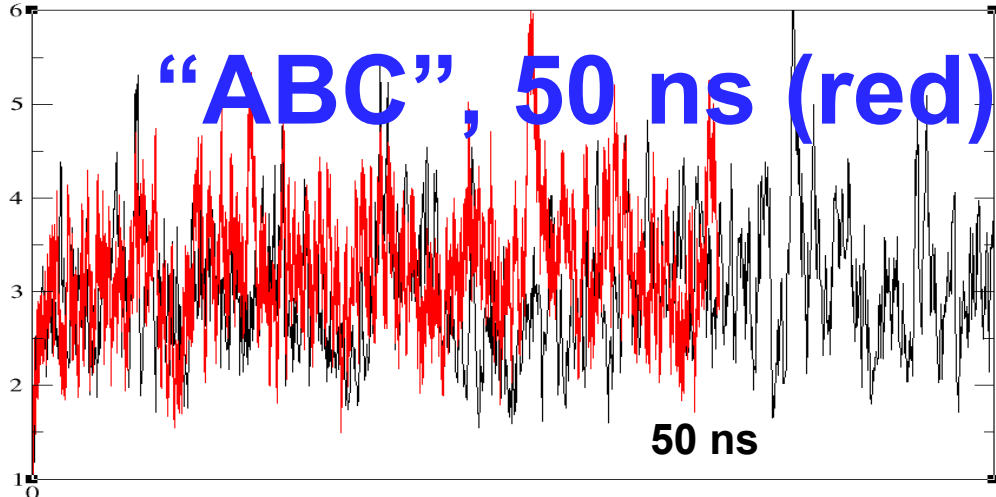
Figure 6. Distribution of CG twist (degrees) as a function of the flanking sequences: CCGA (blue); ACGT (green); ACGA (red).

“ABC”, 50 ns (red), Anton (black)

RMSD (\AA) vs. time

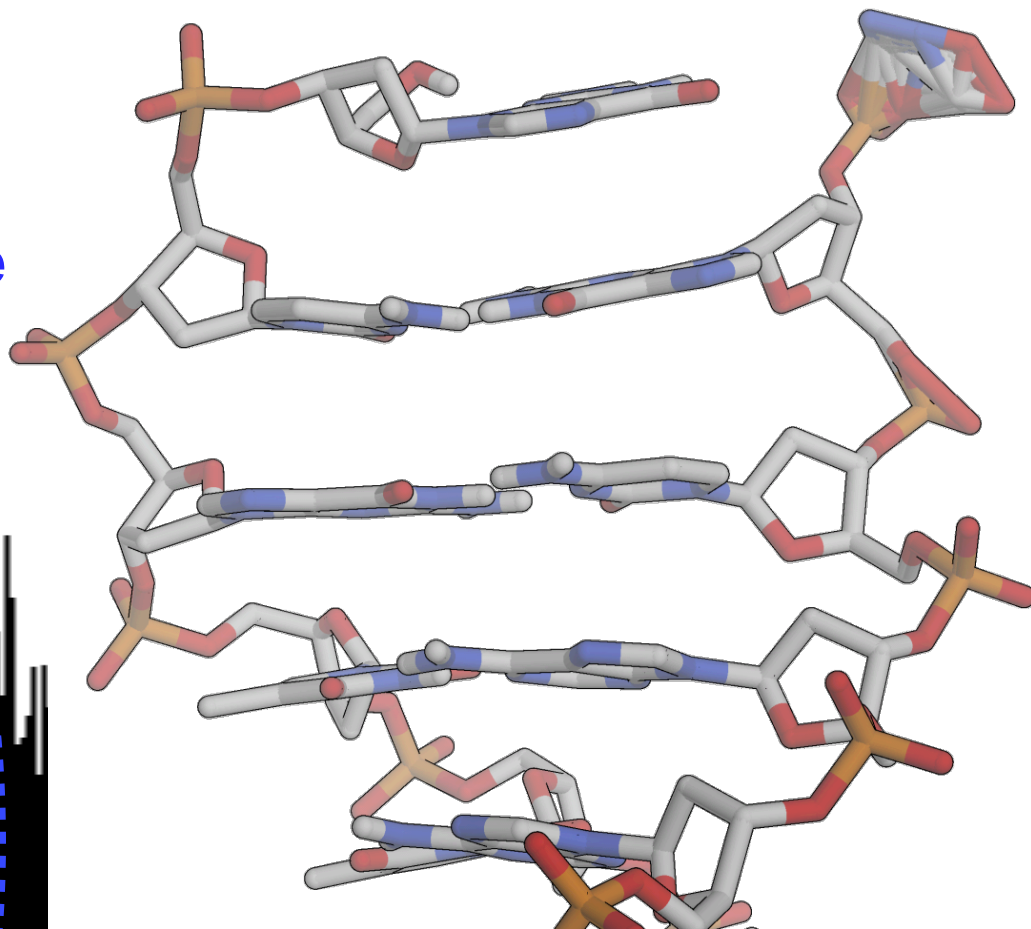
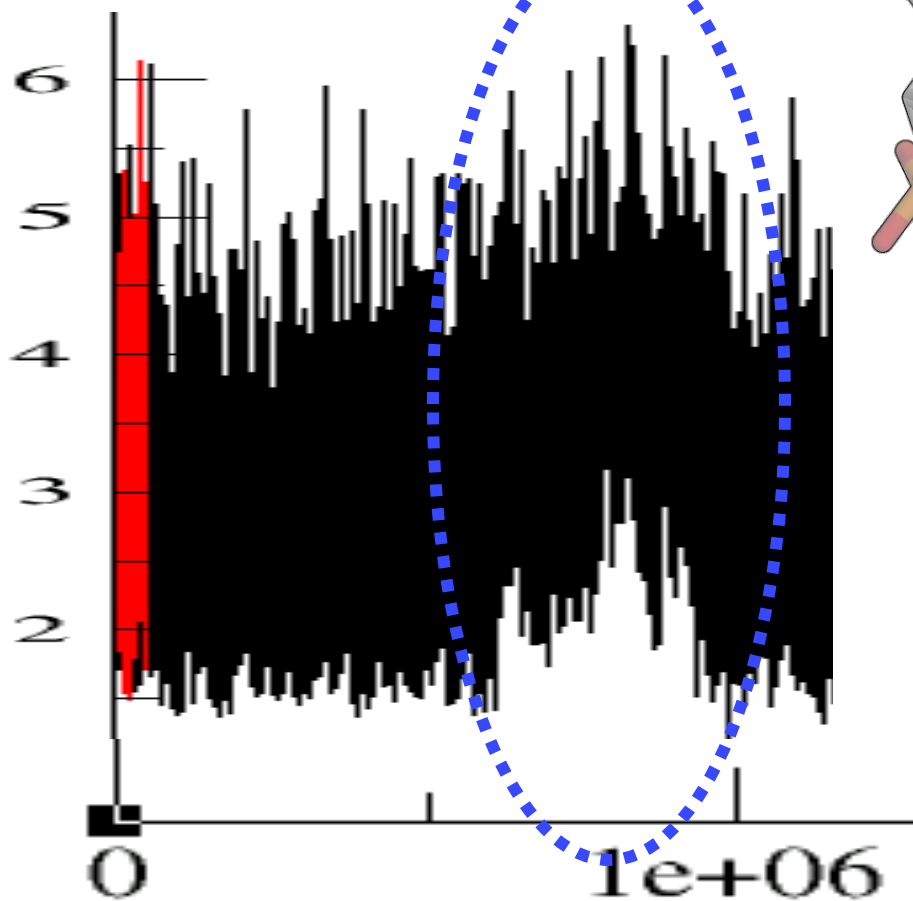
Top: Original ABC work

Below: On Anton...



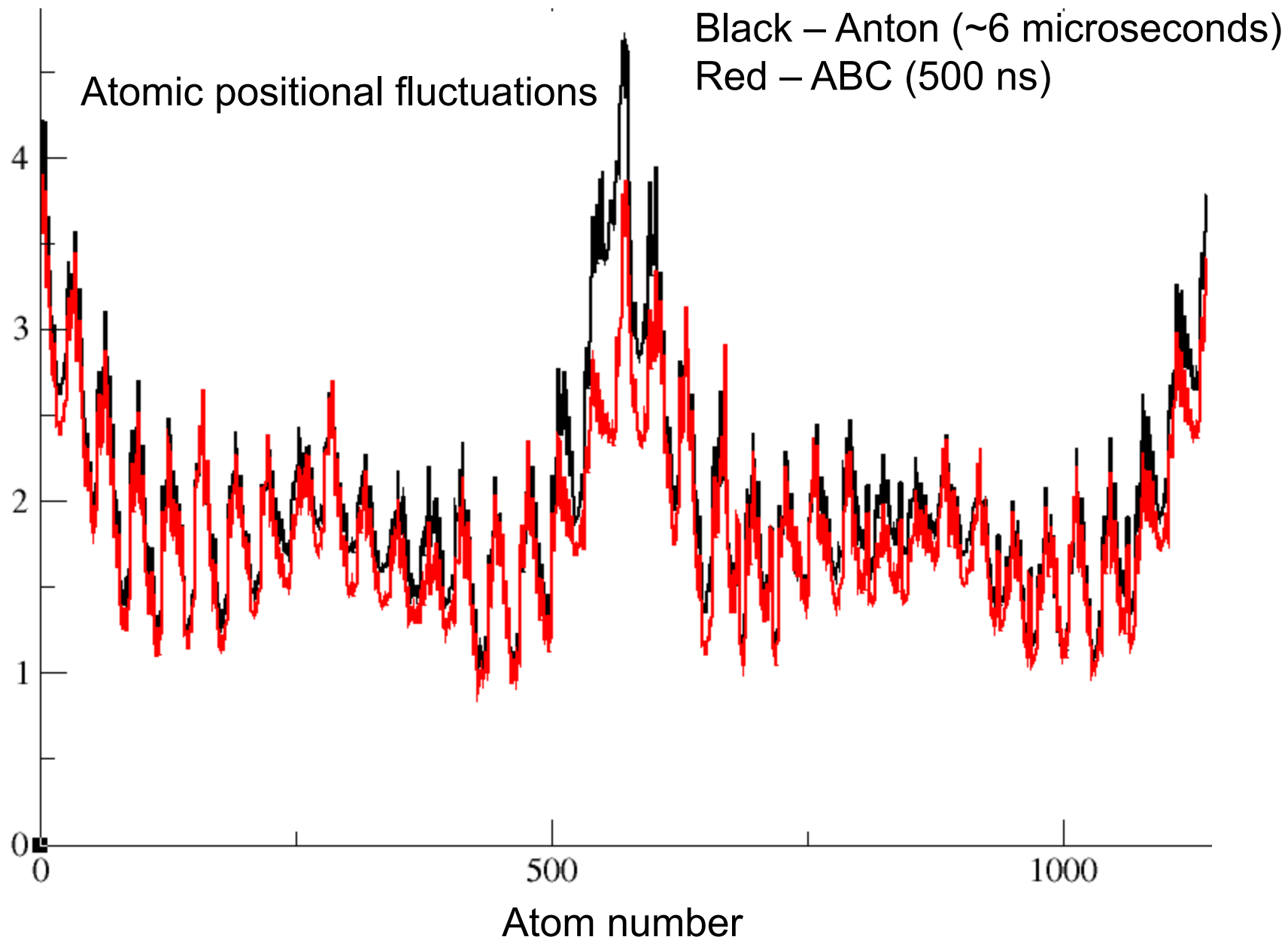
anton

terminal base
pair opening!

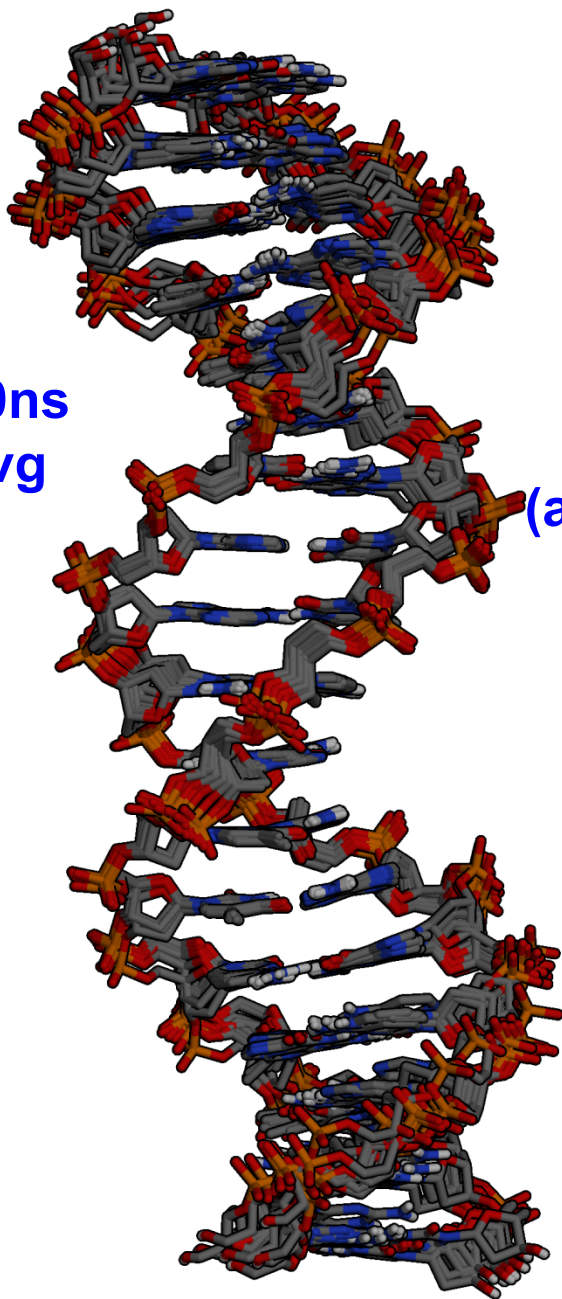


880-890 ns

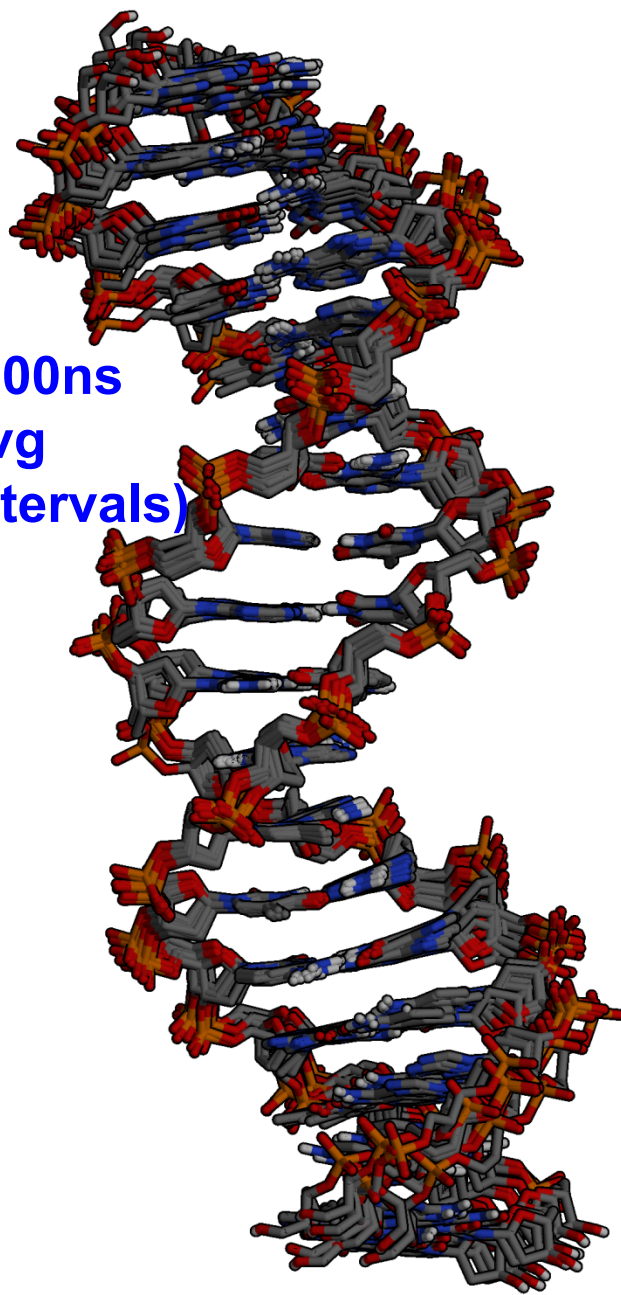
GAAC: GCACGAACGAACGAACGC



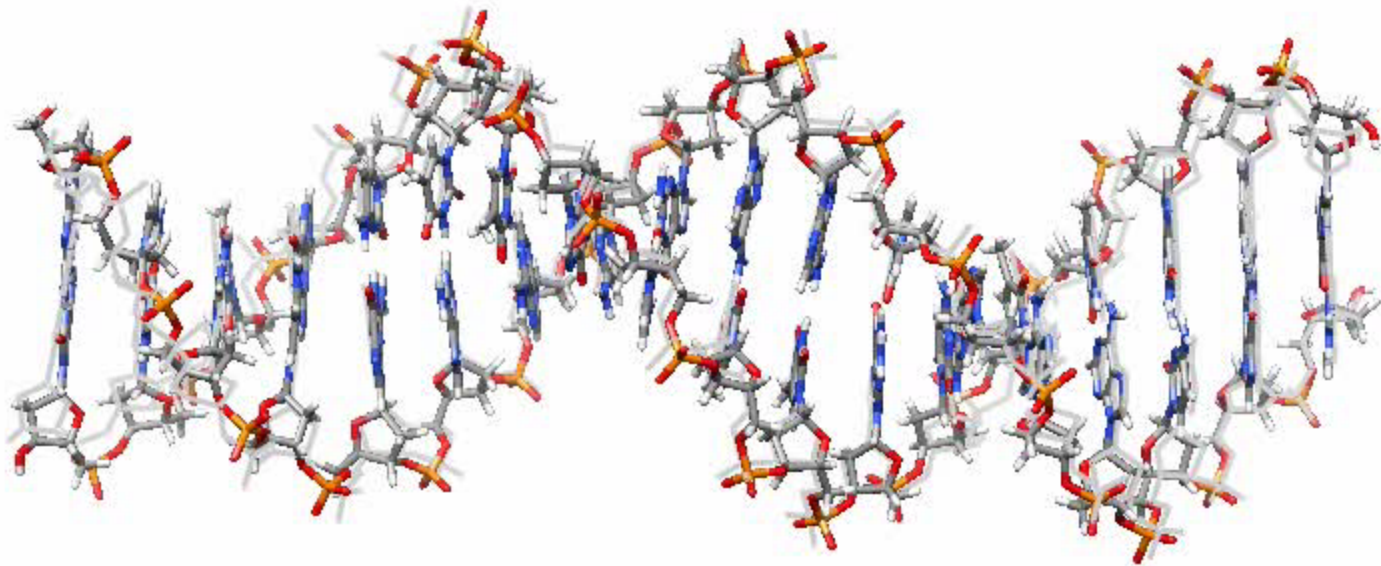
**abc, 50ns
5ns avg**



**anton, 7000ns
5ns avg
(at 500ns intervals)**



Anton run:

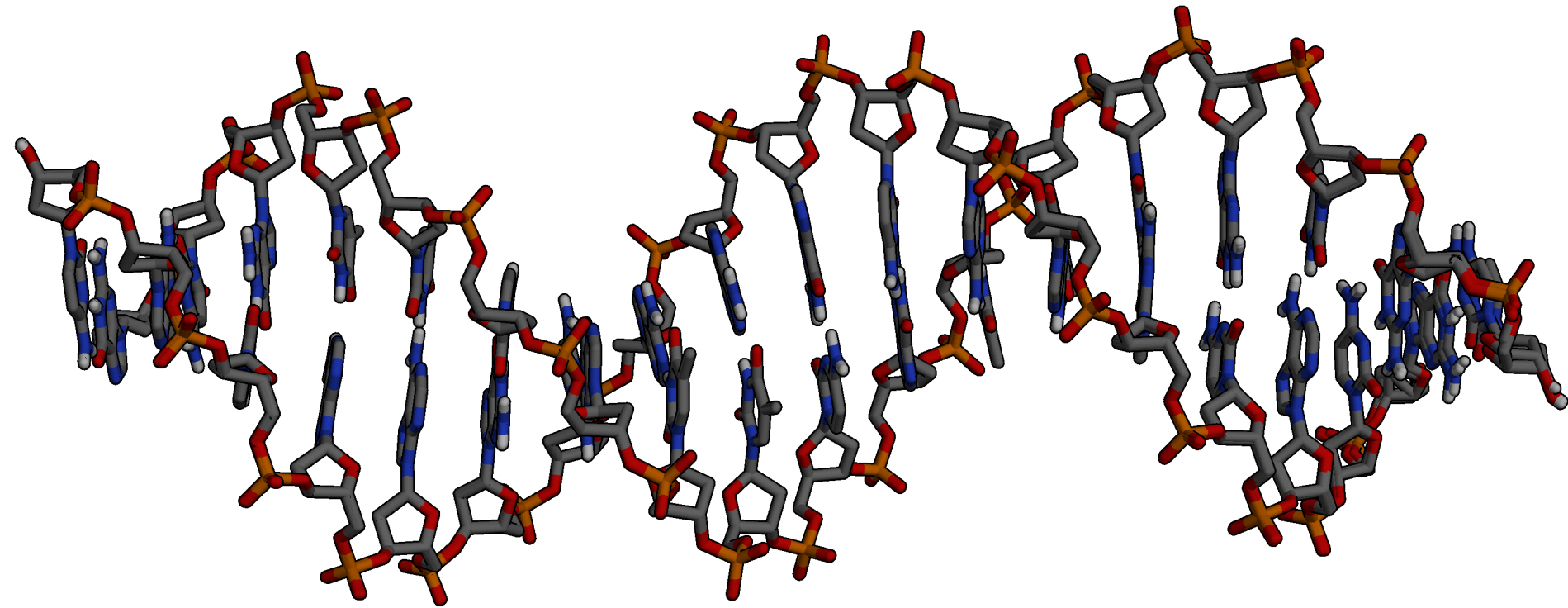


2 ns intervals, 10 ns running average, every 5th frame (~10 us).

~2010-2011

5 “average” structures overlaid @

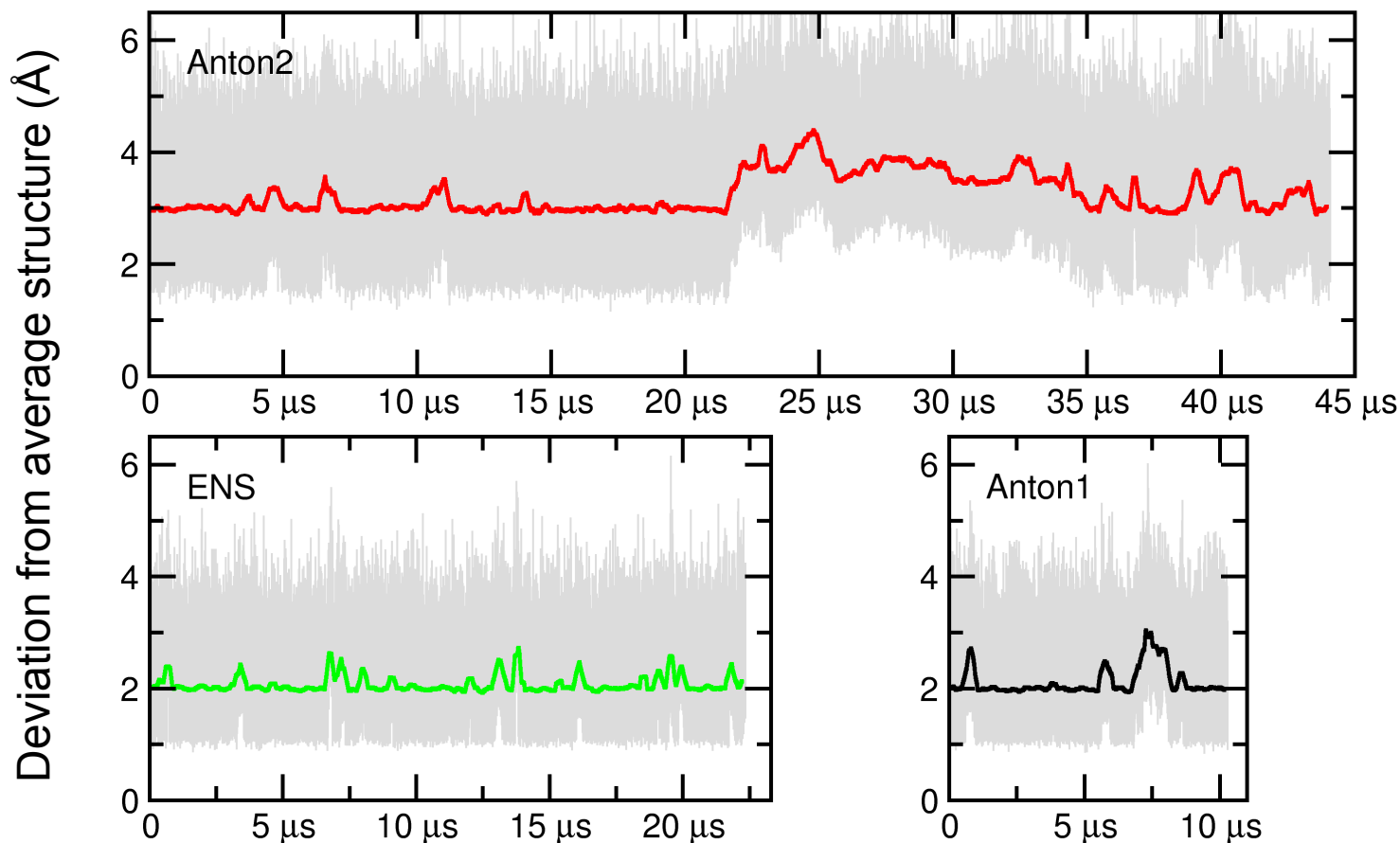
1.0-4.0 μs , 1.5-4.5 μs , 2.0-5.0 μs , 2.5-5.5 μs , 3.0-6.0 μs ...
RMSd (0.028 Å) (0.049 Å) (0.076 Å) (0.160 Å)



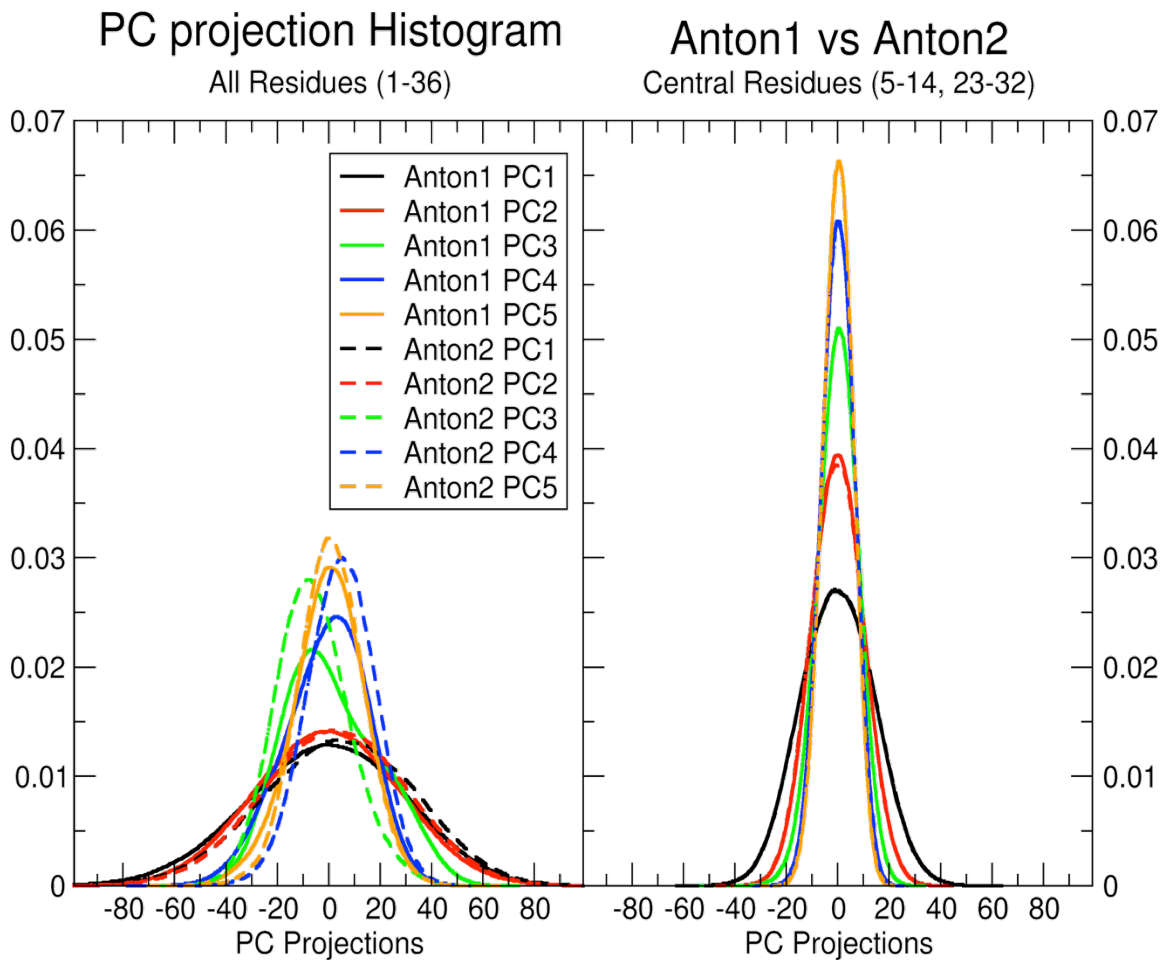
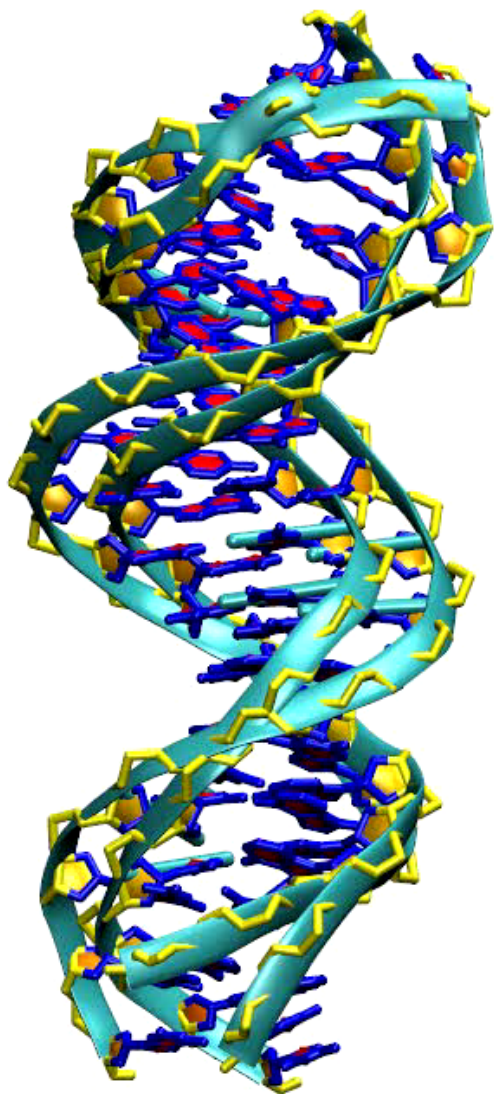
...this cannot be right, can it?
(breathing, bending, twisting, ...)

How to test?

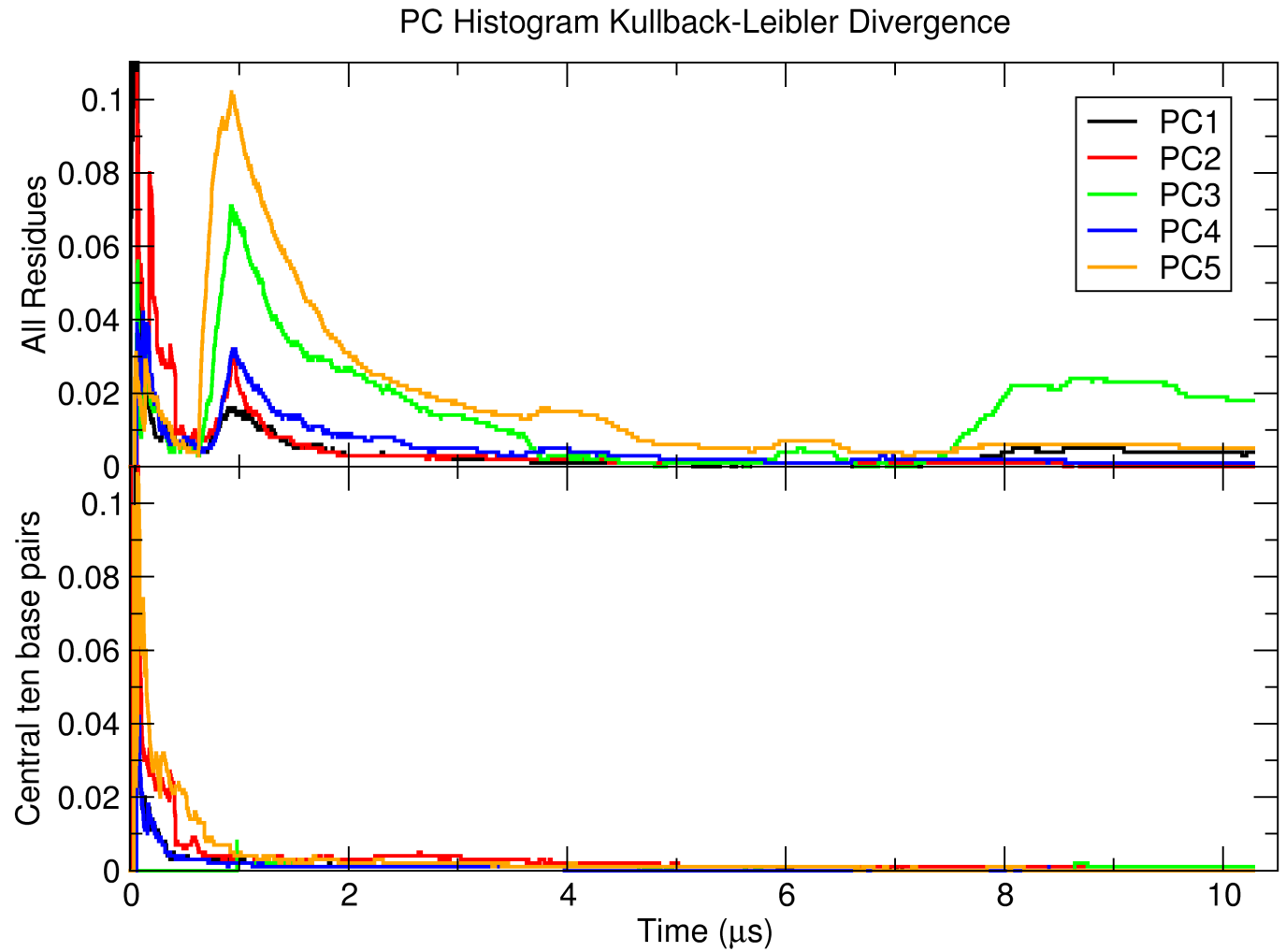
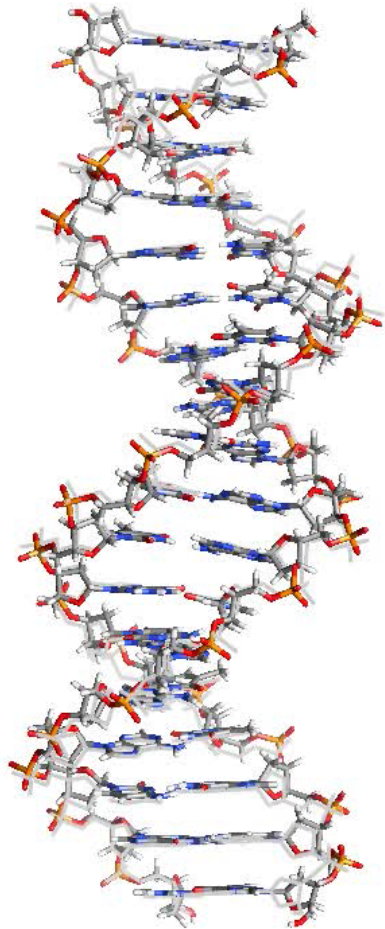
- Do a longer run on Anton (write grant, get grant, run sims, ✓) = 44 μ s
- Run an ensemble of 100 shorter simulations and aggregate = 20 μ s
- Assume Anton is wrong: Run AMBER on CPUs and GPUs (~2 years, and still not long enough, only 2-4 μ s ☹, but results are consistent ☺)



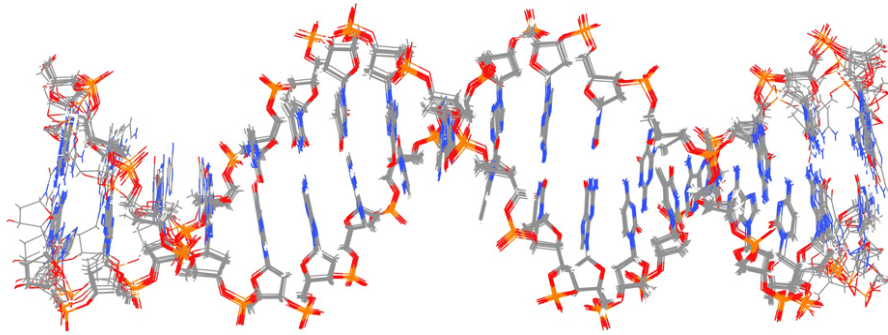
Test for convergence within and between simulations...



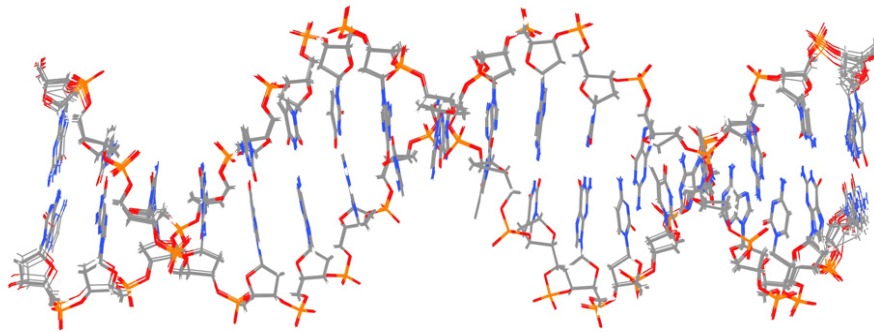
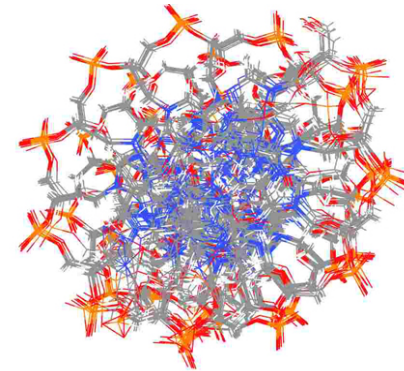
Test for convergence within and between simulations...



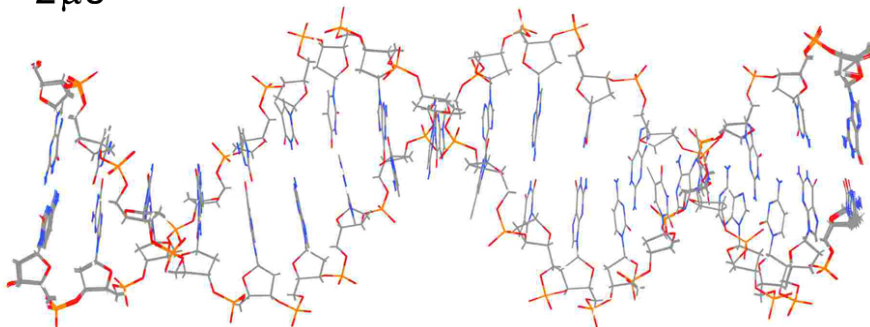
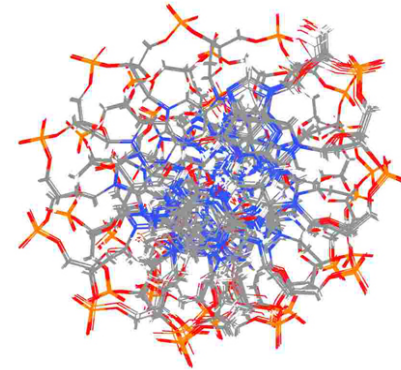
**Test for convergence within and between simulations...
(perform running average over different timescales and cluster,
showing 10 representatives)**



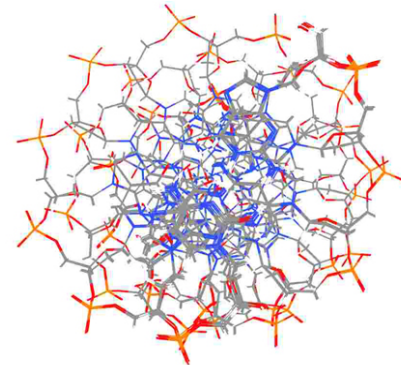
325ns



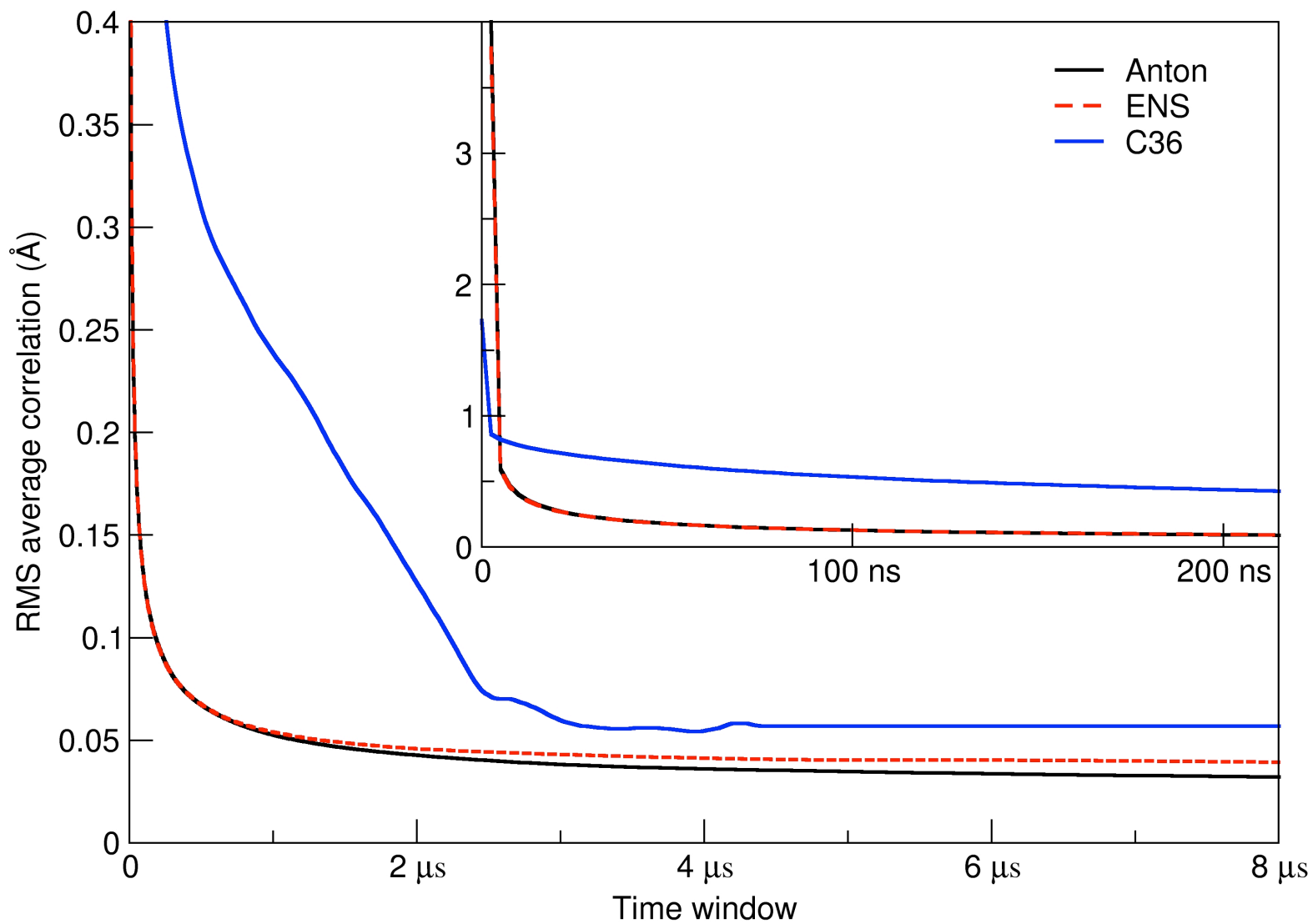
2μs



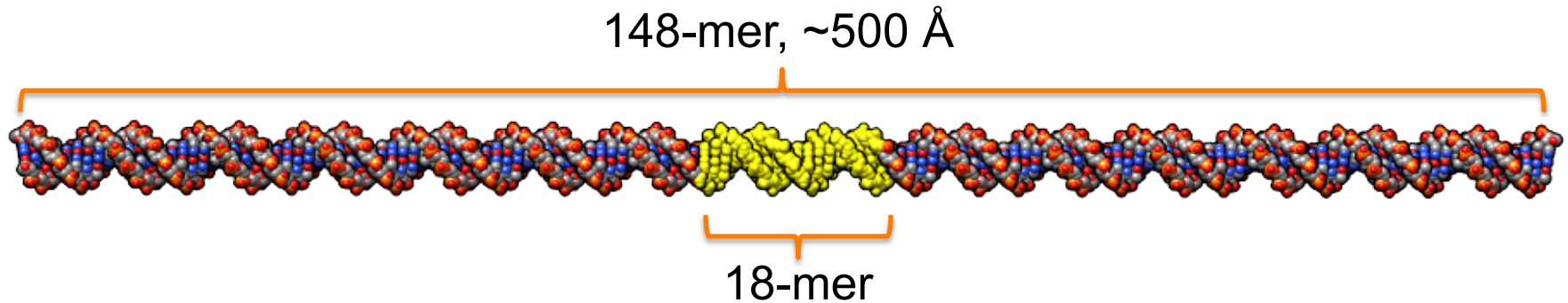
8μs



...alternative force field: CHARMM C36 runs on Blue Waters



DNA helices are relatively rigid, long persistence length



does it make sense for DNA to present consistent structure
and for regular Watson-Crick DNA to be “rigid” – YES!

Questions about recognition:

- conformational selection?
- induced fit / deformability?
- why are mismatches easily recognized?

What do we know about the dynamics of DNA helices?

¹³C & ¹⁵N NMR

triplet anisotropy decay

electron paramagnetic resonance

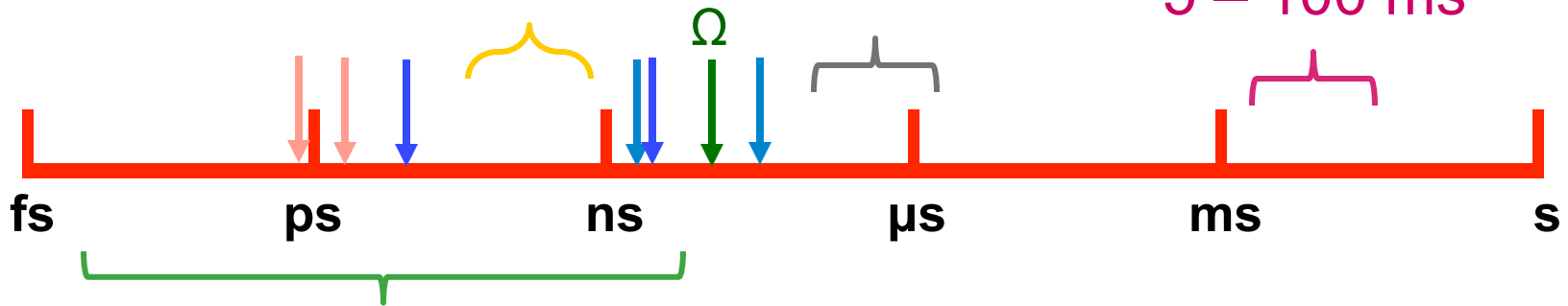
PELDOR

³¹P and/or field-cycling NMR

FT IR difference spectroscopy

⋮	⋮	X ≠ A
⋮	⋮	Y ≠ T
CG	XY	
GC	AT	> 1 ms
CG	XY	
GC	⋮	
⋮	⋮	
⋮	⋮	

NMR: base pair opening
5 – 100 ms



Berg: Dynamic stokes shift (base replaced by dye)
“power law” dynamics over 6 orders of magnitude of time
40 fs – 40 ns

What about longer timescales?

What do we know about the dynamics of DNA helices?

^{13}C & ^{15}N NMR

triplet anisotropy decay

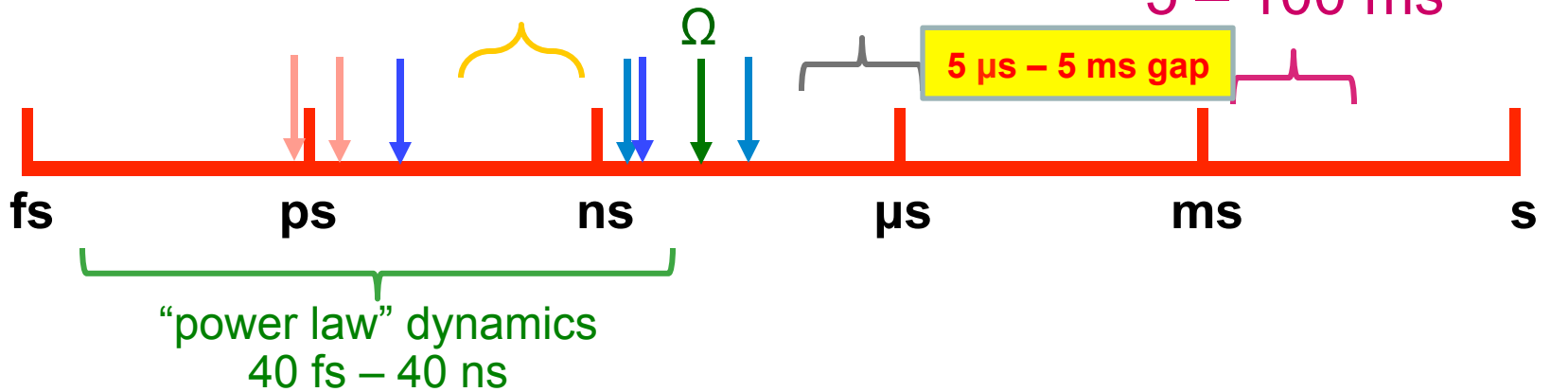
electron paramagnetic resonance

PELDOR

^{31}P and/or field-cycling NMR

FT IR difference spectroscopy

⋮	⋮	X ≠ A
⋮	⋮	Y ≠ T
CG	XY	
GC	AT	> 1 μs
CG	XY	
GC	⋮	
⋮	⋮	
⋮	⋮	



is this “gap” in dynamics real?

What do we know about the dynamics of DNA helices?

¹³C & ¹⁵N NMR

triplet anisotropy decay

electron paramagnetic resonance

PELDOR

³¹P and/or field-cycling NMR

FT IR difference spectroscopy

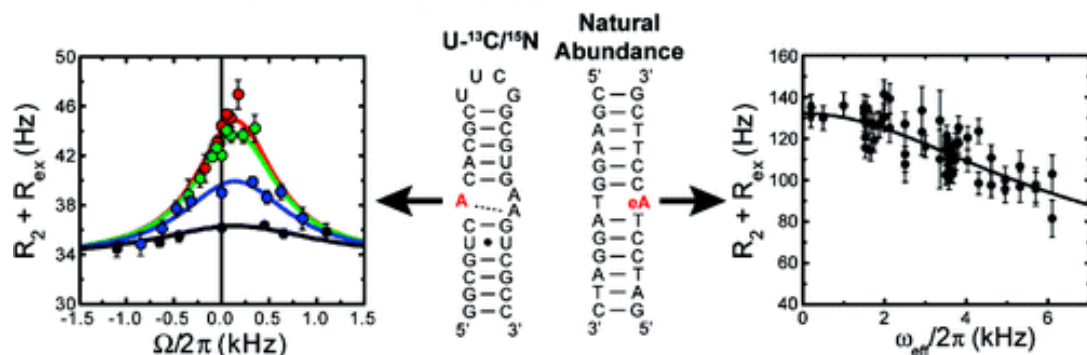
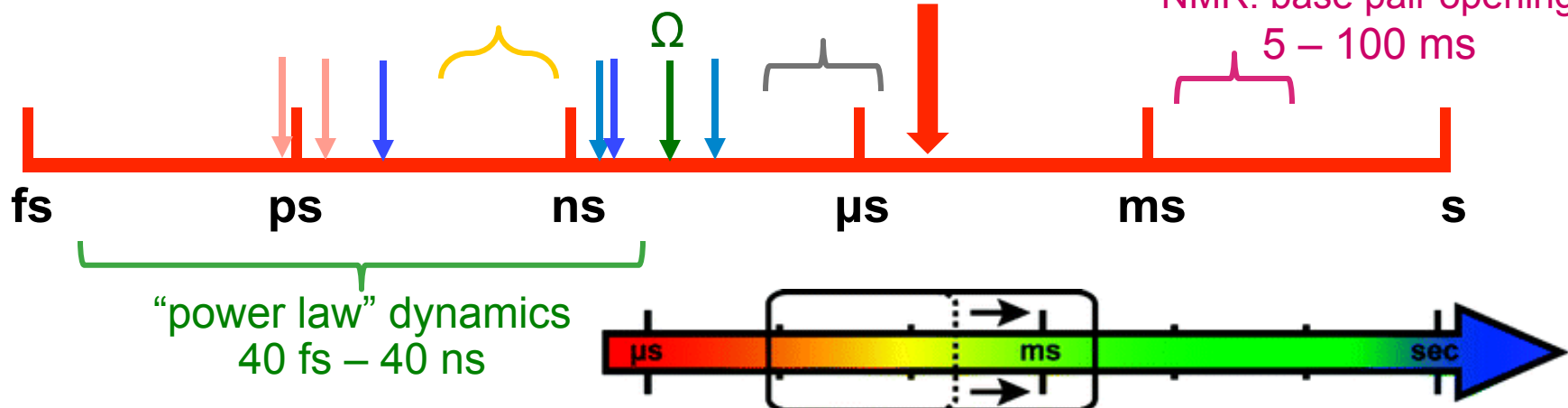
$\begin{matrix} \vdots \\ \vdots \\ \text{CG} \\ \text{GC} \\ \text{CG} \\ \text{GC} \\ \vdots \\ \vdots \end{matrix}$
 $\begin{matrix} \vdots \\ \vdots \\ \text{XY} \\ \text{AT} \\ \text{XY} \\ \vdots \\ \vdots \end{matrix}$
 &
 $\begin{matrix} \vdots \\ \vdots \\ \text{X} \neq \text{A} \\ \text{Y} \neq \text{T} \\ \vdots \\ \vdots \end{matrix}$

 $> 1 \mu\text{s}$

selective off-resonance
 $R_{1\rho}$ carbon relaxation
 $26 \pm 8 \mu\text{s}$

for mismatches!

NMR: base pair opening
 5 – 100 ms



What do we know about the dynamics of DNA helices?

:::: X ≠ A
 :::: Y ≠ T
 CG XY
 GC & AT > 1 μs
 CG XY
 GC ::
 ::::
 ::::

¹³C & ¹⁵N NMR

triplet anisotropy decay

electron paramagnetic resonance

PELDOR

³¹P and/or field-cycling NMR

FT IR difference spectroscopy

selective off-resonance

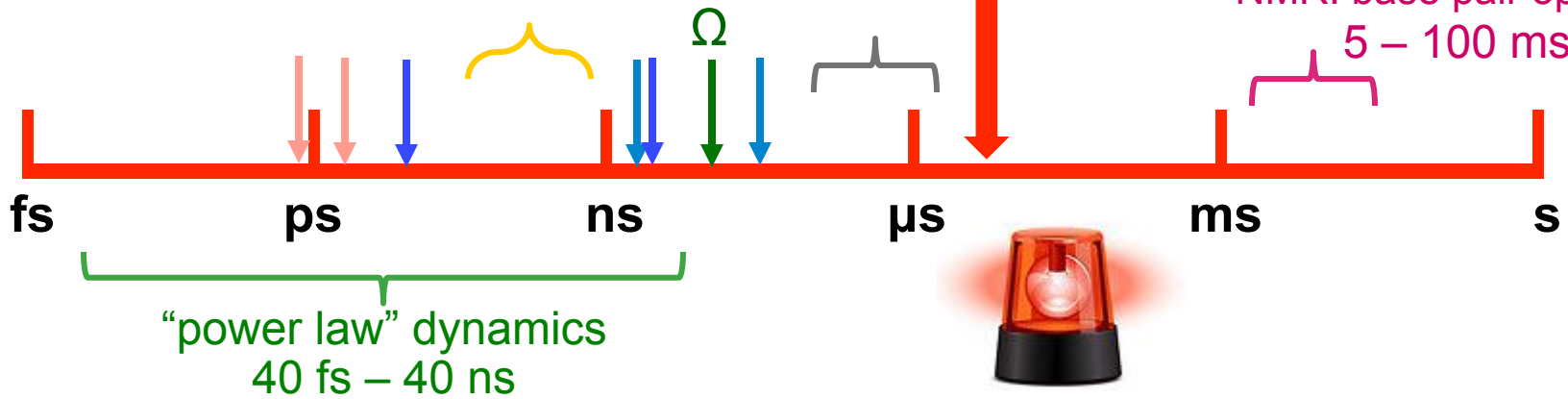
$R_{1\rho}$ carbon relaxation

$26 \pm 8 \mu\text{s}$

for mismatches!

NMR: base pair opening

5 – 100 ms



Questions about recognition:

- conformational selection?
- induced fit / deformability?
- why are mismatches easily recognized?

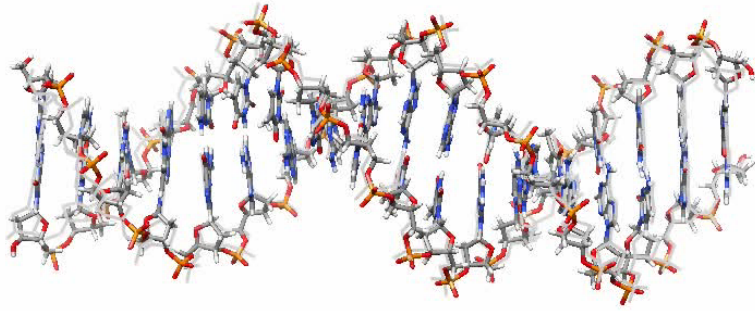
no, decay is too fast!

requires bp opening

timescale mismatch

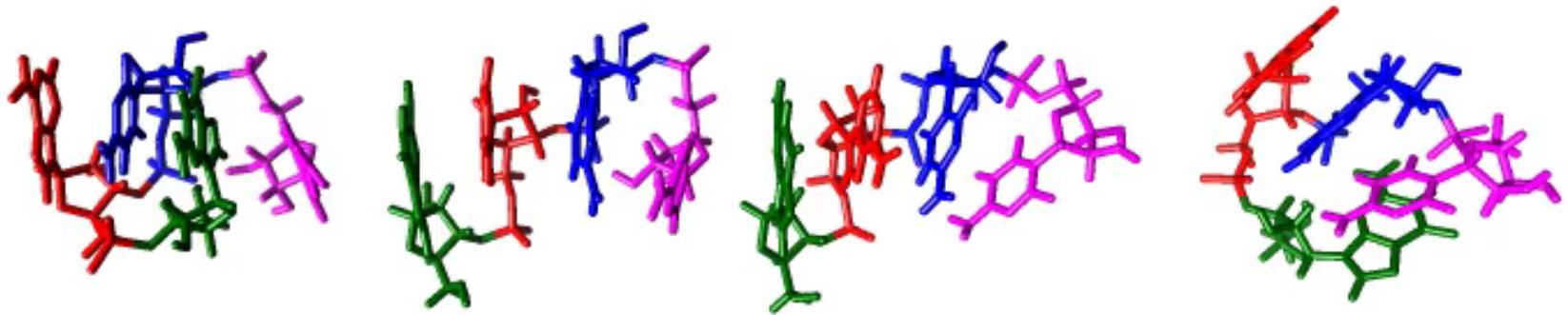
Today: two “long-time-to-develop” short stories...

- ✓ can we converge DNA duplex structure/dynamics?



Anonymous NIH R-01 reviewer in 2005:
“One has to wonder how many relatively short MD simulations have to be performed on short DNA fragments before what can be learned will have been learned...”

- ✓ sampling RNA structure *accurately* is difficult

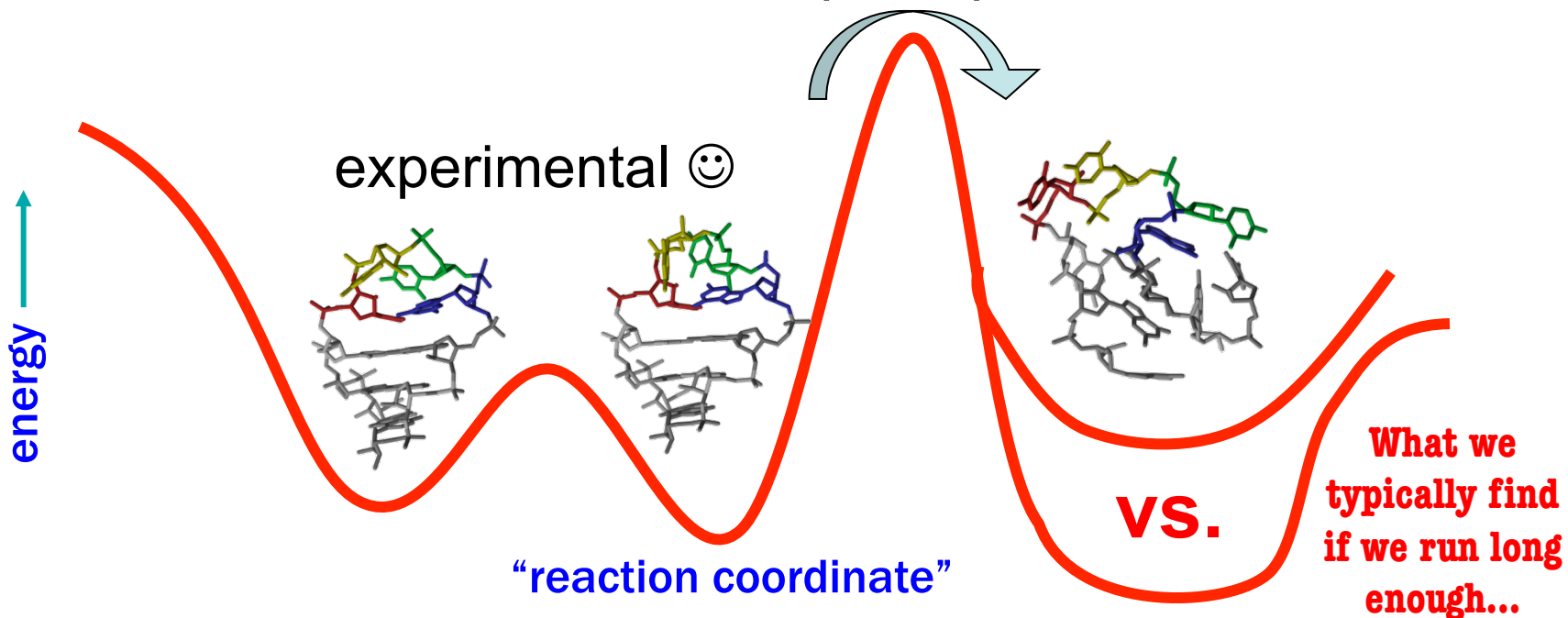


are the force fields reliable?
(free energetics, sampling, dynamics)

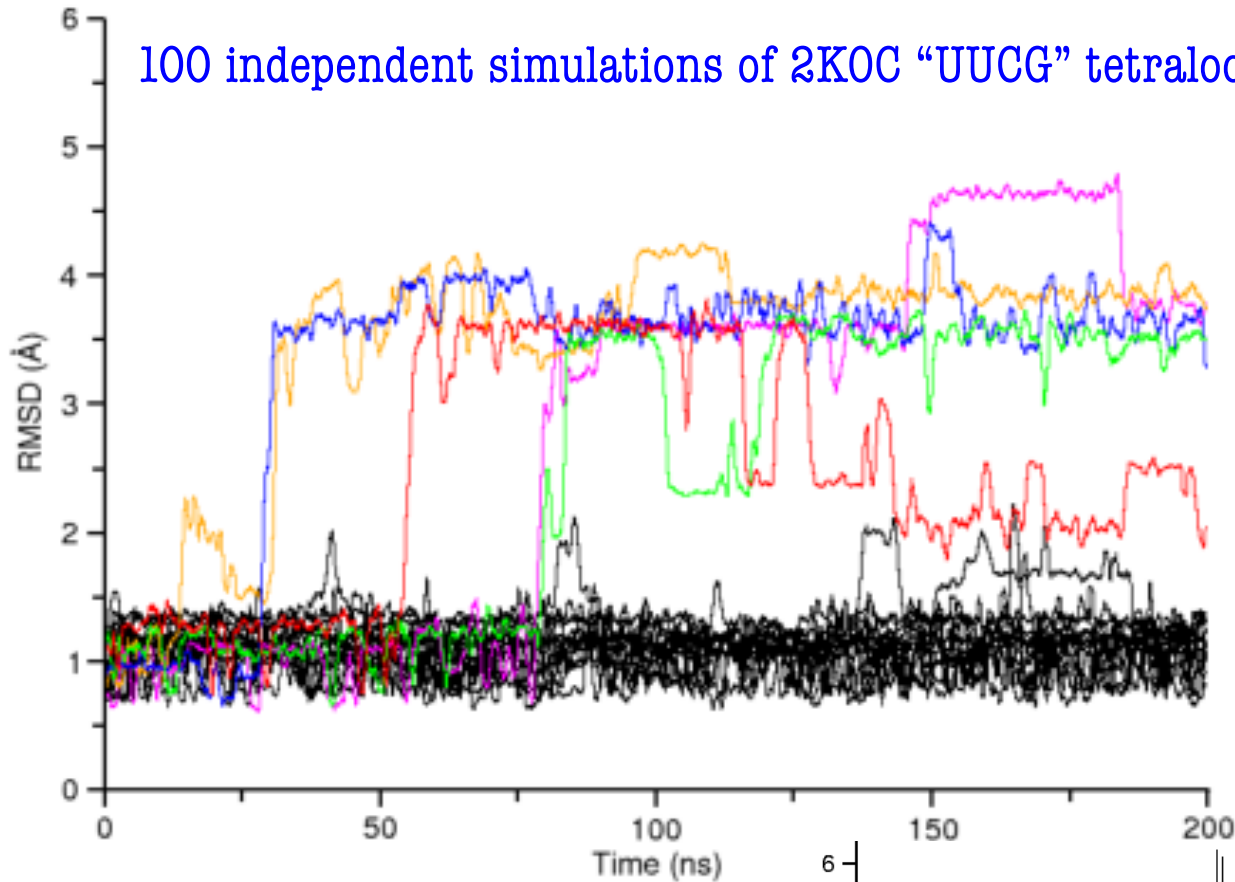
NMR structures of DNA & RNA **all tetraloops** **crystal simulations**

RNA motifs **quadruplexes**
RNA-drug interactions

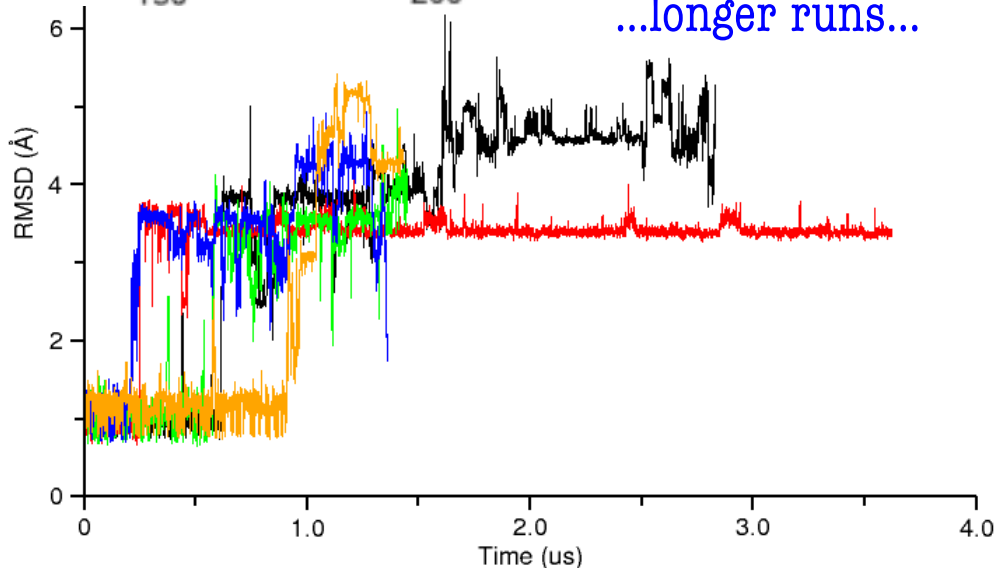
Computer power?



100 independent simulations of 2K0C "UUCG" tetraloop



...longer runs...

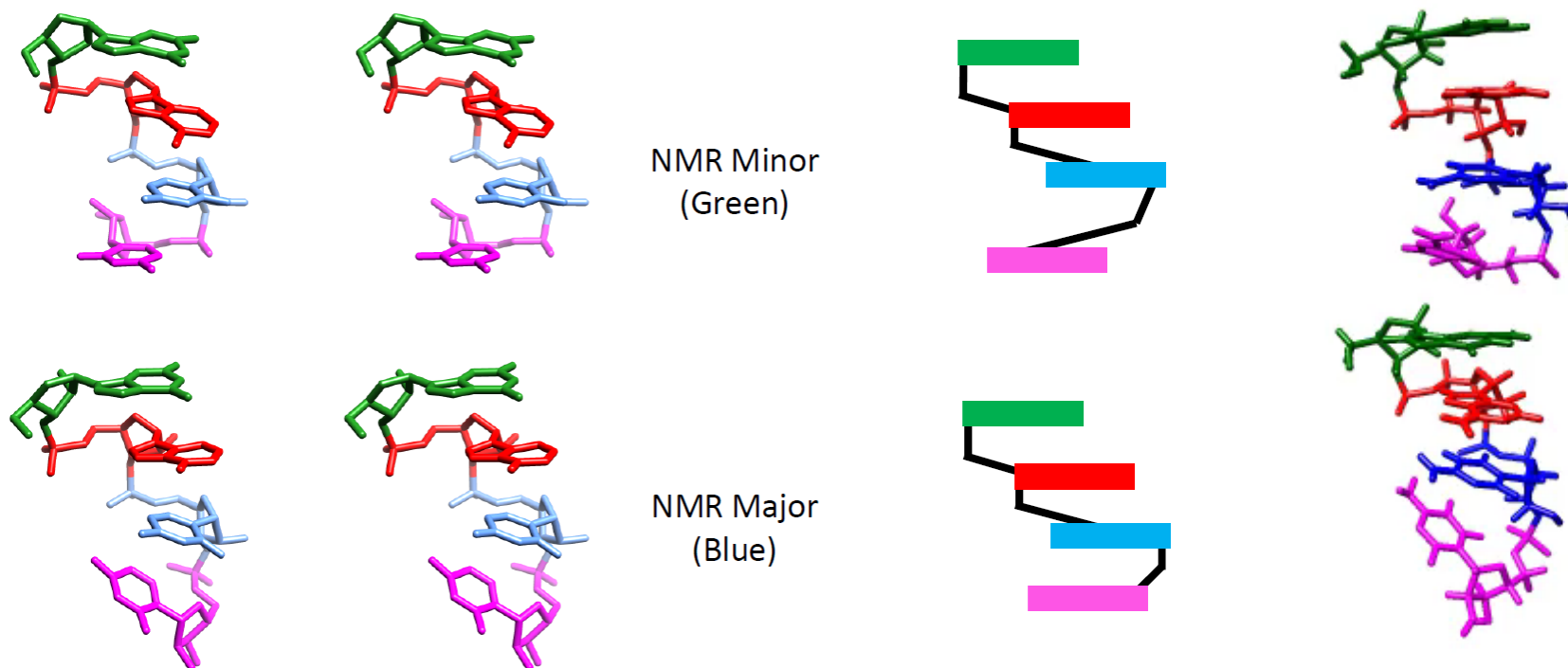


**Limited sampling
& too complex:
Is there a simpler
set of systems?**

...a system where we can get complete sampling

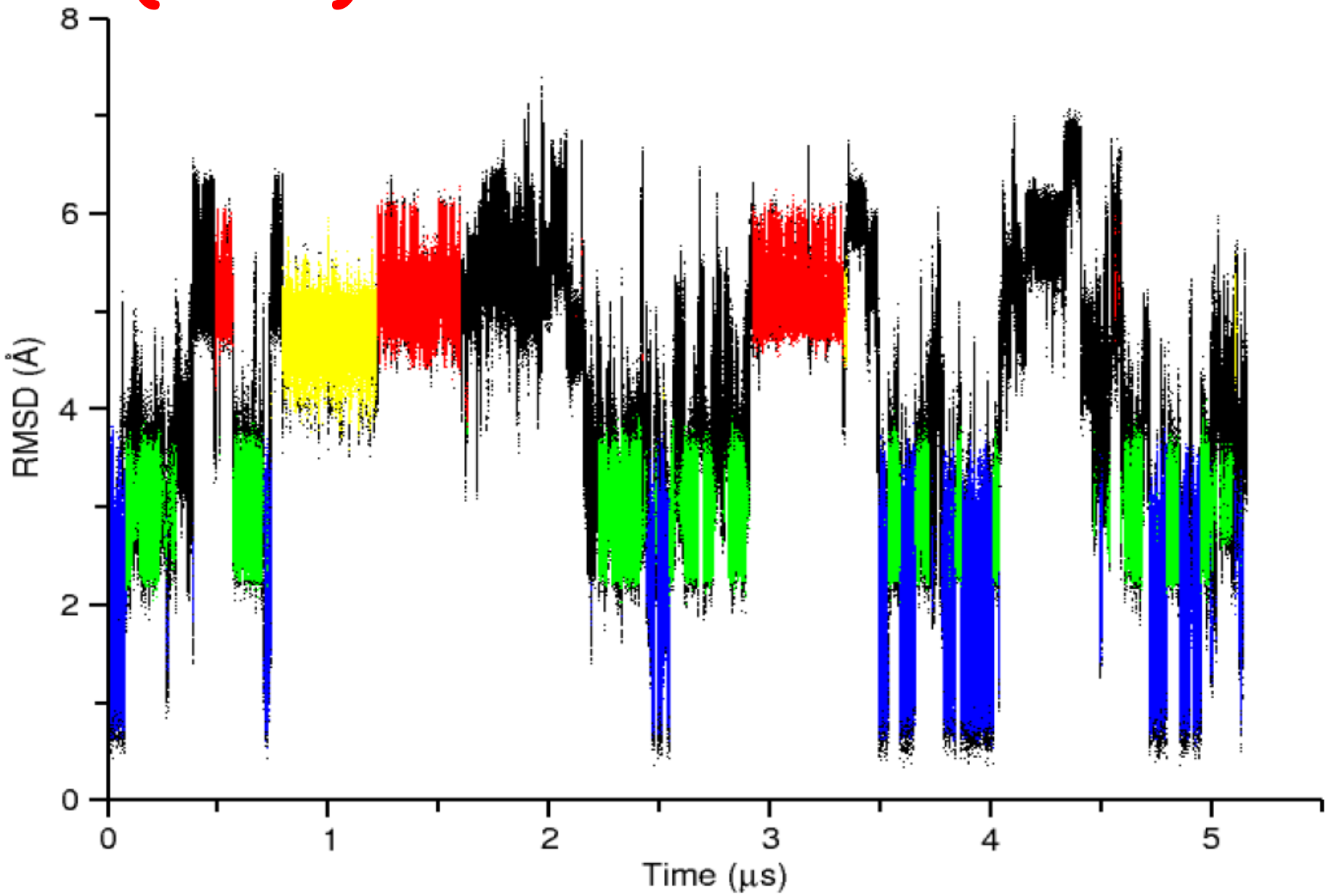
r(GACC) tetranucleotide

[Turner / Yildirim]

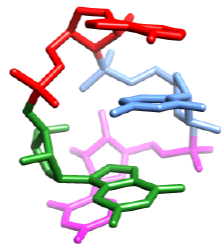
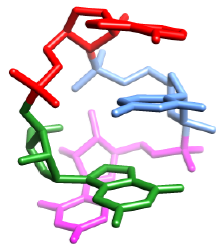


NMR suggests two dominant conformations...
...compare to MD simulations in explicit solvent

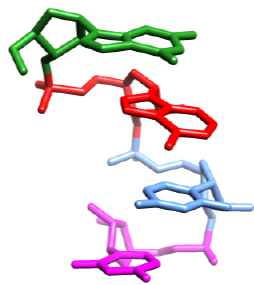
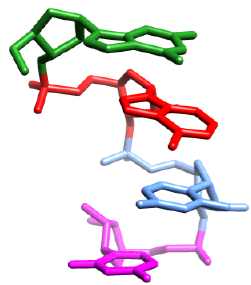
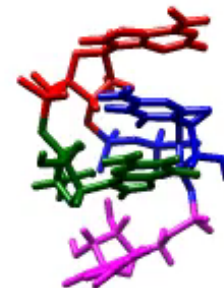
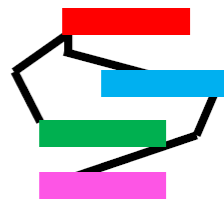
r(GACC) tetranucleotide: AMBER ff12



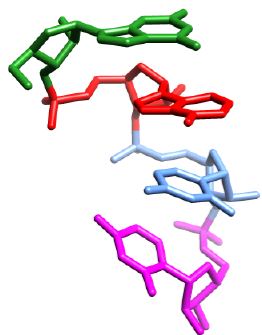
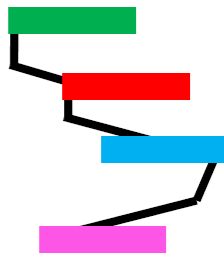
< explicit solvent time-contiguous MD >



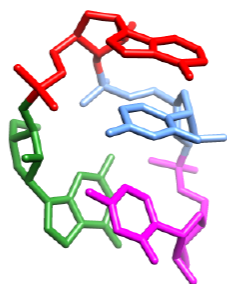
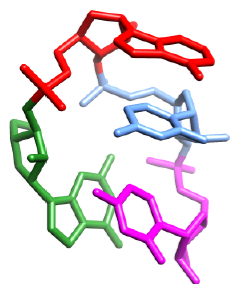
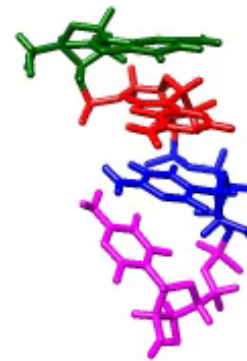
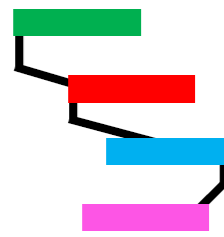
Intercalated
(Red)



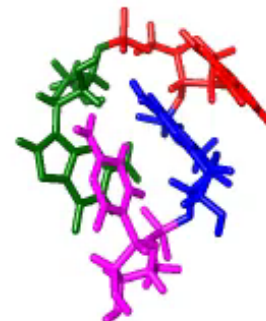
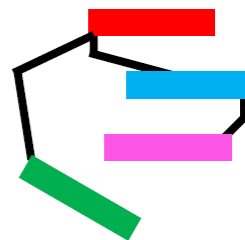
NMR Minor
(Green)



NMR Major
(Blue)



Inverted
(Yellow)



...still need more sampling!

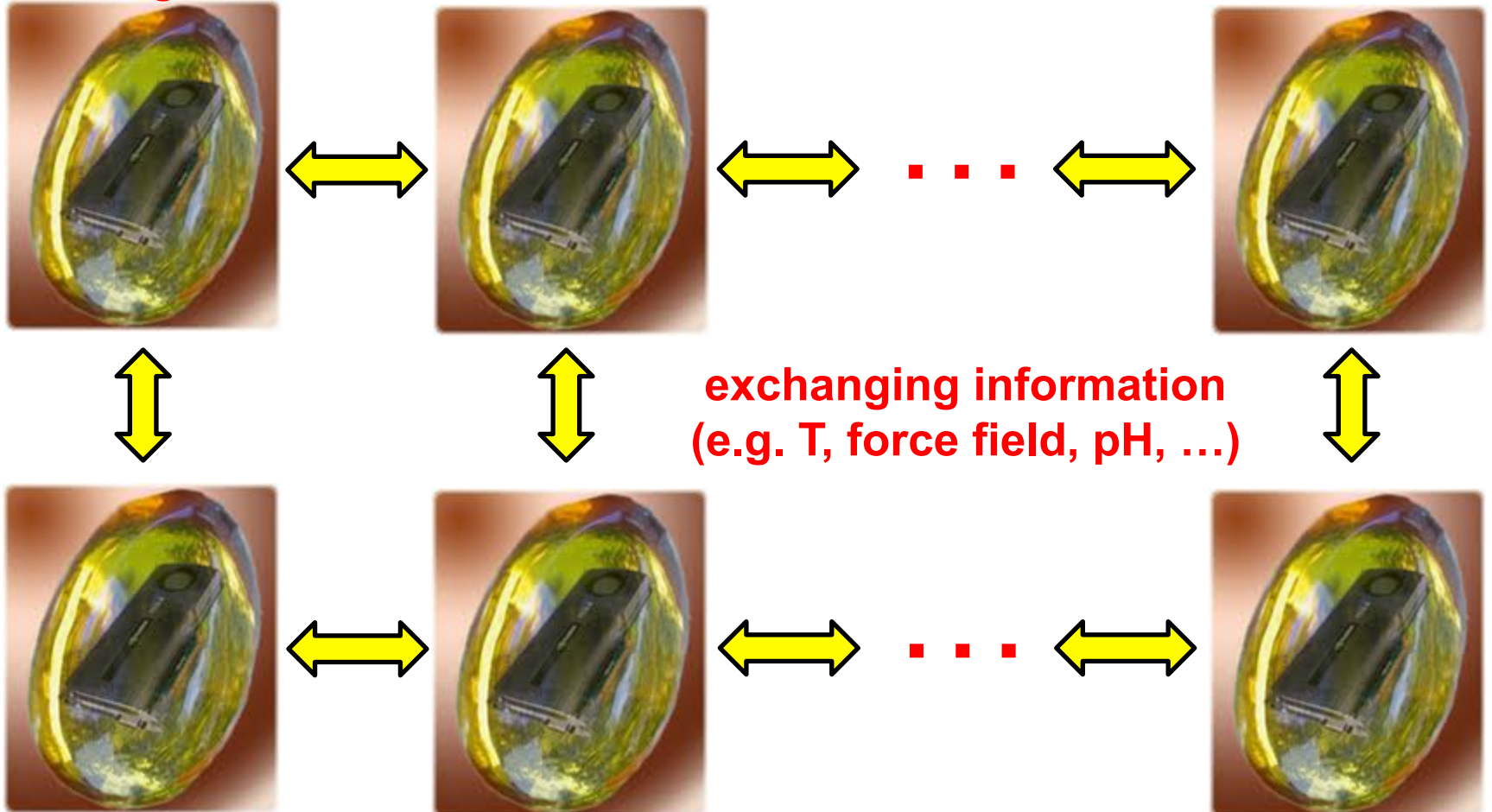
(enablers)

- **strong GPU performance of AMBER/PMEMD**
- **good replica exchange functionality**
- **access to Keeneland, Stampede, Blue Waters, ...**



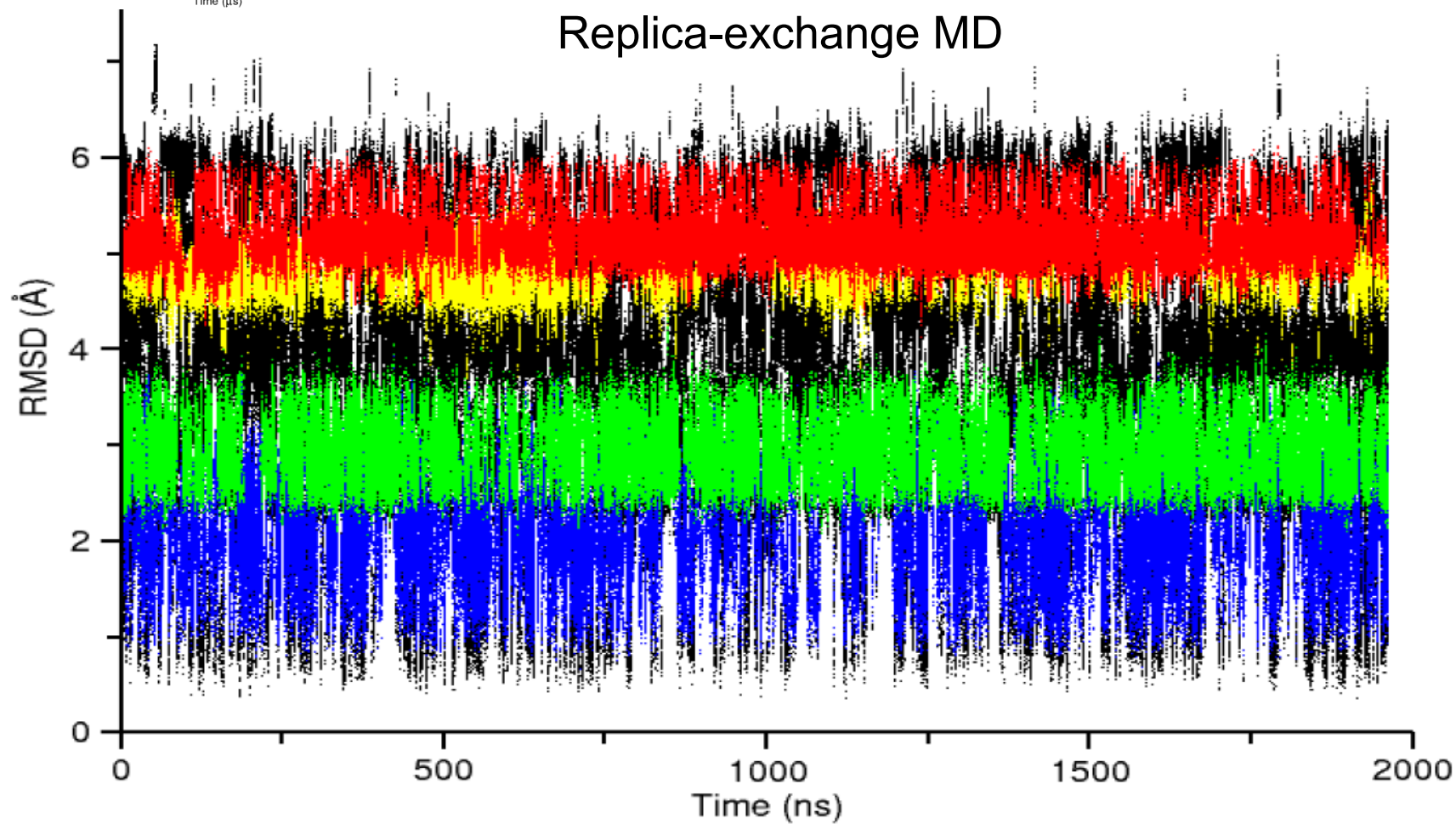
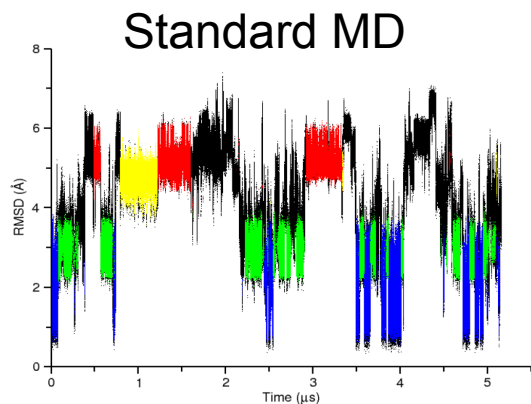
Blue Waters PRAC: The main goals are to hierarchically and tightly couple a series of optimized molecular dynamics engines to fully map out the conformational, energetic and chemical landscape of RNA.

independent ||
MD engines



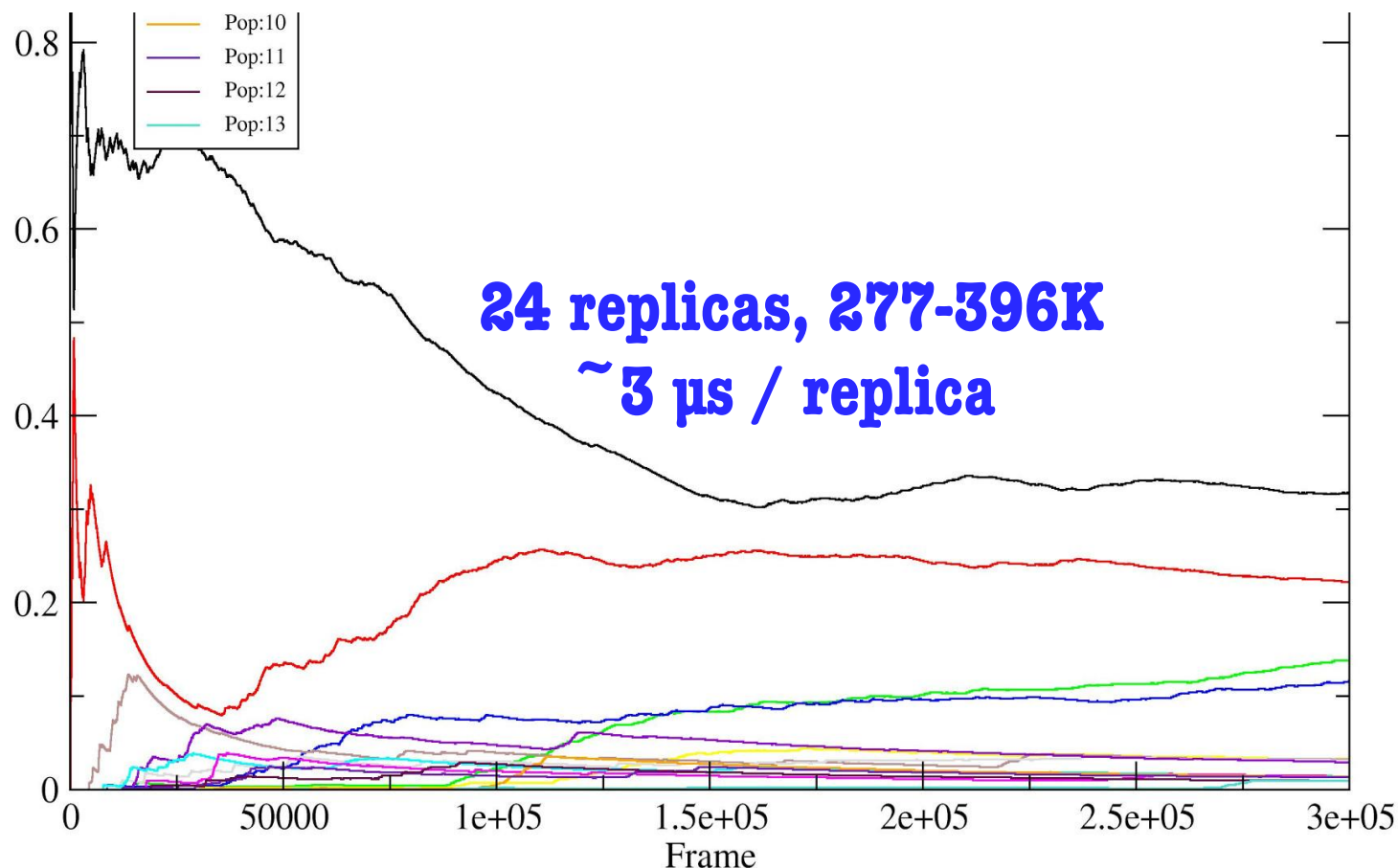
Current players: Cheatham, Roitberg, Simmerling, York, Case

r(GACC) tetranucleotide



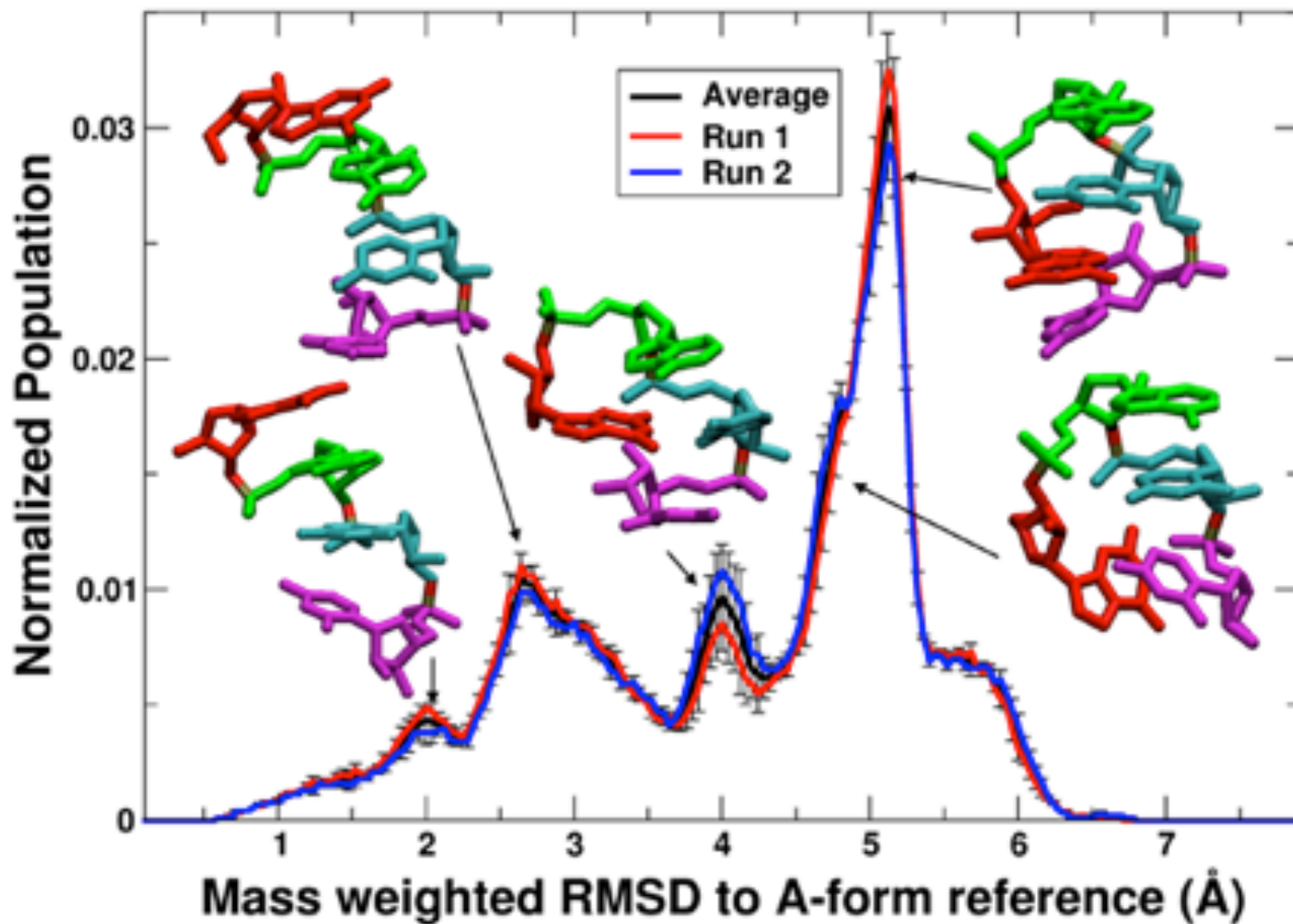
Other issues:

- **T-REMD still not “fully” converged (depending on def.)**

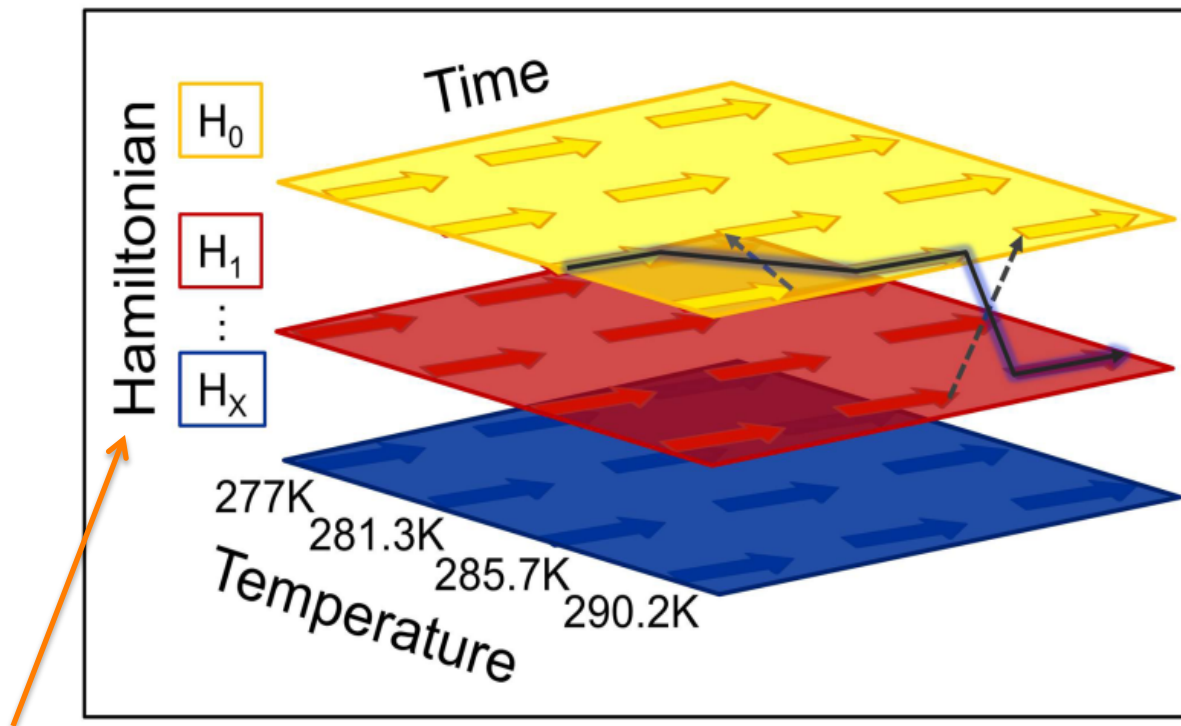


- **Not only are those four conformations populated, more like ~20+ populated > 1%**

RMSD distribution profiles: Distance from A-form reference



multi-D REMD – Bergonzo / Roe, Roitberg / Swails



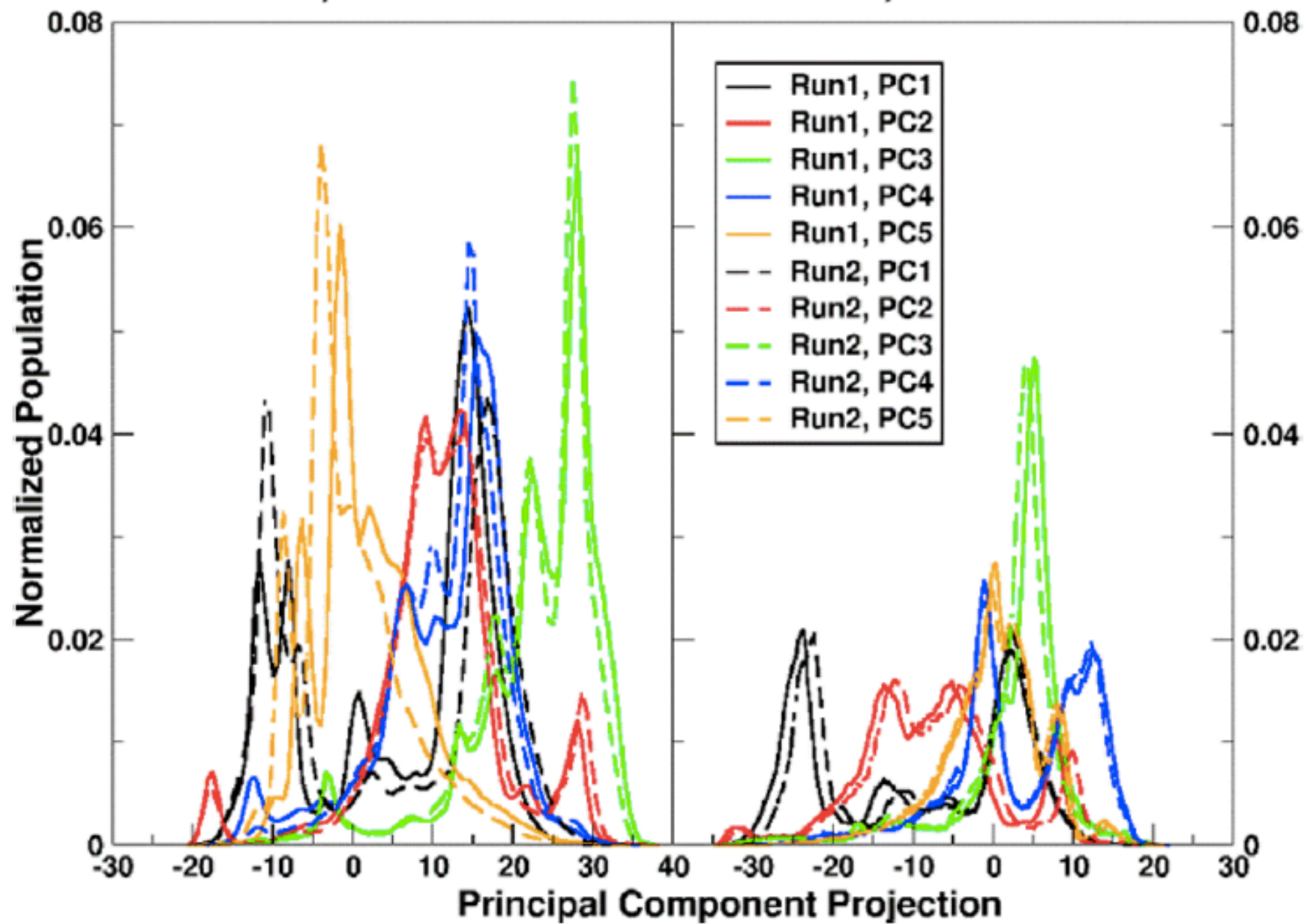
Change in “energy representation”

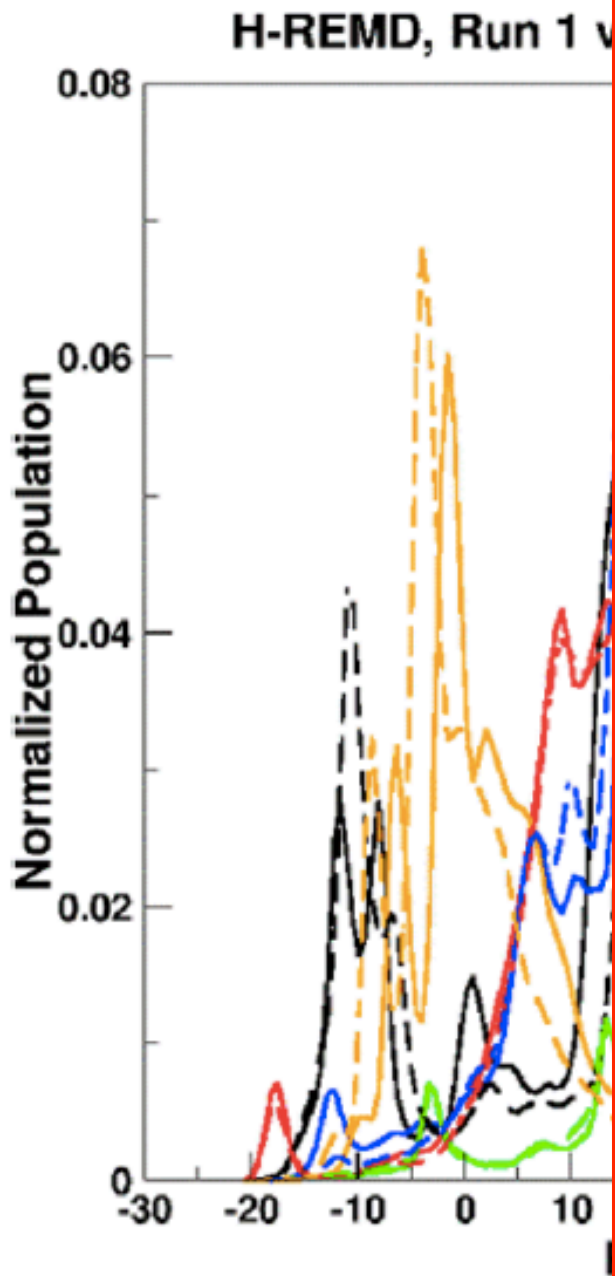
- pH
- restraints, umbrella potentials, ...
- force field / parameter sets
- biasing potentials (aMD)

Fukunishi, H., Wanatabe, O., and Takada, S., J. Chem. Phys. 2002.
Sugita, Y., Kitao, A., and Y. Okamoto, J. Chem. Phys. 2000.

H-REMD, Run 1 vs. Run 2

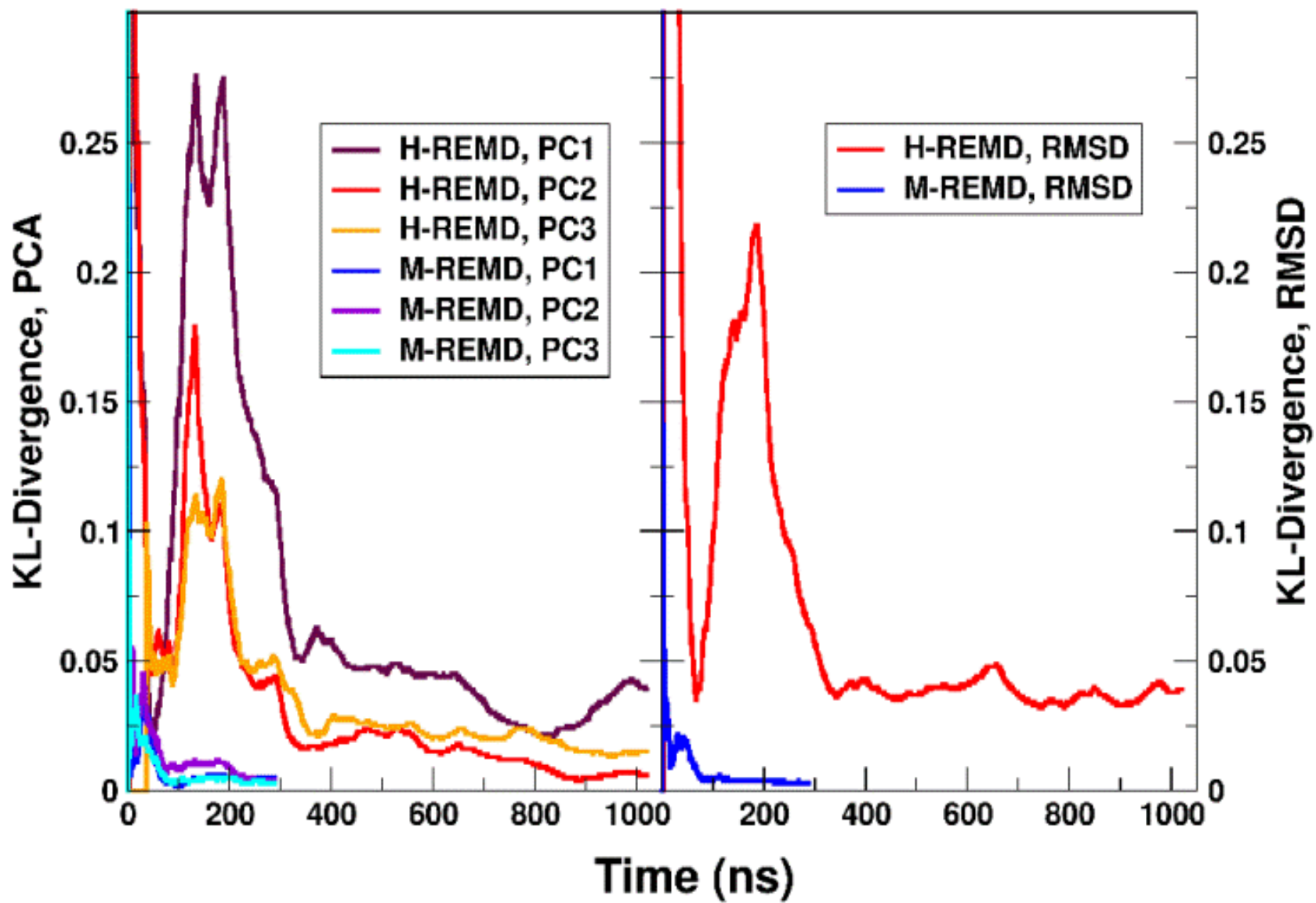
M-REMD, Run 1 vs. Run 2





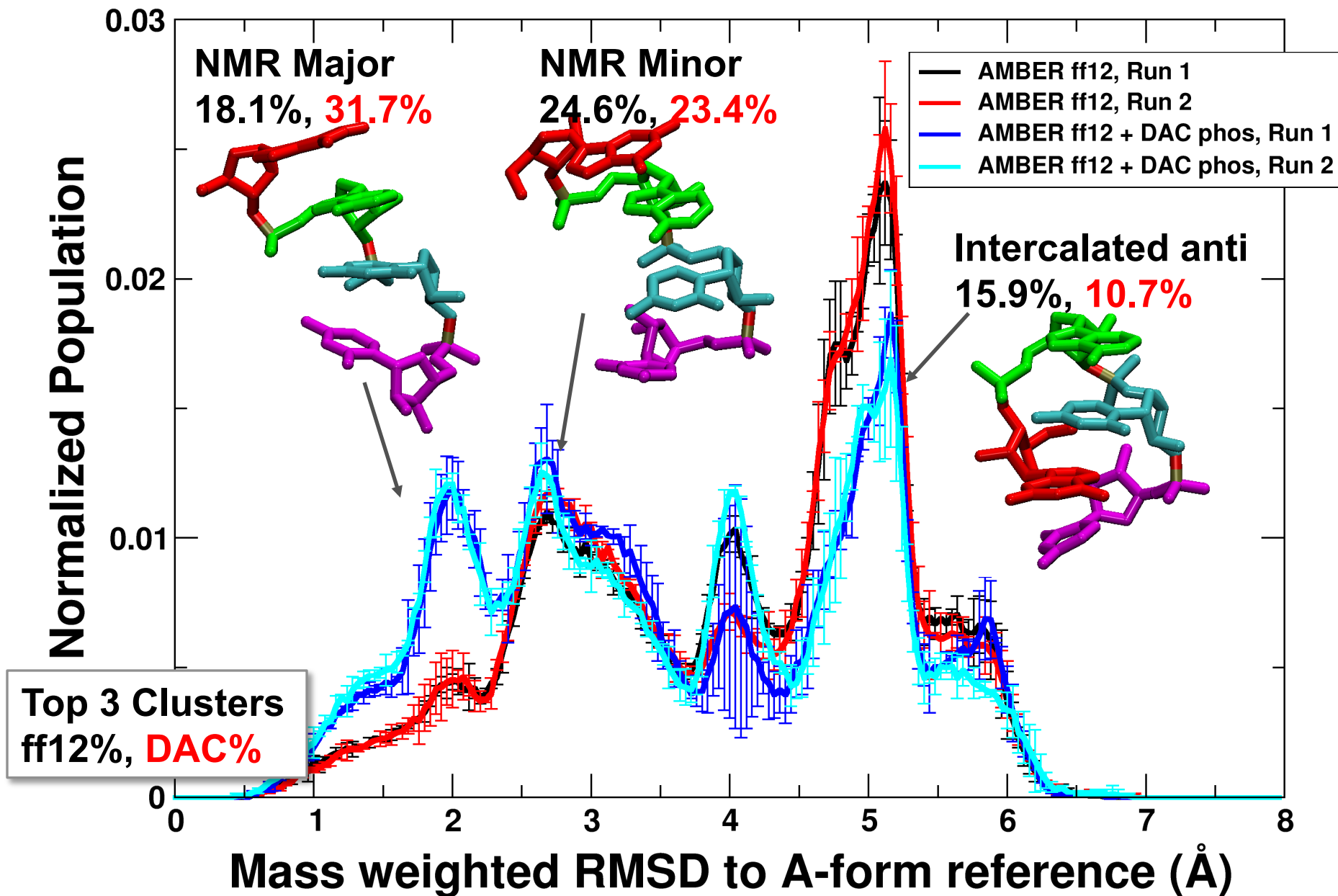
CPPTRAJ in AmberTools

```
# Read in both trajectories
#
trajin traj.run1.nc
trajin traj.run2.nc
# RMS-fit to first frame
#
rms first :1-4&!@H=
# Create an average structure
#
average gaccAvg.rst7 ncrestart
# Save coordinates as 'crd1'
#
createcrd crd1
run
# Fit to average structure
#
reference gaccAvg.rst7.1 [avg]
# RMS-fit to average structure
#
crdaction crd1 rms ref [avg] :1-4&!@H=
# Calculate coordinate covariance matrix
#
crdaction crd1 matrix covar :1-4&!@H= name gaccCovar
# Diagonalize coordinate covariance matrix, first 15 E.vecs
#
runanalysis diagmatrix gaccCovar out evecs.dat vecs 15
# Now create separate projections for each trajectory
#
crdaction crd1 projection P1 modes evecs.dat \
    beg 1 end 15 :1-4&!@H= crdframes 1,$STOP1
crdaction crd1 projection P2 modes evecs.dat \
    beg 1 end 15 :1-4&!@H= crdframes $START2,last
# Now histogram first 5 projections for each
#
hist P1:1,* ,*,*,100 out pca.hist.agr norm name P1-1
hist P1:2,* ,*,*,100 out pca.hist.agr norm name P1-2
```

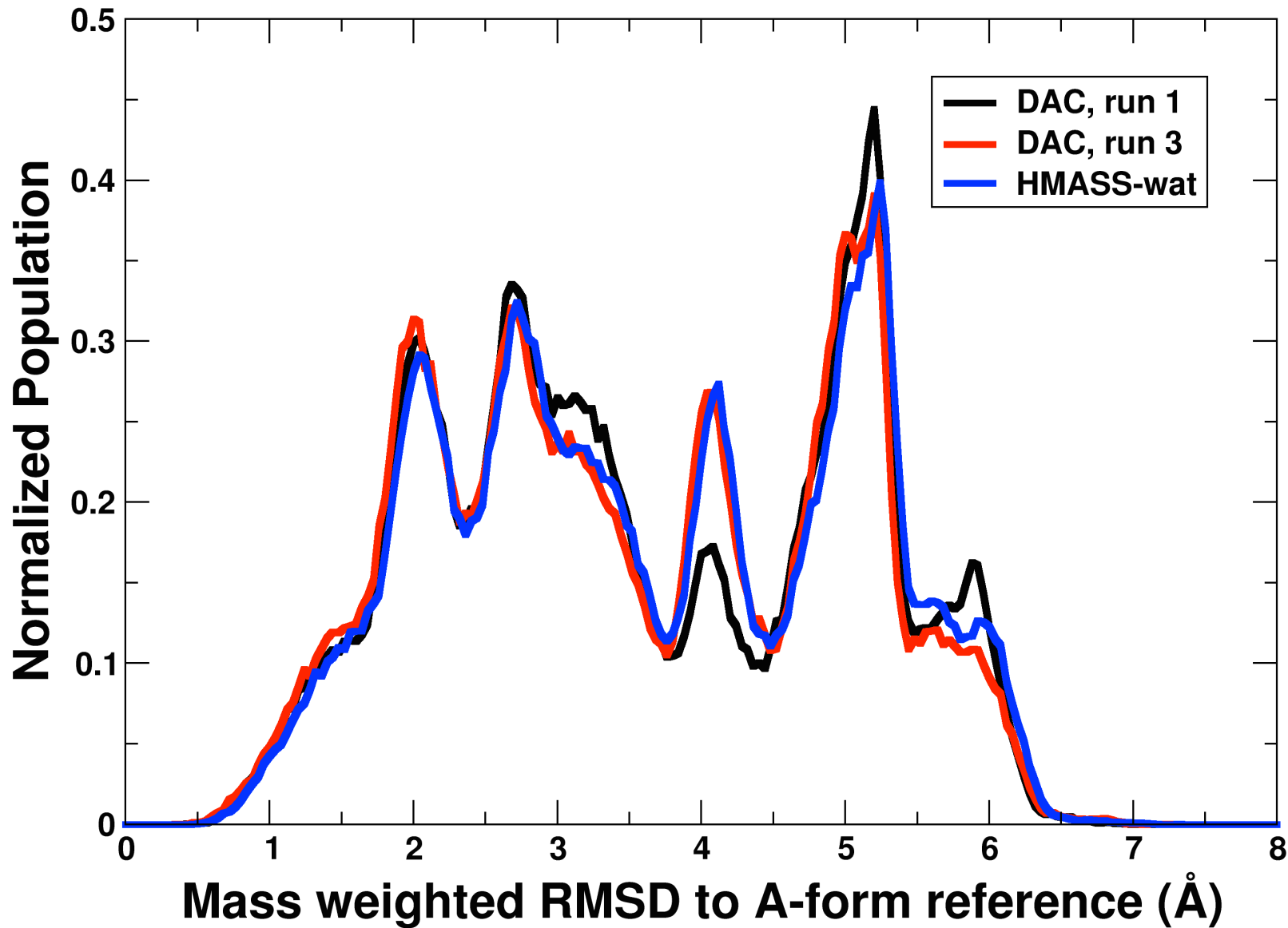
Convergence Analysis, GACC Ensemble

1st half vs. 2nd half, Force Field Comparison

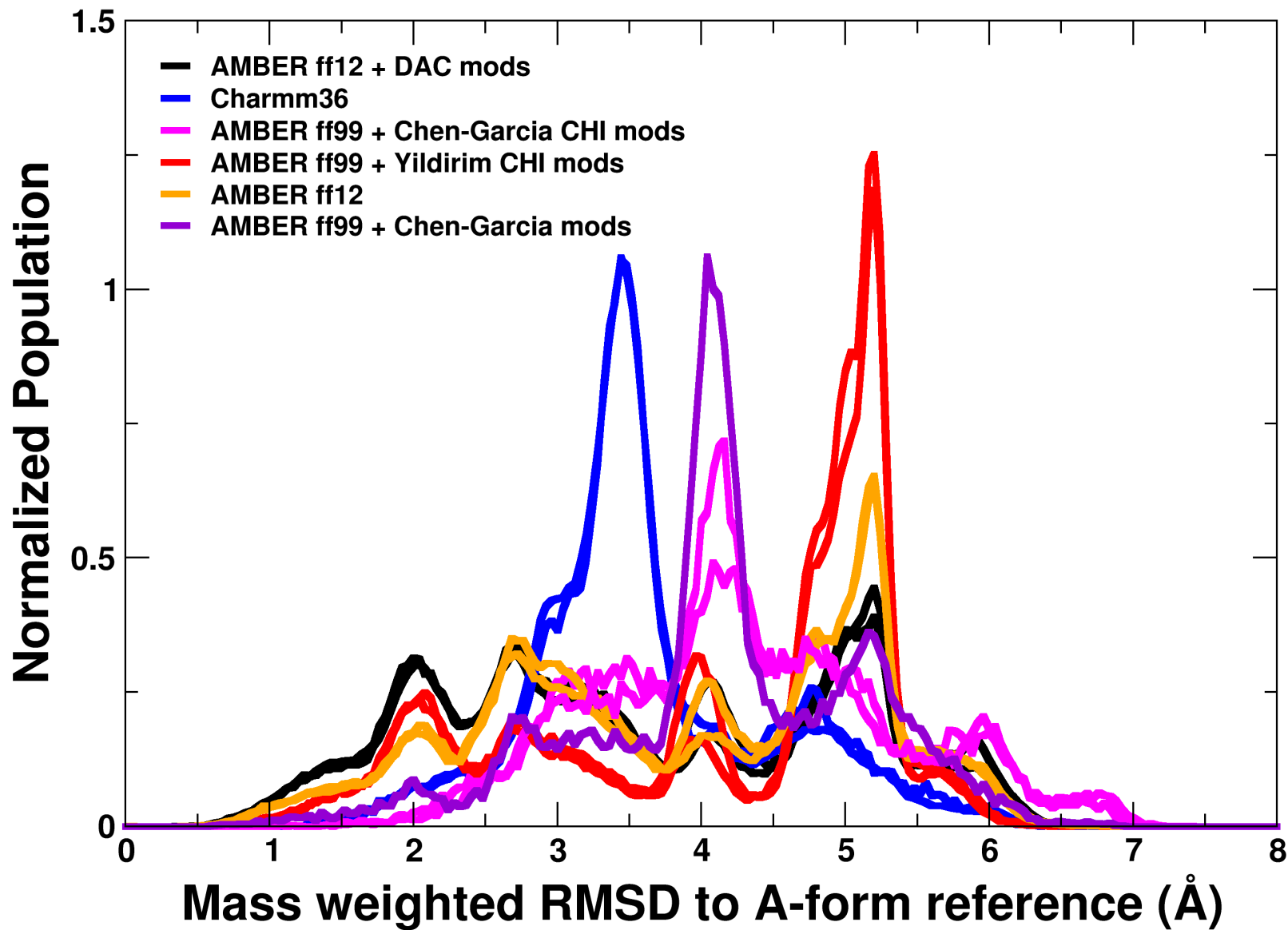


GACC Ensemble, using H-mass Repartitioning

277K replicas

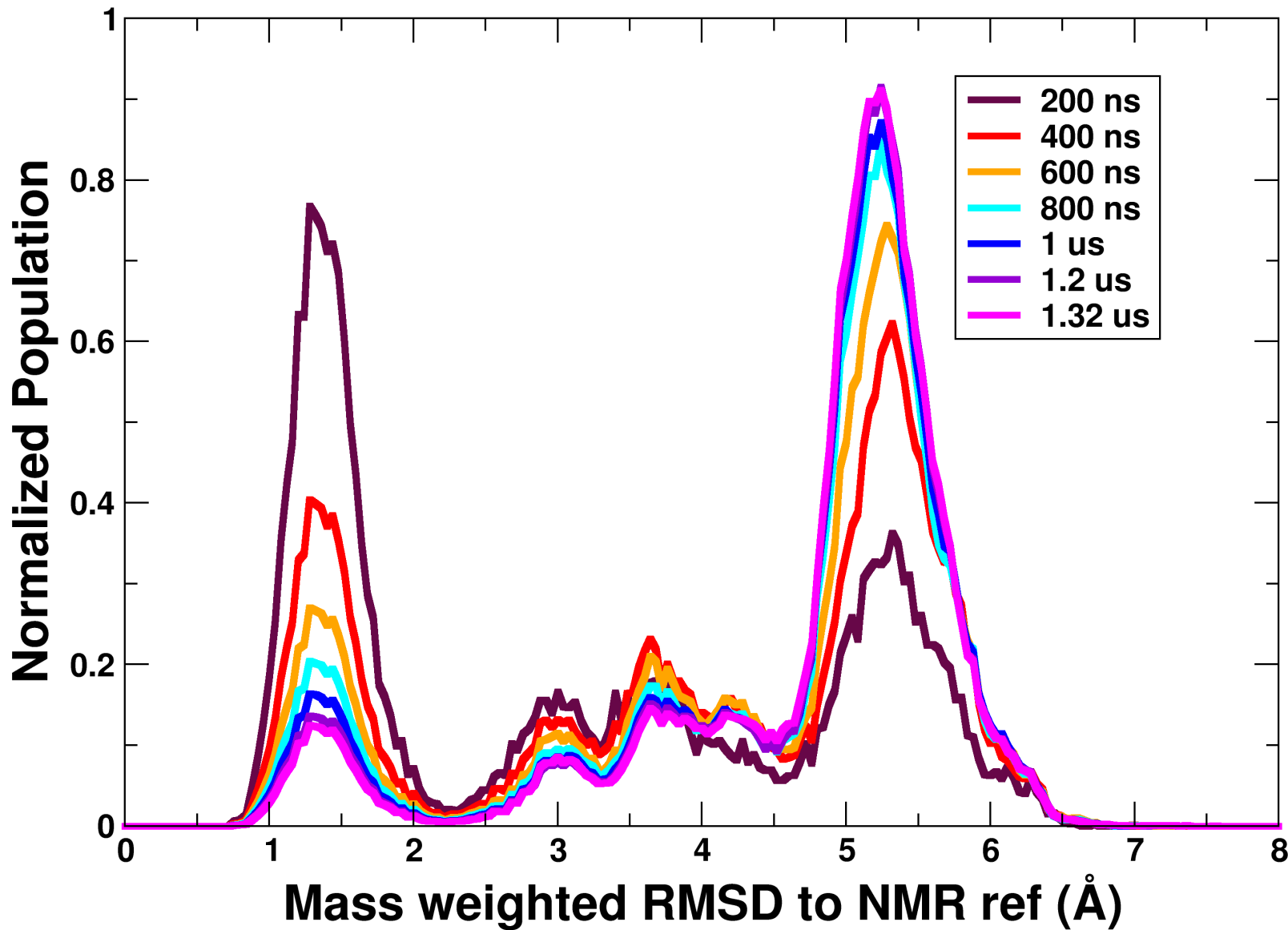


GACC Ensemble, Force Field Comparison



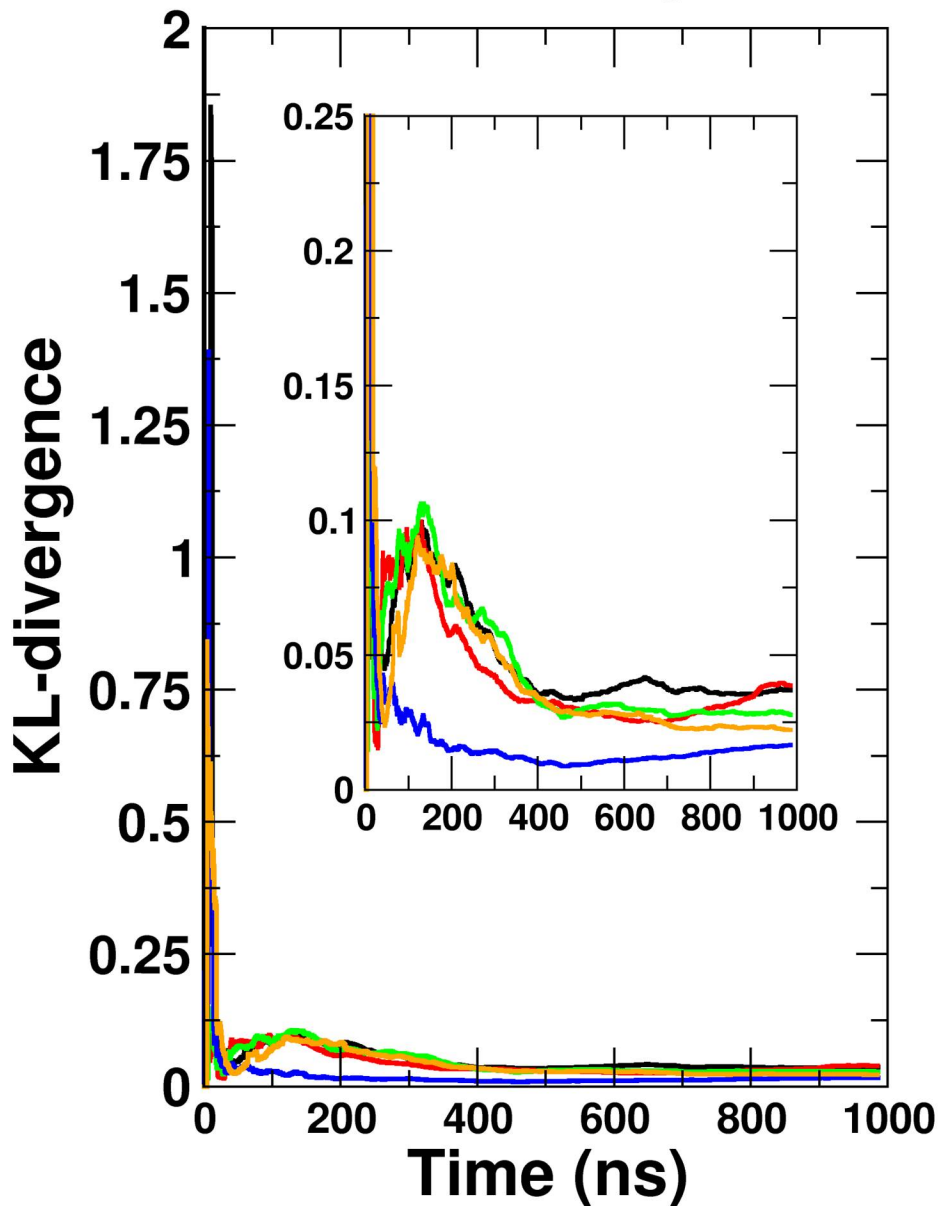
UUCG M-REMD Populations - Convergence Analysis

277K Replica, Truncated - Restrained (low)



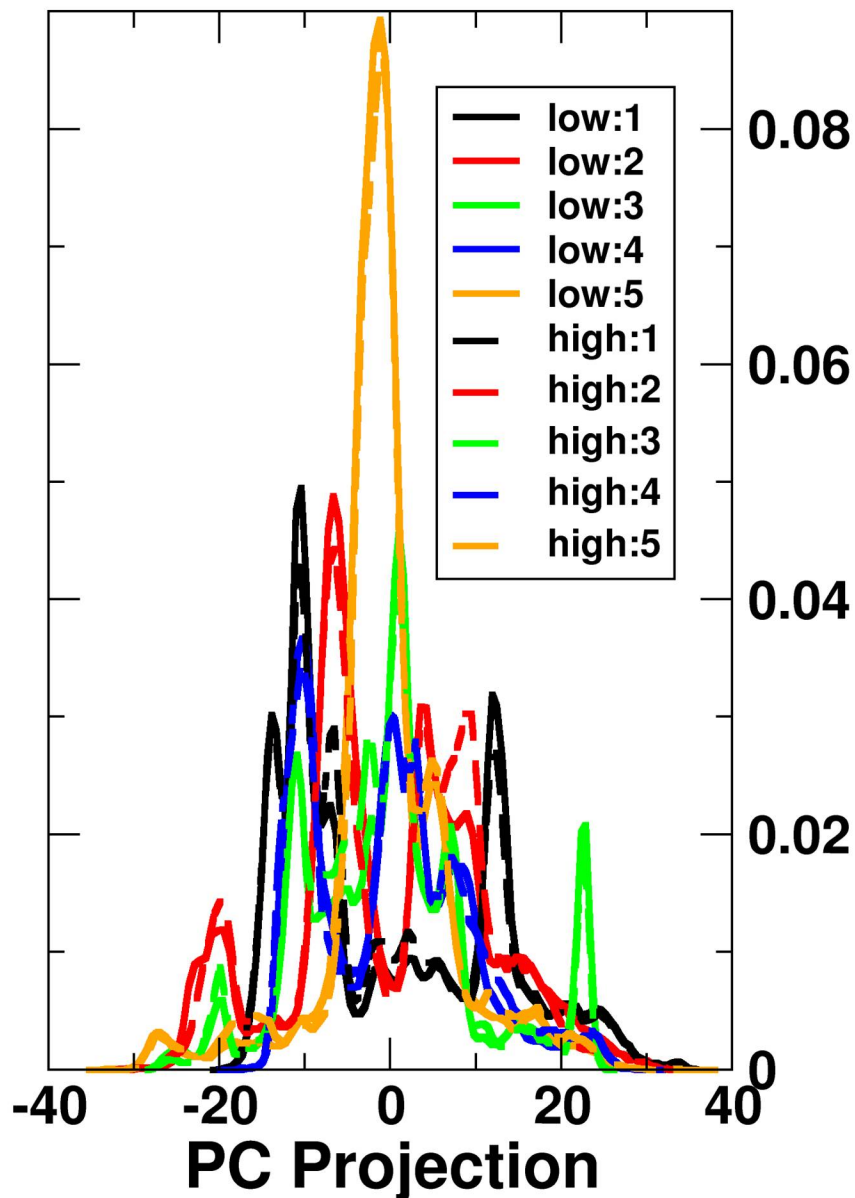
KL Divergence of PCA

Rest-low vs. Rest-high, 277K



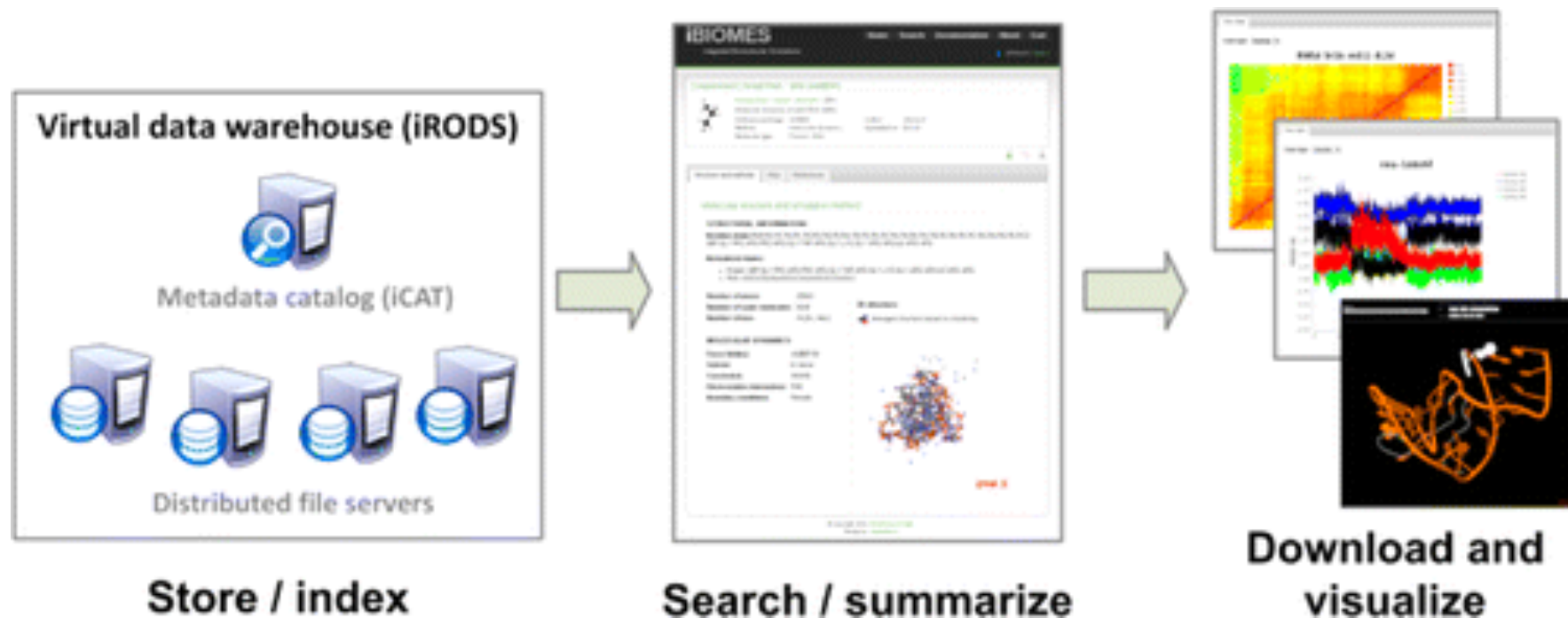
PCA Histogram Analysis

Rest-low vs. Rest-high, 277K



iBIOMES: Managing and Sharing Biomolecular Simulation Data in a Distributed Environment

Julien C. Thibault,[†] Julio C. Facelli,^{†,‡} and Thomas E. Cheatham, III^{*,§}



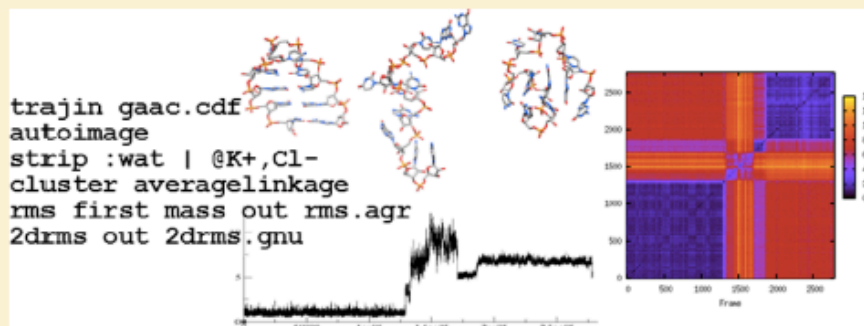
PTRAJ and CPPTRAJ: Software for Processing and Analysis of Molecular Dynamics Trajectory Data

Daniel R. Roe* and Thomas E. Cheatham, III*

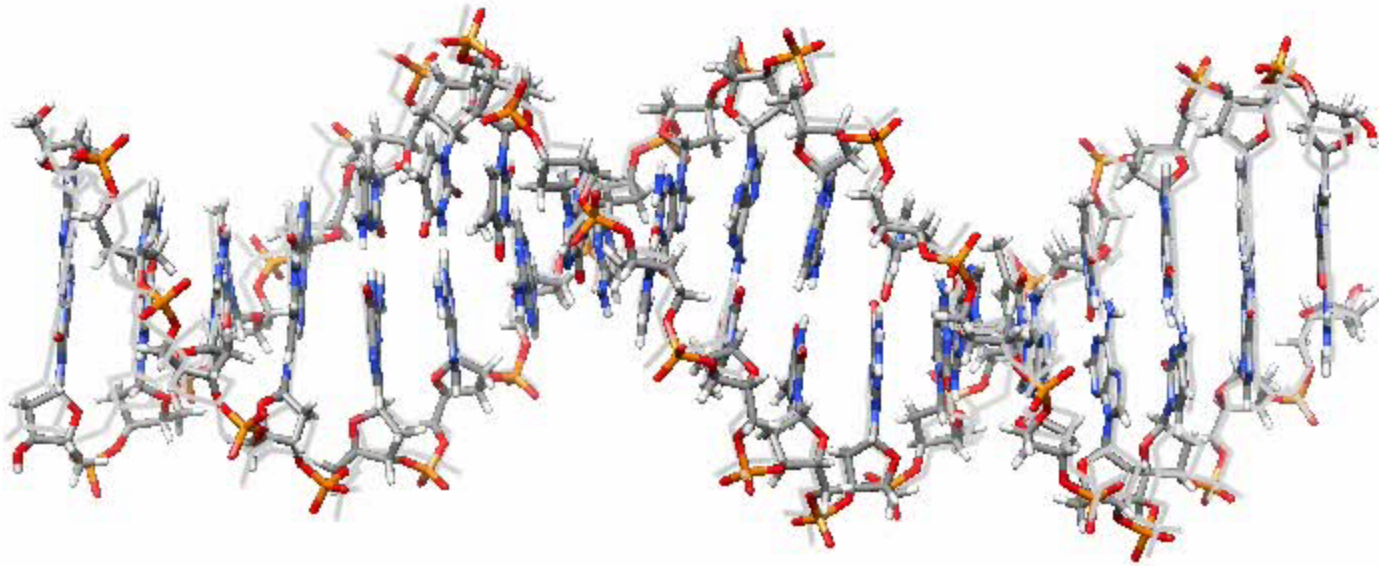
Department of Medicinal Chemistry, College of Pharmacy, 2000 South 30 East Room 105, University of Utah, Salt Lake City, Utah 84112, United States

Supporting Information

ABSTRACT: We describe PTRAJ and its successor CPPTRAJ, two complementary, portable, and freely available computer programs for the analysis and processing of time series of three-dimensional atomic positions (i.e., coordinate trajectories) and the data therein derived. Common tools include the ability to manipulate the data to convert among trajectory formats, process groups of trajectories generated with ensemble methods (e.g., replica exchange molecular dynamics), image with periodic boundary conditions, create average structures, strip subsets of the system, and perform calculations such as RMS fitting, measuring distances, B-factors, radii of gyration, radial distribution functions, and time correlations, among other actions and analyses. Both the PTRAJ and CPPTRAJ programs and source code are freely available under the GNU General Public License version 3 and are currently distributed within the AmberTools 12 suite of support programs that make up part of the Amber package of computer programs (see <http://ambermd.org>). This overview describes the general design, features, and history of these two programs, as well as algorithmic improvements and new features available in CPPTRAJ.



questions?



2 ns intervals, 10 ns running average, every 5th frame (~10 us).