Advanced path sampling of rare events in biomolecular systems

Peter Bolhuis van 't Hoff institute for Molecular Sciences University of Amsterdam, The Netherlands



Outline

- Simulation of biomolecular systems
- Basic TPS
 - shooting algorithms
 - stable state definitions
 - example Photoactive Yellow Protein
 - reaction coordinate analysis
- Advanced path sampling
 - rates by Transition Interface Sampling (TIS)
 - replica exchange TIS and multiple state TIS
 - single replica multiple state TIS
 - example Trp-cage folding network
- Conclusion

Protein self-assembly

understanding the cellular processes

- folding
- structure formation (cytoskeleton)
- complex formation (regulation)
- neurodegenerative/genetic diseases

—

novel self assembling biomaterials

- artificial tissue
- smart packaging
- self healing coatings
- sensors

Challenge

understanding and predicting protein assembly with advanced molecular simulation



All-atom force fields for biomolecules

• Potential energy for protein

$$V(\mathbf{r}) = \sum_{bonds} k_r (r - r_{eq})^2 + \sum_{angles} k_{\theta} (\theta - \theta_{eq})^2 + \sum_{dihedrals} \frac{1}{2} v_n (1 + \cos(n\phi - \phi_0))$$
$$+ \sum_{i < j} \left(\frac{a_{ij}}{r_{ij}^{12}} - \frac{b_{ij}}{r_{ij}^{6}} + \frac{q_i q_j}{\varepsilon r_{ij}} \right)$$

vdW interactions only between non-bonded |i-j|>4

Currently available empirical force fields

- CHARMm (MacKerrel et 96)
- AMBER (Cornell et al. 95)
- GROMOS (Berendsen et al 87)
- OPLS-AA (Jorgensen et al 95)
- ENCAD (Levitt et al 83)

- Subtle differences in improper torsions, scale factors 1-4 bonds, united atom rep.
- Partial charges based on empirical fits to small molecular systems
- Amber & Charmm also include ab-initio calculations
- Not clear which FF is best : top 4 mostly used
- Water models also included in description
 - TIP3P, TIP4P
 - SPC/E

.

• Current limit: 10⁶ atoms, microseconds (with Anton ms)



Molecular Dynamics of proteins

$$\begin{split} V(\mathbf{r}) &= \sum_{bonds} k_r (r - r_{eq})^2 + \sum_{angles} k_{\theta} (\theta - \theta_{eq})^2 + \sum_{dihedrals} \frac{1}{2} v_n (1 + \cos(n\phi - \phi_0)) \\ &+ \sum_{i < j} \left(\frac{a_{ij}}{r_{ij}^{12}} - \frac{b_{ij}}{r_{ij}^6} + \frac{q_i q_j}{\varepsilon r_{ij}} \right) \\ & \mathbf{F} = -\frac{\partial \mathbf{V}}{\partial \mathbf{r}} \end{split}$$



Free energy

MD yields

- **equilibrium statistics:** free energy landscapes, stable structures, transition states, ...
- **kinetics:** rates, mechanisms, transport properties, ...

FE landscape: high dimensional trajectory can be projected onto collective variable λ

$$e^{-\beta F(\lambda)} = \langle \delta(\lambda(\mathbf{r}) - \lambda) \rangle$$

Conformational order parameter λ

Timescales in proteins



Conformational space

Reactions Local flexibility Collective motions **B**ond Larger domain Methyl Loop vibration rotation motion motions Side-chain rotamer fs ps μs ms S ns

- Straightforward MD inefficient
- enhanced sampling: thermodynamic integration, umbrella sampling, hyper dynamics, adaptive biasing force, metadynamics, makes exponential barrier problem linear requires good reaction coordinate

Transition path sampling

Importance sampling of the rare event path ensemble: all dynamical trajectories that lead over (high) barrier and connect stable states.



Shooting moves





 $P_{\text{acc}}[x^{(0)}(\mathcal{T}) \to x^{(n)}(\mathcal{T})] = h_A[x_0^{(n)}]h_B[x_T^{(n)}] \qquad h_A(t) = \begin{cases} 1 & \text{if } x_t \in A \\ 0 & \text{if } x_t \notin A \end{cases}$



Flexible one way shooting



- higher acceptance, better convergence for diffusive transitions and long pathways
- requires some stochastic dynamics, e.g. thermostat
- needs check for decorrelation of paths
- useful for diffusive (bio)systems

Path sampling indicators



Least changed path

Path length distribution

Spring shooting for asymmetric barriers



Protein association/dissociation

System

- beta-lac dimer (2AKQ)
- ~65000 atoms
- AMBER99SB
- T=300K, P=1 atm
- dissociation $\Delta G=30$ kJ/mol



TPS of extremely asymmetric barrier

- spring shooting algorithm
- spring const 5 kT, dmax = 200
- 50 ps between frames
- max path length 50 ns
- acceptance 30%



Definition of stable states













TPS of proton transfer

- 28244 atoms
- CPMD/QMMM
- BLYP functional
- Electronic mass 750 au
- QM region: pCA, Glu46, Tyr42, Thr50, Arg52
- Gromos96 force field
- TPS: two way shooting, perturbation temp 35 K
- 160 paths/ 50% acceptance
- average path length 0.5-1.5 ps
- reaction time microseconds



stable states	pR (reaction)	pB' (product)
pCA-Glu46(H)	> 1.60 A	< 0.98 A
OX2-Tyr42	> 3.70 A	< 1.80 A
OXI-Tyr42	> 5.30 A	< 1.80 A

Transition path sampling of partial unfolding



Table 1. Statistics of the TPS ensembles. The average path length is a weighted average over the whole ensemble. Decorrelated pathways have lost the memory of the previous decorrelated pathway. The aggregate time is the ensemble aggregate length

	$pB' - I_{\alpha}$	$U_{\alpha} - S_E$	$U_{\alpha} - S_X$	$S_E - pB$
acceptance	41%	25%	38%	44%
avg. path length	105 ps	1.8 ns	1.5 ns	1.7 ns
accepted paths	3847	305	584	311
decorr. paths	180	18	7	29
aggregate time (µs)	1.0	2.3	2.3	1.2

Vreede, Juraszek, PGB, PNAS 2010



Transition states by committor



W.E, E. Vanden-Eijnden, J. Stat.Phys, **123** 503 (2006)

Reaction coordinate analysis

- Each TPS shot is a committor attempt. Use this information to optimise model of reaction coordinate r
- The probability that structure *x* with rc *r* is on a transition path (for diffusive dynamics)

 $p(TP|r) = 2p_B(r)(1 - p_B(r))$

• Assume committor function to be

$$p_B(x) = \frac{1}{2} + \frac{1}{2} \tanh[r(q(x))]$$

• parametrize *r* as linear combination of q

$$r(\mathbf{x}) = \sum_{i} \alpha_{i} q(\mathbf{x}) + \alpha_{0}$$

• best *r* is maximizing likelihood

$$L(\alpha) = \prod_{i=1}^{N_B} p_B(r(q(\mathbf{x}_i^{(B)}))) \prod_{i=1}^{N_A} (1 - p_B(r(q(\mathbf{x}_i^{(B)}))))$$



Included order parameters

.

Number of waters around		Distance between center of mass of side chains	
pCA	riw _x	pCA - Tyr42	dOCK-Manager
Ty:42	CW.y	pCA - Glu46	dXE ^{rm}
Glu46	COW 2	pCA - Phe62	-dOCFY-com
Distance between atoms		pCA - Phe96	dXF2 ^{tom}
pCA-OH - Tyr42-OH	dXY	pCA – lie49	dXP ^{rom}
pCA-O4' - Glu46-CD	OXE	Lys64 - Thr70	dK7tim
pCA-01 - Cys69-N	dOaC	Distance between center of mass of groups of residues	
Glu46-CD - Thr50-OG1	der	13-17 - 114-116	dN-loop
Glu46-CD - 7#42-OH	OYE	35-37 - 98-101	g100pt1
Arg52-CZ - Asp97-CG	dRD	35-37 - 114-116	dioops2
Lys54-NZ - Thr70-OG1	dK7	med	200120
pCA-O1 - Asp97-N	diClaN	11-15	rmsdari
pCA-O4" - Ile49-N	000	19-23	rmsdag
pCA-O4' - Thr50-N	dXT	43-51	rmsd.
pCA-O4' - Arg52-N	dXth	62-68	rmaden
pCA-O4' - Asp97-N	dX/V1	75-86	rmsdcr
oCA-O4' - Aso97-CG	dX0V2	111-116	rmsdaup
Ala44-N - Pro54-CG	dPA	Dihedral angles in pCA	433583
Gly47-CA - Ara52-O	dGR	N-CA-CB-SG	dificace
Glu46-CD - Avo43-ND2	d'EN	CA-CB-SG-C1	din _{ceso}
Glu46-CD - Gly51-N	dEC	Other	
Ash43-0 - Gh47-H	dh61	Number of hydrogen bonds in #3	nhb
Ala44-O - Aup48-H	dhb2	Cosines of dihedral angles ϕ in $\alpha 3$	40-40
Ala45-O - 1e49-H	dhb3	Cosines of dihedral angles ψ in ω	WH2 - W10
Glu46-O - Thr50-H	dhb4		
Gly47-O - Gly51-H	dhb/5		
Asp2D-CG - Lys55-NZ	dDK		
Asp24-CG - Lys55-NZ	dDK2		
Glu9-CD - Lys110-MZ	dEK		
Glu12-CD - Lys190-NZ	d5K2		
K111-NZ - Glu116-CD	ake		

Reaction coordinate of $helix_{\alpha 3}$ unfolding



Reaction coordinate by likelihood maximization (Peters & Trout, JCP 2006)

Order Parameters involved (out of 78): $RMSD_{\alpha}$ nwY42: water molecules around Tyr42 dPA: distance Ala44(N) - Pro54(CY) dhb2: distance Ala44(O) - Asp48(H)

 δ Lmin = 4.17

n	In L	RC
I	-2117	3.89–29.10 × rmsdα
2	-2098	3.88–26.35 × rmsdα – 0.19 × nwY42
3	-2085	5.11–16.81 × rmsdα – 4.68 × dhb2 – 2.55 × dPA

Reaction coordinate $pB' \rightarrow I_{\alpha}$

 $r = 5.11 - 16.28 \text{ rmsd}_{\alpha 3} - 4.68 \text{ dhb}_{2-} 2.55 \text{ d}_{PA}$



Solvent exposure transitions





Outline

- Simulation of biomolecular systems
- Basic TPS
 - shooting algorithms
 - stable state definitions
 - example Photoactive Yellow Protein
 - reaction coordinate analysis
- Advanced path sampling
 - rates by Transition Interface Sampling (TIS)
 - replica exchange TIS and multiple state TIS
 - single replica multiple state TIS
 - example Trp-cage folding network
- Conclusion

Rough free energy landscapes



WW domain folding

J. Juraszek and PGB, Biophys. J.98, 646 (2010).

PYP signal transduction

J.Vreede J, J. Juraszek and PGB PNAS 107 2397(2010)

Markov state model

molecular dynamics trajectory

coarse grained trajectory





integrate equations of motion

time step $\Delta t \approx fs$

$$\frac{dp_i(t)}{dt} = \sum_{j \neq i} k_{ji} p_j(t) - \sum_{j \neq i} k_{ij} p_i(t)$$

master equation, solve analytically or by KMC time step set by rates

Transition interface sampling



Transition interface sampling

Introduce set of interfaces λi



TIS : for each interface i sample pathways that cross λi

 $P_A(\lambda|\lambda_i)$ = probability that path crossing λ_i for first time after leaving A reaches λ before A

Sample with shooting, replica exchange, reversal, first/last interface move T.S. van Erp, PRL **98**, 268301 (2007) P.G. Bolhuis, JCP **129**, 114108 (2008)

Pros and cons

- Advantages TPS/TIS
 - correct rate (recrossings are counted)
 - no reaction coordinate needed
 - access to the entire path space by reweighting
 - mechanistic insight through committor analysis
- disadvantage of TPS/TIS
 - stable states need to be carefully defined: core sets
 - not easy to implement
 - computationally expensive
- challenges and convergence issues of TPS/TIS
 - can get trapped in intermediate metastable states
 - multiple channels not easily sampled (addressed by RETIS)

Challenges for path sampling

Multiple channels

- multiple channels are not sampled properly with shooting
- Replica exchange TIS



T.S. van Erp, PRL **98**, 268301 (2007) PGB, JCP **129**, 114108 (2008)

Presence of intermediates

- paths become very long because of intermediates
- Multiple state TIS

J. Rogal, PGB, J. Chem. Phys. (2008).



Path replica exchange



T.S. van Erp, PRL **98**, 268301 (2007) P.G. Bolhuis, JCP **129**, 114108 (2008)

Include paths starting in B



$$P_{acc}(i \leftrightarrow j) = \min\left(1, \frac{g_{\lambda i}[\mathbf{x}^{(j)}(L^{(j)})]g_{\lambda j}[\mathbf{x}^{(i)}(L^{(i)})]}{g_{\lambda i}[\mathbf{x}^{(i)}(L^{(i)})]g_{\lambda j}[\mathbf{x}^{(j)}(L^{(j)})]}\right)$$

 $g_{\lambda}[\mathbf{x}(L)] = \begin{cases} 1 & \text{if path crosses } \lambda \\ 0 & \text{otherwise} \end{cases}$

Samples AA, AB, BA and BB paths -shooting move -time reversal move

-exchanges

Replica Exchange Transition Interface Sampling



Shooting Move



T.S. van Erp, PRL **98**, 268301 (2007) P.G. Bolhuis, JCP **129**, 114108 (2008) First Interface Move



Exchange Move







Center of Mass Autocorrelation



Necessity of exchange

• illustrated on two channels with different barrier height.



Challenges for path sampling

Multiple channels

- multiple channels are not sampled properly with shooting
- Replica exchange TIS



T.S. van Erp, PRL **98**, 268301 (2007) PGB, JCP **129**, 114108 (2008)

Presence of intermediates

- paths become very long because of intermediates
- Multiple state TIS

J. Rogal, PGB, J. Chem. Phys. (2008).



Multiple state transition interface sampling



 $P_A(\lambda_{(s+1)A}|\lambda_{(s+1)A})$ = probability path crossing s for first time after leaving A reaches s+1 before A



rates can be used in Markov state model

Single replica MSTIS



Problem: interfaces close to stable states will be favored Solution: bias with e.g. Wang Landau scheme

Single replica Wang-Landau path sampling

- a single replica walks along set of interfaces
- each interface of i for state **J** has a density of paths g_i and histogram h_i
- each time a interface is sampled $g_i = g_i^* f$ and $h_i = h_i + I$
- a swap between interfaces is accepted with:

$$P_{acc}^{wl}[\lambda_i \leftrightarrow \lambda_{i+1}] = \min\left[1, \frac{g_i}{g_{i+1}}\right]$$

- when path switches from $J \rightarrow J$ to $J \rightarrow K$, allow switch to new set of interfaces for K
- if histogram is "flat" then
 - reset histogram $h_i = 0$
 - reduce factor f = \sqrt{f} , continue
- weights g_i reflect ratio of pathways on each interface = crossing probability

advantage: only one replica needed disadvantage: need to wait until histogram is flat only correct in limit of $f \rightarrow 0$

F.Wang and D.P.Landau, PRL 86, 2050 (2001) W. Du and PGB, JPC, 139, 044105, (2013)

Langevin dynamics on 2D potential





- β =8, γ =2.5, Δt =0.05
- number of interfaces 20
- metric: distance to center of state
- shooting, time reversals and swaps
- reweighting gives FE and committor



Improve on Wang-Landau

- WL converges slowly
- improve by imposing fixed bias
- DOP turns out to converge to crossing probability $P_1(\lambda)$
 - because $P_i(\lambda_i)/P_i(\lambda_j)$ is probability to reach λ_j from λ_i



• perfect bias for flat histogram is crossing probability $P_I(\lambda)$

Convergence fixed bias vs WL

use equilibrium population to assess convergence

$$\mathbf{p}_{eq} = \lim_{t \to \infty} \mathbf{p}^T(t) = \mathbf{p}^T(0) \lim_{t \to \infty} \exp(\mathbf{K}t)$$



Trp-cage folding



Kinetics from rate matrix

		, PN	SN	Mg	meta	Pd	LN	LSN	Lm	Lo	Ι	W	other state	U
Ν		3.75×10^{-3}	2.33×10^{-4}	4.67×10^{-4}	1.65×10^{-2}	5.35×10^{-3}	2.43×10^{-3}		1.04×10^{-4}		1.00×10^{-5}	2.12×10^{-7}	9.08×10^{-5}	2.35×10^{-5}
PN	6.68×10^{-1}		6.73×10^{-4}	3.66×10^{-4}	8.61×10^{-3}	3.48×10^{-3}	2.21×10^{-3}		7.16×10^{-5}		2.02×10^{-4}		1.70×10^{-3}	4.92×10^{-5}
SN	1.18×10^{-3}	1.91×10^{-5}		4.48×10^{-6}	2.88×10^{-4}	8.16×10^{-4}	2.85×10^{-5}	8.81×10^{-4}		2.55×10^{-5}	1.10×10^{-4}	2.58×10^{-8}	1.05×10^{-3}	2.26×10^{-4}
Mg	4.47×10^{-1}	1.97×10^{-3}	8.50×10^{-4}		3.45×10^{-1}		8.25×10^{-2}		3.57×10^{-5}			2.37×10^{-6}	1.49×10^{-3}	
meta	7.65×10^{-1}	2.24×10^{-3}	2.64×10^{-3}	1.67×10^{-2}		3.68×10^{-3}	7.85×10^{-3}	2.19×10^{-5}	3.42×10^{-4}		1.50×10^{-4}	8.59×10^{-7}	1.07×10^{-3}	9.01×10^{-5}
\mathbf{Pd}	4.87×10^{-1}	1.78×10^{-3}	1.47×10^{-2}		7.22×10^{-3}		8.42×10^{-5}	1.01×10^{-4}		1.61×10^{-4}	1.46×10^{-4}	2.56×10^{-6}	4.79×10^{-3}	8.32×10^{-5}
LN	1.01×10^{-1}	5.16×10^{-4}	2.35×10^{-4}	3.59×10^{-3}	7.06×10^{-3}	3.85×10^{-5}		6.35×10^{-4}	2.16×10^{-3}		6.42×10^{-5}	7.31×10^{-6}		5.52×10^{-4}
LSN			3.23×10^{-2}		8.77×10^{-5}	2.06×10^{-4}	2.83×10^{-3}			3.68×10^{-3}	9.89×10^{-5}	3.96×10^{-7}	1.41×10^{-3}	1.08×10^{-3}
Lm	6.05×10^{-2}	2.34×10^{-4}		2.17×10^{-5}	4.29×10^{-3}		3.02×10^{-2}					2.71×10^{-6}		
Lo			2.27×10^{-3}			7.98×10^{-4}		8.95×10^{-3}			4.04×10^{-4}	1.74×10^{-6}	5.14×10^{-2}	8.69×10^{-3}
Ι	1.27×10^{-2}	1.44×10^{-3}	2.76×10^{-2}		4.10×10^{-3}	2.04×10^{-3}	1.95×10^{-3}	6.74×10^{-4}		1.13×10^{-3}		3.77×10^{-6}	1.25×10^{-2}	6.50×10^{-3}
W	1.00×10^{-2}		2.42×10^{-4}	1.17×10^{-4}	8.77×10^{-4}	1.33×10^{-3}	8.30×10^{-3}	1.01×10^{-4}	2.21×10^{-4}	1.83×10^{-4}	1.41×10^{-4}		1.05×10^{-5}	1.97×10^{-1}
other	9.16×10^{-3}	9.63×10^{-4}	2.10×10^{-2}	1.57×10^{-4}	2.34×10^{-3}	5.31×10^{-3}		7.65×10^{-4}		1.15×10^{-2}	1.00×10^{-3}	2.24×10^{-8}		2.94×10^{-3}
U	8.42×10^{-5}	9.92×10^{-7}	1.60×10^{-4}		6.98×10^{-6}	3.28×10^{-6}	4.75×10^{-5}	2.09×10^{-5}		6.91×10^{-5}	1.84×10^{-5}	1.50×10^{-5}	1.04×10^{-4}	





Transition path theory analysis

Analysis of large kinetic matrix for more insight

E, Vanden-Eijnden, J. Stat. Phys 2006 Noe et al., PNAS 2009

compute forward committor q_i for each state i

$$q_i^+ = \sum_{k \in \mathbf{I}} T_{ik} q_k^+ + T_{iU} \qquad \mathbf{T} = \exp(\mathbf{K}\tau)$$

Compute flux from i to j

$$f_{ij} = \pi_i q_i^- T_{ij} q_j^+$$

compute effective flux

$$f_{ij}^{+} = \max\{0, f_{ij} - f_{ji}\}$$
$$k_{NU} = \frac{\sum_{i \neq N} \pi_N T_{Ni} q_i^{+}}{\tau \sum_{i \in \mathbf{M}} \pi_i q_i^{-}}$$

overall rate constant

$$k_{NU}$$
 is 1.01 × 10⁻⁴ns⁻¹ \approx (9.9µs)⁻¹,
 k_{UN} = 4.17×10⁻⁴ns⁻¹ \approx (2.4µs)⁻¹

$$k_{NU} \exp = (12\mu s)^{-1}$$

 $k_{UN} \exp = \approx (4.1\mu s)^{-1}$

TPT flux analysis



158 μ s MD, around 70000 trajectories, represents around 15 ms of time rate matrix and flux are non sparse, many pathways possible

The Reweighted Path Ensemble



Projection of Reweighted Path Ensemble

3

-10

RPE can be used to project the conditional path dependent population density

$$\rho_{ij}(\mathbf{q}) = \langle h_i(\mathbf{x}_0) h_j(\mathbf{x}_L) \delta(\mathbf{q}(\mathbf{x}_k) - \mathbf{q}) \rangle_{RPE}$$
$$\rho(\mathbf{q}) = \rho_{AA}(\mathbf{q}) + \rho_{AB}(\mathbf{q}) + \rho_{BA}(\mathbf{q}) + \rho_{BB}(\mathbf{q}) = \langle \delta(\mathbf{q}(\mathbf{x}_k) - \mathbf{q}) \rangle_{RPE}$$

and thus the free energy landscape

$$F(\mathbf{q}) = -k_B T \ln\left(\rho(\mathbf{q})\right) + const,$$

and the (averaged) committor

$$p_A(\mathbf{q}) = \frac{\rho_{AA}(\mathbf{q}) + \rho_{BA}(\mathbf{q})}{\rho(\mathbf{q})}$$

$$p_B(\mathbf{q}) = \frac{\rho_{AB}(\mathbf{q}) + \rho_{BB}(\mathbf{q})}{\rho(\mathbf{q})}.$$

and the path density

$$n_{ij}(\mathbf{q}) = \langle h_i(\mathbf{x}_0) h_j(\mathbf{x}_L) h_{\mathbf{q}}(\mathbf{x}^L) \rangle_{RPE}$$

$$h_{\mathbf{q}}(\mathbf{x}^{L}) = \begin{cases} 1 & \text{if path visits } \mathbf{q} \\ 0 & \text{otherwise} \end{cases}$$





W Lechner, PGB, J. Stat. Phys. 145 841 (2011).



RPE Free energy for Trp-cage



- The I state looks like a state, but it is not stable
- Free energy projection can be misleading

Summary



- Transition path sampling gives unbiased paths of protein conformational changes
 - millisecond light-induced unfolding mechanism in PYP
 - association/dissociation of protein dimers
 - reaction coordinate requires advanced data analysis of simulations



- Single replica MSTIS samples equilibrium network of Trp Cage
 - asynchronous sampling scheme (in contrast to parallel replica exchange)
 - corroborates experimental evidence for near native intermediate
 - structure and correct time scales predicted by single replica MSTIS

• Path sampling gives quantitative insight in rough free energy landscapes, kinetics and mechanism simultaneously