

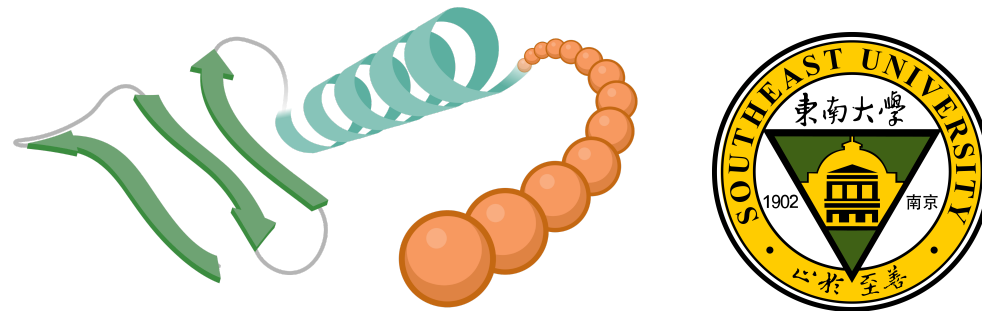
Final presentation

Research on the mechanism of protein folding dynamics

Reporter: **Zhenyu Wei**

Mentor: **Yunfei Chen**

School of mechanical engineering
Southeast University



Content

- **Research background**
- **Multi-scale modeling**
- **Enhanced sampling**
- **Free energy of ion interaction in solution**
- **Free energy matrix**
- **OpenPD**
- **Future Proposal**

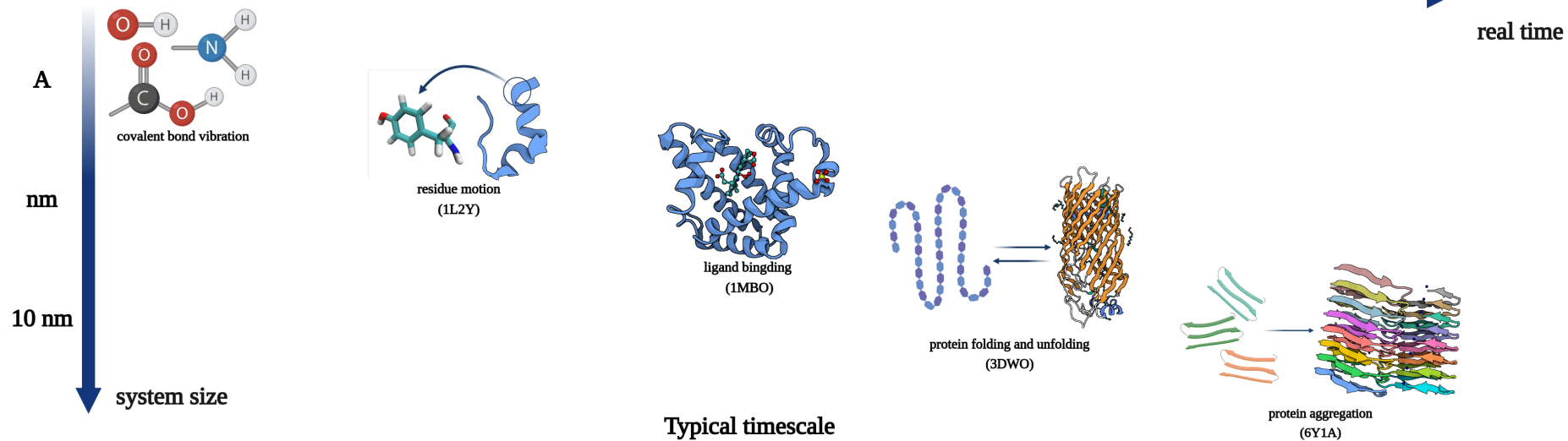
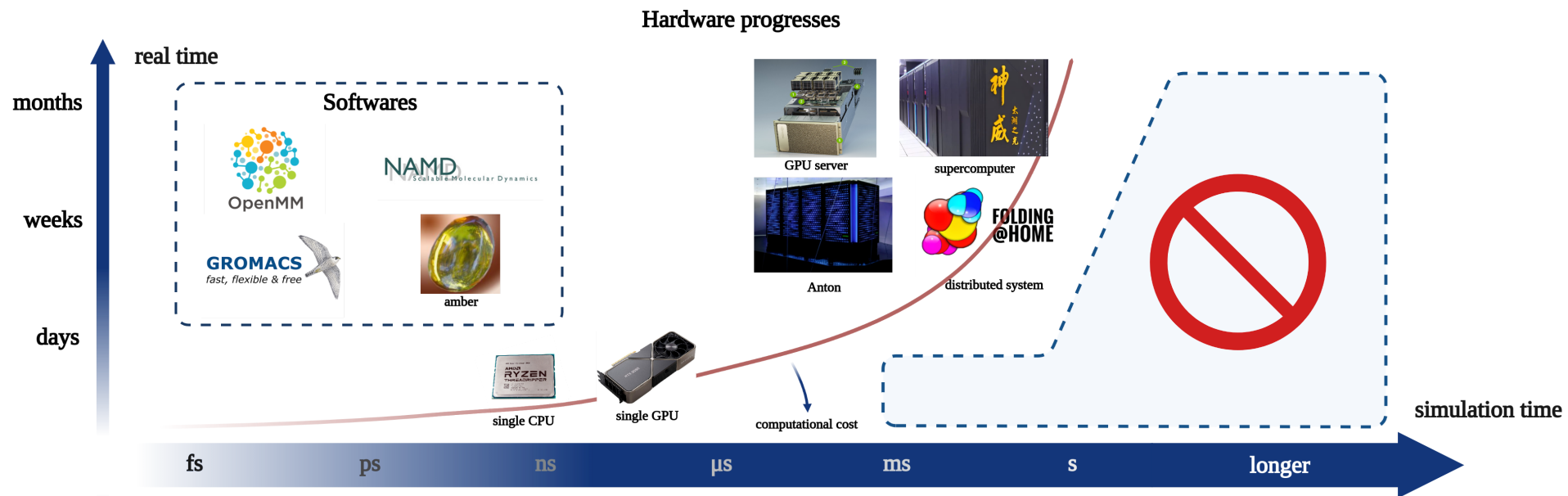
Research background

Importance of protein folding and misfolding

- Protein is essential in a wide range of fields.
 - **Membrane protein** play a vital role in the inter-cellular substance exchange and communication
 - **Protein enzyme** is the most efficient catalyst, maintaining gazillion of chemical reaction in live system
- Protein misfolding leads to many diseases:
 - Creutzfeldt-Jakob disease (CJD)
 - Type 2 diabetes
 - Many neurodegenerative diseases

- The **sequence** of peptide chain determine the final structure of protein **completely**
- The protein can fold to the **native structure** rapidly and precisely
- Levinthal's paradox¹:

The protein with 100 peptides may misfold into a maximum of 3^{198} different conformations. Therefore, it would require a time longer than the age of the universe to arrive at its **correct native conformation** if protein random searched all the possible conformation.



Multi-scale modeling

Fundamentals

- When we sample in NVT ensemble, the probability density follows the **Boltzmann distribution**:

$$p(\mathbf{q}_a, \mathbf{q}_{a-}) = \frac{\exp\left[-\frac{U(\mathbf{q}_a, \mathbf{q}_{a-})}{k_B T}\right]}{\int_{\Omega} \prod_{i=1}^{3N_a} dq_a^{(i)} \prod_{j=1}^{3N_{a-}} dq_{a-}^{(j)} \exp\left[-\frac{U(\mathbf{q}_a, \mathbf{q}_{a-})}{k_B T}\right]}$$

- We can write **experimental observations** in the form of ensemble average:

$$\langle A \rangle = \int_{\Omega} \prod_{i=1}^{3N_a} dq_a^{(i)} \prod_{j=1}^{3N_{a-}} dq_{a-}^{(j)} A(\mathbf{q}_a) p(\mathbf{q}_a, \mathbf{q}_{a-}) = \int_{\Omega_a} \prod_{i=1}^{3N_a} dq_a^{(i)} A(\mathbf{q}_a) \int_{\Omega_{a-}} \prod_{j=1}^{3N_{a-}} dq_{a-}^{(j)} p(\mathbf{q}_a, \mathbf{q}_{a-})$$

- Now we define a **marginal** probability density:

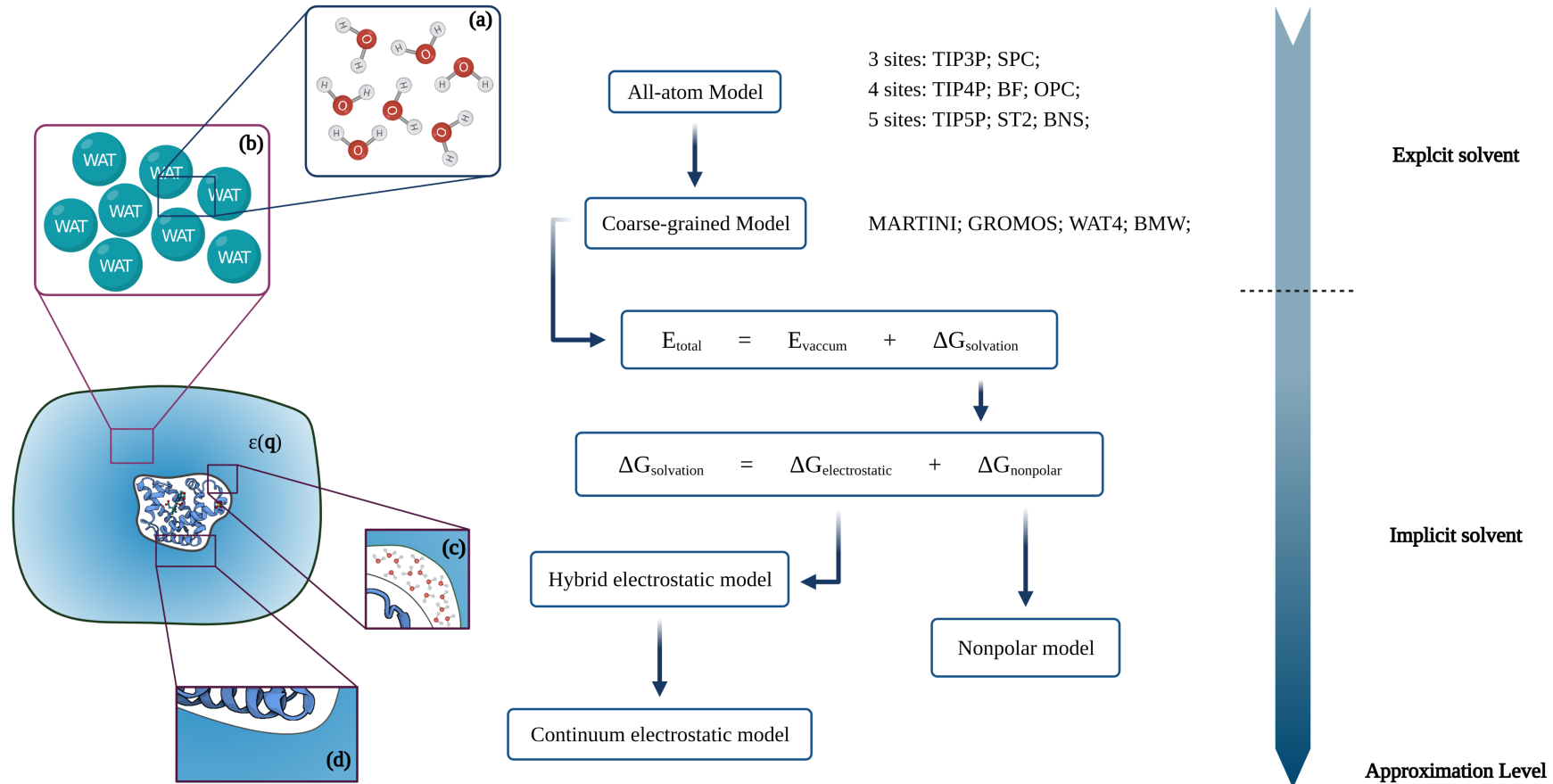
$$\bar{p}(\mathbf{q}_a) = \frac{\int_{\Omega_{a-}} \prod_{i=1}^{3N_{a-}} dq_{a-}^{(i)} \exp\left[-\frac{U(\mathbf{q}_a, \mathbf{q}_{a-})}{k_B T}\right]}{\int_{\Omega} \prod_{i=1}^{3N_a} dq_a^{(i)} \prod_{j=1}^{3N_{a-}} dq_{a-}^{(j)} \exp\left[-\frac{U(\mathbf{q}_a, \mathbf{q}_{a-})}{k_B T}\right]} = \frac{\exp\left[-\frac{F(\mathbf{q}_a)}{k_B T}\right]}{\int_{\Omega_a} \prod_{i=1}^{3N_a} dq_a^{(i)} \exp\left[-\frac{F(\mathbf{q}_a)}{k_B T}\right]}$$

- Then we can rewrite the ensemble average as:

$$\langle A \rangle = \int_{\Omega_a} \prod_{i=1}^{3N_a} dq_a^{(i)} A(\mathbf{q}_a) \bar{p}(\mathbf{q}_a)$$

- , where the DoFs of $a-$ have been integrated out.

Implicit solvent model



- Implicit solvent model treat solvent as a **continuum dielectrics**.

Electrostatic background

- Consider a electrostatic potential ϕ :

$$\Delta G_{\text{ele}} = q [\phi - \phi_{\text{vac}}]$$

- Combine the definition of ϕ and Gauss Law, we give the **Poisson equation**:

$$\nabla \cdot [\epsilon_r \nabla \phi] = -\frac{\rho}{\epsilon_0}$$

- And its special case for spherically symmetric system:

$$\frac{1}{r^2} \frac{d}{dr} \left[r^2 \frac{d\phi}{dr} \right] = -\frac{\rho(r)}{\epsilon_r \epsilon_0}$$

- Consider mean-field theory and Debye-Huckel theory:

$$\frac{1}{r^2} \frac{d}{dr} \left[r^2 \frac{d\phi}{dr} \right] = -\frac{e_0}{\epsilon_r \epsilon_0} \sum_i^{N_i} z_i c_i^0 \exp \left[-\frac{z_i e_0 \phi(r)}{k_B T} \right]$$

- This equation is referred to as **Poisson-Boltzmann Equation**. In the limit of $k_B T \gg z_i e_0 \phi(r)$, we have:

$$\frac{1}{r^2} \frac{d}{dr} \left[r^2 \frac{d\phi}{dr} \right] = -\frac{e_0}{\epsilon_r \epsilon_0} \sum_i^{N_i} z_i c_i^0 \left[1 - \frac{z_i e_0 \phi(r)}{k_B T} \right] = - \left[\frac{e_0^2}{\epsilon_r \epsilon_0 k_B T} \sum_i^{N_i} z_i^2 c_i^0 \right] \phi(r)$$

- , which is called **Linearized Poisson-Boltzmann Equation**.

Generalized Born Model

- Still et al^[1] invent Generalized Born Model follows:

$$G_{\text{ele}} = \frac{1}{2\epsilon} \sum_{i=1}^{N_c} \sum_{j=1, j \neq i}^{N_c} \frac{q_i q_j}{r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^{N_c} \frac{q_i^2}{R_i}$$

- This equation can be rearranged as:

$$G_{\text{ele}} = \frac{1}{2} \sum_{i=1}^{N_c} \sum_{j=1, j \neq i}^{N_c} \frac{q_i q_j}{r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^{N_c} \sum_{j=1, j \neq i}^{N_c} \frac{q_i q_j}{r_{ij}} - \frac{1}{2} \left(1 - \frac{1}{\epsilon}\right) \sum_{i=1}^{N_c} \frac{q_i^2}{R_i}$$

[1]: W Clark Still, Anna Tempczyk, Ronald C Hawley, and Thomas Hendrickson. Semianalytical treatment of solvation for molecular

mechanics and dynamics. Journal of the American Chemical Society, 112(16):6127–6129, 1990

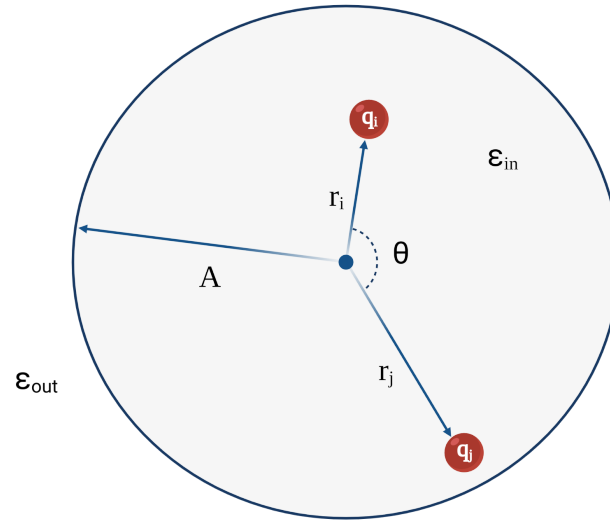
- As $\Delta G_{\text{ele}} = q [\phi - \phi_{\text{vac}}]$, we have:

$$\Delta G_{\text{ele}} = -\frac{1}{2} \left(1 - \frac{1}{\epsilon} \right) \sum_{i=1}^{N_c} \sum_{j=1}^{N_c} \frac{q_i q_j}{f_{GB}}$$

- , where f_{GB} stands for the combination of last two terms. As this expression is similar to the Born energy, the model is referred to as **generalized Born model**. And Still states that:

$$f_{GB} = \sqrt{r_{ij}^2 + R_i R_j \exp \left[-\frac{r_{ij}^2}{\gamma R_i R_j} \right]}$$

- Consider a **two dielectric model** shown below:



- This model is spherically symmetric, describing several charges **embedded** in a dielectric with **low permittivity** ϵ_{in} which is surrounded by a different dielectric with **higher permittivity** ϵ_{out}

- For such a discrete distributed charge system, we can rewrite **Poisson equation** with **Green function**:

$$\nabla^2 \mathbf{G}(\mathbf{r}_i, \mathbf{r}_j) = -\frac{\delta(\mathbf{r}_i - \mathbf{r}_j)}{\epsilon_{\text{in}} \epsilon_0}$$

- , solution of which has the form:

$$\mathbf{G}(\mathbf{r}_i, \mathbf{r}_j) = \frac{1}{4\pi\epsilon_{\text{in}}\epsilon_0} \frac{1}{\|\mathbf{r}_i - \mathbf{r}_j\|} + \mathbf{F}(\mathbf{r}_i, \mathbf{r}_j)$$

- Obviously, $\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j)$ stands for the complex solute-solvent interaction:

$$\Delta G_{\text{ele}} = \frac{1}{2} \sum_{i=1}^{N_c} \sum_{j=1}^{N_c} \mathbf{F}(\mathbf{r}_i, \mathbf{r}_j) q_i q_j$$

- Analytical solution of Poisson Equation for this model exists^[1]:

$$\mathbf{F}(\mathbf{r}_i, \mathbf{r}_i) = -\frac{1}{A} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \sum_{n=0}^{\infty} \frac{t_{ii}^n}{1 + \frac{n}{n+1} \beta}$$

$$\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j) = -\frac{1}{A} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \sum_{n=0}^{\infty} \frac{t_{ij} P_n(\cos\theta)}{1 + \frac{n}{n+1} \beta}$$

- , where $t_{ij} = r_i r_j / A^2$, $\beta = \epsilon_{\text{in}} / \epsilon_{\text{out}}$, and $P_n(\cos\theta)$ are associated Legendre functions.

[1]: John G Kirkwood. Theory of Solutions of Molecules Containing Widely Separated Charges with Special Application to Zwitterions. The

Journal of Chemical Physics, 2(7):351–361, 1934.

- First we consider the solution of $\mathbf{F}(\mathbf{r}_i, \mathbf{r}_i)$. In the limit of $\beta \rightarrow 0$:

$$\Delta G_{\text{ele}}^{ii} = \frac{q_i^2}{2} \mathbf{F}(\mathbf{r}_i, \mathbf{r}_i) = -\frac{1}{2} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \frac{q_i^2}{A(1 - t_{ii})} = -\frac{1}{2} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \frac{q_i^2}{A - \frac{r_i^2}{A}}$$

- As the final result is similar to the Born energy, we name the denominator $R_i = A - r_i^2/A$ after **effective Born radius**.

- Next consider $\mathbf{F}(\mathbf{r}_i, \mathbf{r}_j)$, we first replace $n/(n+1)$ with a constant α , for all $n \geq 1$, rewriting the series as:

$$1 + \frac{1}{1 + \alpha\beta} - \frac{1}{1 + \alpha\beta} + \sum_{n=1}^{\infty} \frac{t_{ij}^n P_n(\cos\theta)}{1 + \alpha\beta} = \frac{1}{1 + \alpha\beta} \left[\sum_{n=0}^{\infty} \left[\frac{t_{ij}^n P_n(\cos\theta)}{1 + \alpha\beta} \right] + \alpha\beta \right]$$

- , where $\sum_{n=0}^{\infty} t^n P_n(\cos\theta)$ can be simplified as $1/\sqrt{1 - 2t_{ij}\cos\theta + t_{ij}^2}$:

$$\Delta G_{\text{ele}}^{ij} = \frac{q_i q_j}{2} \mathbf{F}(\mathbf{r}_i, \mathbf{r}_j) = -\frac{1}{2} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \frac{q_i q_j}{1 + \alpha\beta} \left[\frac{1}{\sqrt{r_{ij}^2 + R_i R_j}} + \frac{\alpha\beta}{A} \right]$$

- , which yields $f_{GB} = \sqrt{r_{ij}^2 + R_i R_j}$

- Now we need to estimate the effective Born radius. In classical electrostatic:

$$G_{\text{ele}} = \frac{1}{2} \int_{\Omega} \rho \phi dV = \frac{1}{8\pi} \int_{\Omega} \mathbf{E} \cdot \mathbf{D} dV$$

- Then we introduce a essential approximation, which state \mathbf{D} has the same form as Coulomb field $\mathbf{D}_i \approx \frac{q_i}{r^2} \hat{\mathbf{r}}_i$:

$$G_{\text{ele}}^{ii} = \frac{1}{8\pi} \left[\int_{\Omega_{\text{in}}} \frac{q_i q_i}{\epsilon_{\text{in}} r^4} dV + \int_{\Omega_{\text{out}}} \frac{q_i q_i}{\epsilon_{\text{out}} r^4} dV \right]$$

$$\Delta G_{\text{ele}}^{ii} = G_{ii} - \frac{1}{8\pi} \int_{\Omega} \frac{q_i^2}{\epsilon_{\text{in}} r^4} dV = -\frac{1}{8\pi} \left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \int_{\Omega_{\text{out}}} \frac{q_i^2}{r^4} dV$$

- Recall the expression of $\Delta G_{\text{ele}}^{ii}$, we have:

$$\frac{1}{R_i} = \frac{1}{4\pi} \int_{\Omega_{\text{out}}} \frac{1}{r^4} dV$$

- As the boundary of Ω_{in} is usually complex, we rewrite $\Delta G_{\text{ele}}^{ii}$ as:

$$\frac{1}{R_i} = \frac{1}{4\pi} \left[\int_0^k \frac{1}{r^2} 4\pi r^2 dr - \int_{\Omega_{\text{in}} - \Omega_s} \frac{1}{r^4} dV \right] = \frac{1}{k} - \frac{1}{4\pi} \int_{\Omega_{\text{in}} - \Omega_s} \frac{1}{r^4} dV$$

- , where Ω_s is a sphere volume with radius k .

- Finally we can consider the influence of mobile ion based on the Debye-Huckel theory by replacement:

$$\left(\frac{1}{\epsilon_{\text{in}}} - \frac{1}{\epsilon_{\text{out}}} \right) \rightarrow \left(\frac{1}{\epsilon_{\text{in}}} - \frac{\exp[-\kappa f_{GB}]}{\epsilon_{\text{out}}} \right)$$

- , where κ is the Debye screening radius:

$$\kappa = \frac{e_0^2}{\epsilon_r \epsilon_0 k_B T} \sum_i^{N_i} z_i^2 c_i^0$$

Enhanced sampling

Fundamentals

- Consider a mapping:

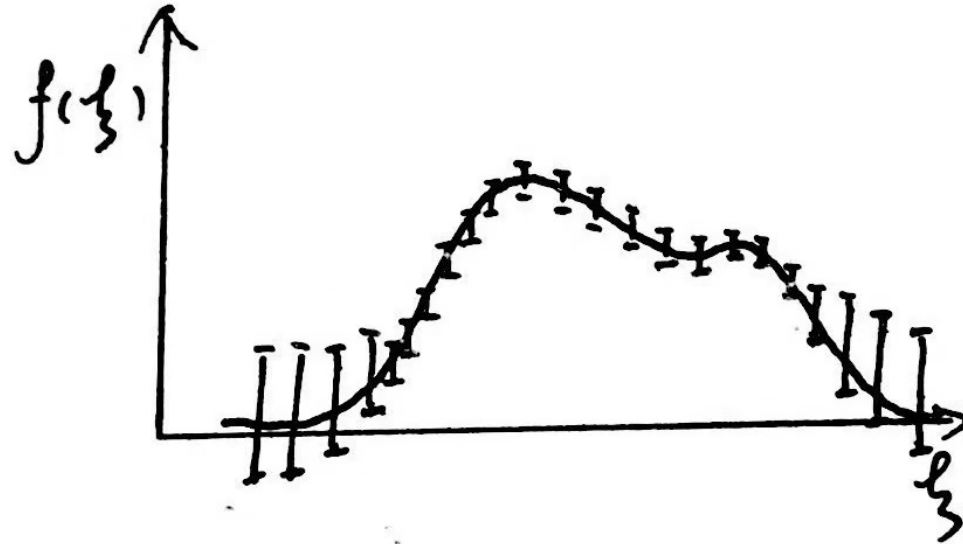
$$\theta : \Omega \in \mathbb{R}^{3N} \rightarrow \Theta \in \mathbb{R}^n, n \ll 3N$$

- Then we can analyze free energy based on the distribution of θ :

$$F(\theta') = -k_B T \ln [H(\theta')]$$

- Advantages:
 - Sampling as a $3N - n$ dimensional hyper-sphere with high efficiency;
 - It's easier to manipulate CV, overcoming the energy barrier;
 - It's easier to compare to the experimental observations;

Umbrella sampling



- In simulation, PDF can be derived from histogram on collective variables. However, the estimation on the **rare events** is accompanied with **high statistical error**

- Simulation with bias potential:

$$\tilde{U}(\mathbf{R}) = U(\mathbf{R}) + \Delta U(\mathbf{R})$$

where three terms corresponds to new potential, origin potential and bias potential, respectively.

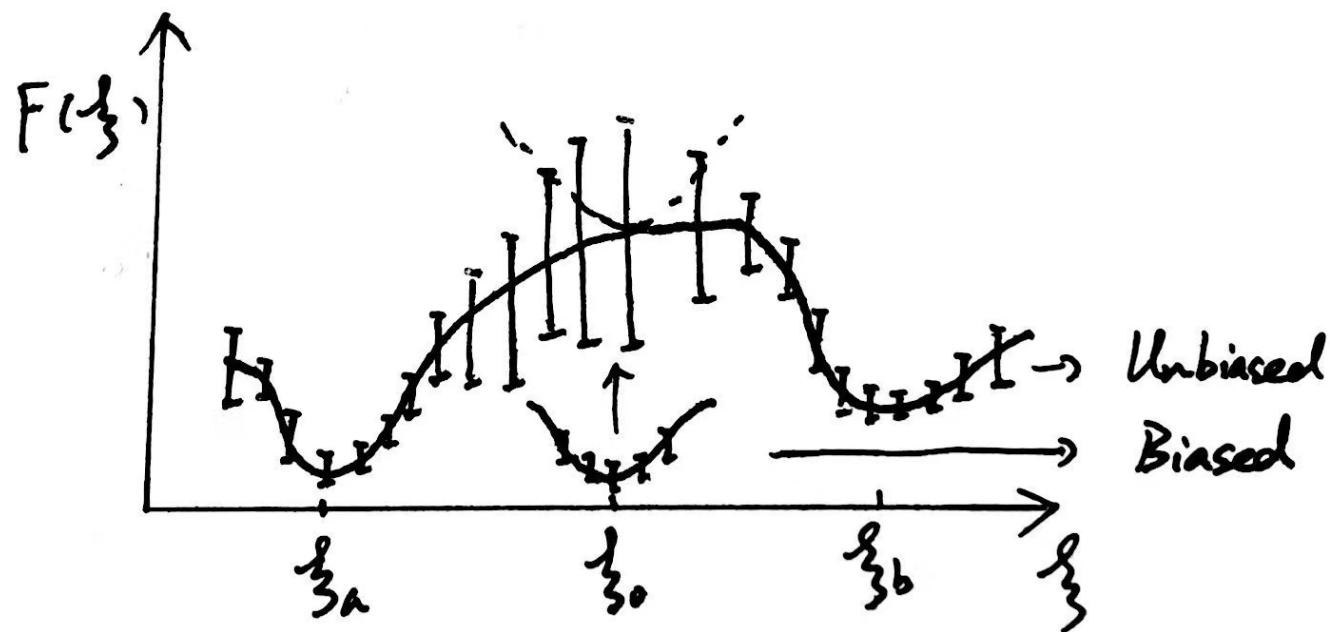
- So we have:

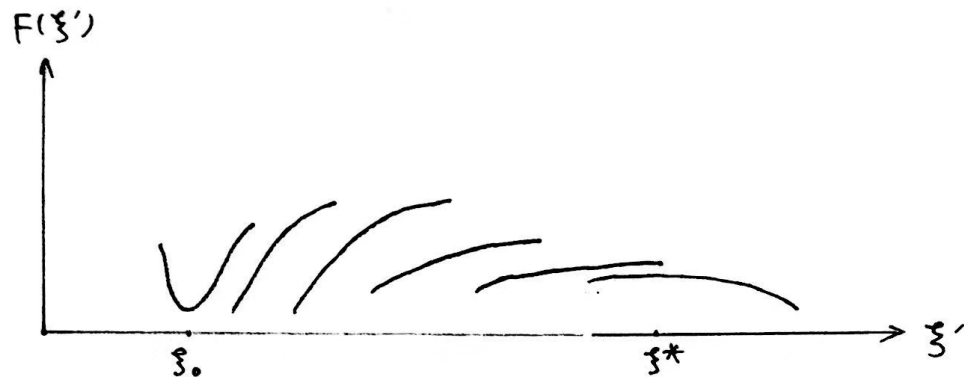
$$\begin{aligned}\tilde{F}(\xi') &= -k_B T \ln \left(\frac{1}{Z} \int d^{3N} R e^{-\beta U(\mathbf{R})} e^{-\beta \Delta U(\mathbf{R})} \delta(\xi(\mathbf{R}) - \xi') \right) \\ &= F(\xi') + \Delta U(\xi') + C\end{aligned}$$

- Usually, a potential like harmonic oscillator is a good choice:

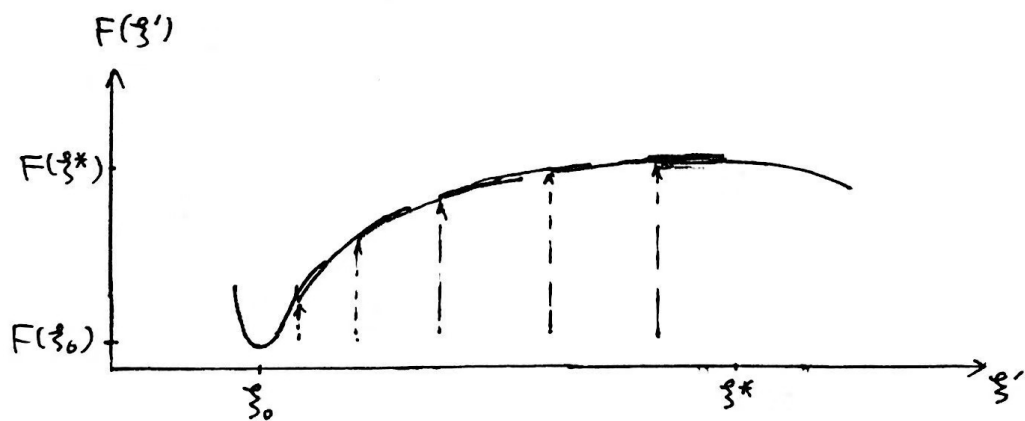
$$\Delta U(\mathbf{R}) = \frac{1}{2} K_{bias} (\xi(\mathbf{R}) - \xi_0)^2$$

where K_{bias} and ξ_0 are two hyper parameters for umbrella sampling





- Suppose we perform a series of simulation over different bias potential
- Each histogram is overlapped with others (like umbrella)
- We can reconstruct the free energy with a series curve with little error



Weighted Histogram Analysis Method

- Initial data:
 - A collective variable $\xi : \Omega \rightarrow \Xi \in \mathbb{R}^d, d \ll 6N$
 - M independent trajectories (usually with a biased potential), $S_i, i \in [1, M]$
 - The histogram and probability distribution (or probability density) of each trajectories on reaction coordinate: $h_i(\xi')$ and $p_i(\xi'), i \in [1, M]$
- Final goal:
 - The probability distribution (or probability density) on reaction coordinate of unbiased simulation

- **Initial Assumption:** the biased probability distribution $p_i(\xi')$ is related to the unbiased one $p_0(\xi')$:

$$p_i(\xi') = f_i c_i(\xi') p_0(\xi')$$

where $c_i(\xi')$ is the biasing factor corresponds to the influence of bias potential:

- $c_i(\xi') = \exp[-(\beta_i - \beta_0)E_0(\xi')]$ for temperature biasing
- $c_i(\xi') = \exp[-\beta\Delta U(\xi')]$ for coordinate biasing

- And f_i is a normalization factor:

$$\begin{aligned}
 f_i^{-1} &= \int_{\Xi} c_i(\xi') p_0(\xi') d\xi' \\
 &= \int_{\Xi} e^{-\beta \Delta U(\xi')} \frac{1}{Q_0} \int_{\Omega} e^{-\beta U(\mathbf{q})} \delta[\xi(\mathbf{q}) - \xi] d\mathbf{q} d\xi \\
 &= \frac{1}{Q_0} \int_{\Xi} \int_{\Omega} e^{-\beta [U(\mathbf{q}) + \Delta U(\mathbf{q})]} \delta[\xi(\mathbf{q}) - \xi] d\mathbf{q} d\xi \\
 &= \frac{1}{Q_0} \int_{\Omega} e^{-\beta [U(\mathbf{q}) + \Delta U(\mathbf{q})]} d\mathbf{q} = \frac{Q_i}{Q_0}
 \end{aligned}$$

- As we know:

$$\Delta A_{ij} = -k_B T \ln \frac{Q_j}{Q_i} = -k_B T \ln \frac{\int_{\Gamma_j} \exp \left[-\frac{U_j(\mathbf{q})}{k_B T} \right] d\mathbf{q}}{\int_{\Gamma_i} \exp \left[-\frac{U_i(\mathbf{q})}{k_B T} \right] d\mathbf{q}}$$

- So we have:

$$\beta^{-1} \ln f_i = \Delta A_{0i}$$

- Based on the initial assumption:

$$p_i(\xi') = f_i c_i(\xi') p_0(\xi')$$

- We can define the estimation of unbiased result from i trajectory:

$$p_0^{(i)}(\xi') = \frac{p_i(\xi)}{f_i c_i(\xi')} = \frac{h(\xi')}{N_i f_i c_i(\xi')}$$

where $h(\xi')$ is the histogram value and N_i is the total samples of trajectory i .

- And we give the final estimation of $p_0(\xi')$ as a weighted average:

$$p_0^{est}(\xi') = \sum_{i=1}^M w_i p_0^{(i)}(\xi')$$

- Obviously, w_i satisfy the normalization condition: $\sum_{i=1}^M w_i = 1$.
- Clearly, if we want to get a good estimation, we want the variance of each data point to be the minimum:

$$\begin{aligned}
 \text{var}(p_0^{est}(\xi')) &= \langle (p_0^{est}(\xi') - \langle p_0^{est}(\xi') \rangle)^2 \rangle \\
 &= \left\langle \left(\sum_{i=1}^M w_i p_0^{(i)}(\xi') - \left\langle \sum_{i=1}^M w_i p_0^{(i)}(\xi') \right\rangle \right)^2 \right\rangle \\
 &= \left\langle \left(\sum_{i=1}^M w_i \left(p_0^{(i)}(\xi') - \langle p_0^{(i)}(\xi') \rangle \right) \right)^2 \right\rangle
 \end{aligned}$$

We define $\delta p_0^{(i)}(\xi') = \left(p_0^{(i)}(\xi') - \langle p_0^{(i)}(\xi') \rangle \right)$. Then we have:

$$\begin{aligned}
 \text{var} (p_0^{\text{est}}(\xi')) &= \left\langle \left(\sum_{i=1}^M w_i \delta p_0^{(i)}(\xi') \right)^2 \right\rangle \\
 &= \left\langle \sum_{i=1}^M w_i^2 (\delta p_0^{(i)}(\xi'))^2 + \sum_{j=1}^M \sum_{k \neq j}^M w_j w_k \delta p_0^{(j)}(\xi') \cdot \delta p_0^{(k)}(\xi') \right\rangle \\
 &= \sum_{i=1}^M w_i^2 \left\langle (\delta p_0^{(i)}(\xi'))^2 \right\rangle + \sum_{j=1}^M \sum_{k \neq j}^M w_j w_k \left\langle \delta p_0^{(j)}(\xi') \cdot \delta p_0^{(k)}(\xi') \right\rangle
 \end{aligned}$$

- Always, we can make the assumption of independent sampled trajectory, which means $\langle \delta p_0^{(j)}(\xi') \cdot \delta p_0^{(k)}(\xi') \rangle = 0$ and we have $\langle (\delta p_0^{(i)}(\xi'))^2 \rangle = \text{var}(p_0^{(i)}(\xi'))$, so we can rewrite equation as:

$$\text{var}(p_0^{est}(\xi')) = \sum_{i=1}^M w_i^2 \text{var}(p_0^i(\xi'))$$

- As we know $\text{var}(ax) = \langle a^2 x^2 \rangle - \langle ax \rangle^2 = a^2 \text{var}(x)$, so recall the definition of $p_0^{(i)}(\xi') = \frac{p_i(\xi)}{f_i c_i(\xi')} = \frac{h(\xi')}{N_i f_i c_i(\xi')}$, we have:

$$\text{var}(p_0^{est}(\xi')) = \sum_{i=1}^M \frac{w_i^2 \text{var}(h(\xi))^2}{N_i^2 c_i^2(\xi) f_i^2}$$

- If we have N independent samples, the probability and variance of having n counts in one bin will follow the **binomial distribution**
 - Mean: $np_i(\xi')$
 - Variance: $np_i(\xi')(1 - p_i(\xi'))$
- In the limit of large N and small $p_i(\xi')$ (exactly what happened for a long simulation and a detailed histogram), the binomial distribution can be treated as a **Poisson distribution**:

$$P(n) = \exp(-Np_i(\xi')) \frac{(Np_i(\xi'))^n}{n!}$$

- Means and variance: $Np_i(\xi')$

- Substituting the variance so we have:

$$\text{var} (p_0^{est}(\xi')) = \sum_{i=1}^M \frac{w_i^2 N_i f_i c_i(\xi') p_0(\xi')}{N_i^2 c_i^2(\xi) f_i^2} = \sum_{i=1}^M \frac{w_i^2 p_0(\xi')}{N_i c_i(\xi) f_i}$$

- We want to minimize equation with the constraint $\sum_{i=1}^M w_i = 1$, so we define a Lagrange multiplier λ , giving a goal function:

$$G = \sum_{i=1}^M \frac{w_i^2 p_0(\xi')}{N_i c_i(\xi) f_i} + \lambda \sum_{i=1}^M w_i$$

- To minimize the goal function, first take the partial derivative:

$$\frac{\partial G}{\partial w_i} = \frac{2w_i p_0(\xi')}{N_i c_i(\xi') f_i} + \lambda$$

- And setting it to zero, obtain:

$$w_i = -\frac{N_i c_i(\xi') f_i}{2p_0(\xi')} \lambda$$

- Substitute to the constraints:

$$\sum_{i=1}^M -\frac{N_i c_i(\xi') f_i}{2p_0(\xi')} \lambda = 1$$

- so we have:

$$\lambda = -\frac{2p_0(\xi')}{\sum_{i=1}^M N_i c_i(\xi') f_i}$$

- and obtain the weight:

$$w_i = \frac{N_i c_i(\xi') f_i}{\sum_{i=1}^M N_i c_i(\xi') f_i}$$

Substitute this weight and finally we get:

$$p_0^{est}(\xi') = \sum_{i=1}^M w_i \frac{h(\xi')}{N_i f_i c_i(\xi')} = \frac{\sum_{i=1}^M h_i(\xi')}{\sum_{i=1}^M N_i f_i c_i(\xi')}$$

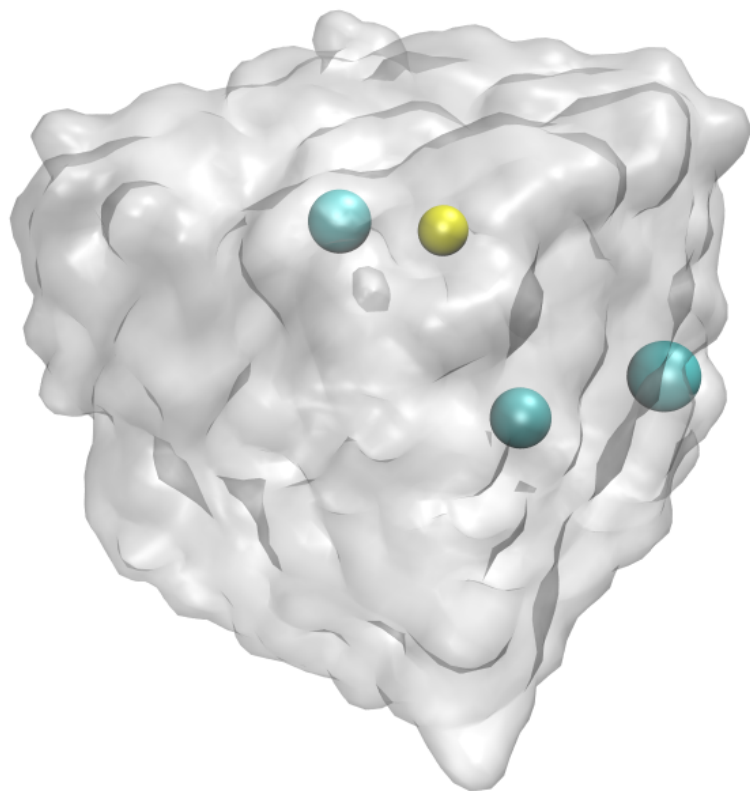
This equation and the normalization condition:

$$f_i^{-1} = \int_{\Xi} c_i(\xi') p_0^{est}(\xi') d\xi'$$

are collectively known as the **Weighted Histogram Analysis Method (WHAM)** equations.

Free energy of ion interaction in solution

Model detail

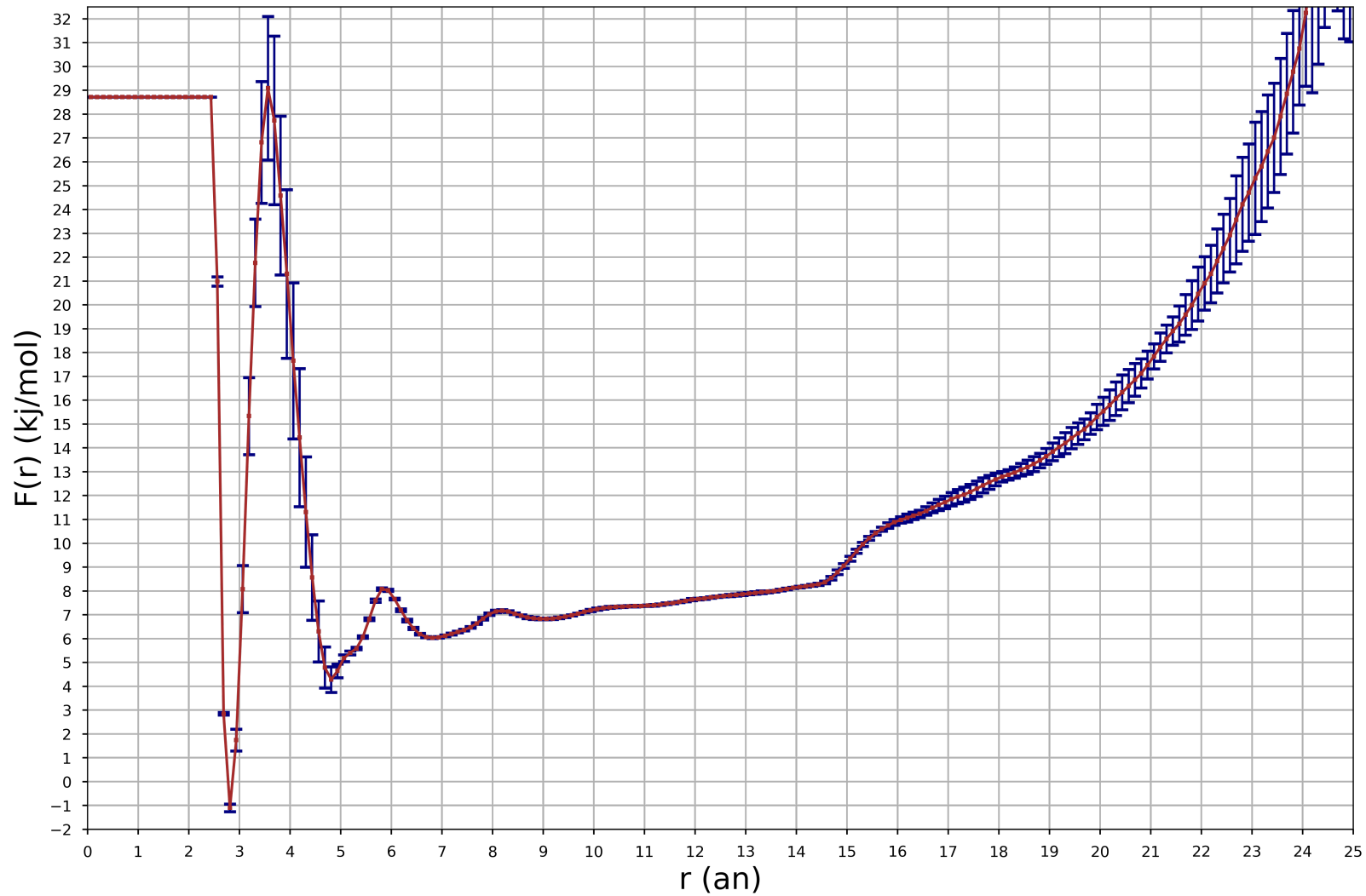


- Simulation tool: OpenMM
- Number of atoms: 2072
- Thermostat: Andersen Thermostat
- Barostat: MonteCarlo Barostat
- Time step: 1 femtosecond
- Equilibrium temperature: 300K
- Equilibrium pressure: 1 Bar

Umbrella sampling settings

- **10** replicas
- Each has **200** trajectories
- Each run for **200ns**
- Each executed for **10 hours**

SCHM Result



MDAnalyser Demo

```
from mdanalyser.analyser import WHAMUmbrellaSamplingAnalyser

path_current = sys.path[0]
path_out_cv = path_current + '/../output/cv_files'

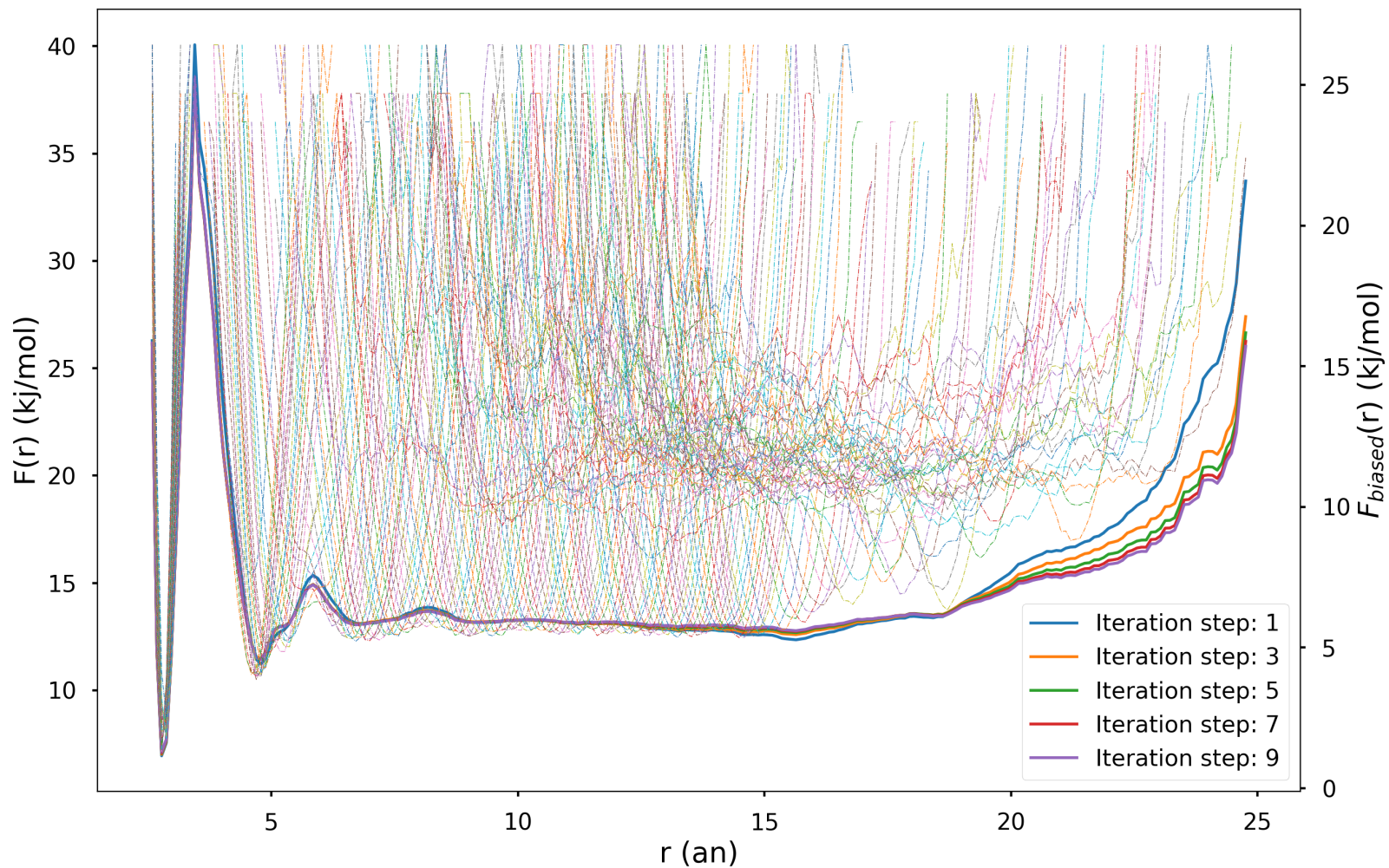
analyzer = WHAMUmbrellaSamplingAnalyser()

analyzer.setDirPath(path_out_cv)
analyzer.setFileIdRangeFromStartToEnd(0, 99)
analyzer.setBinRangeFromStartToEnd(1, 15, 250)
analyzer.setConstantK(100)
analyzer.setTemperature(300)

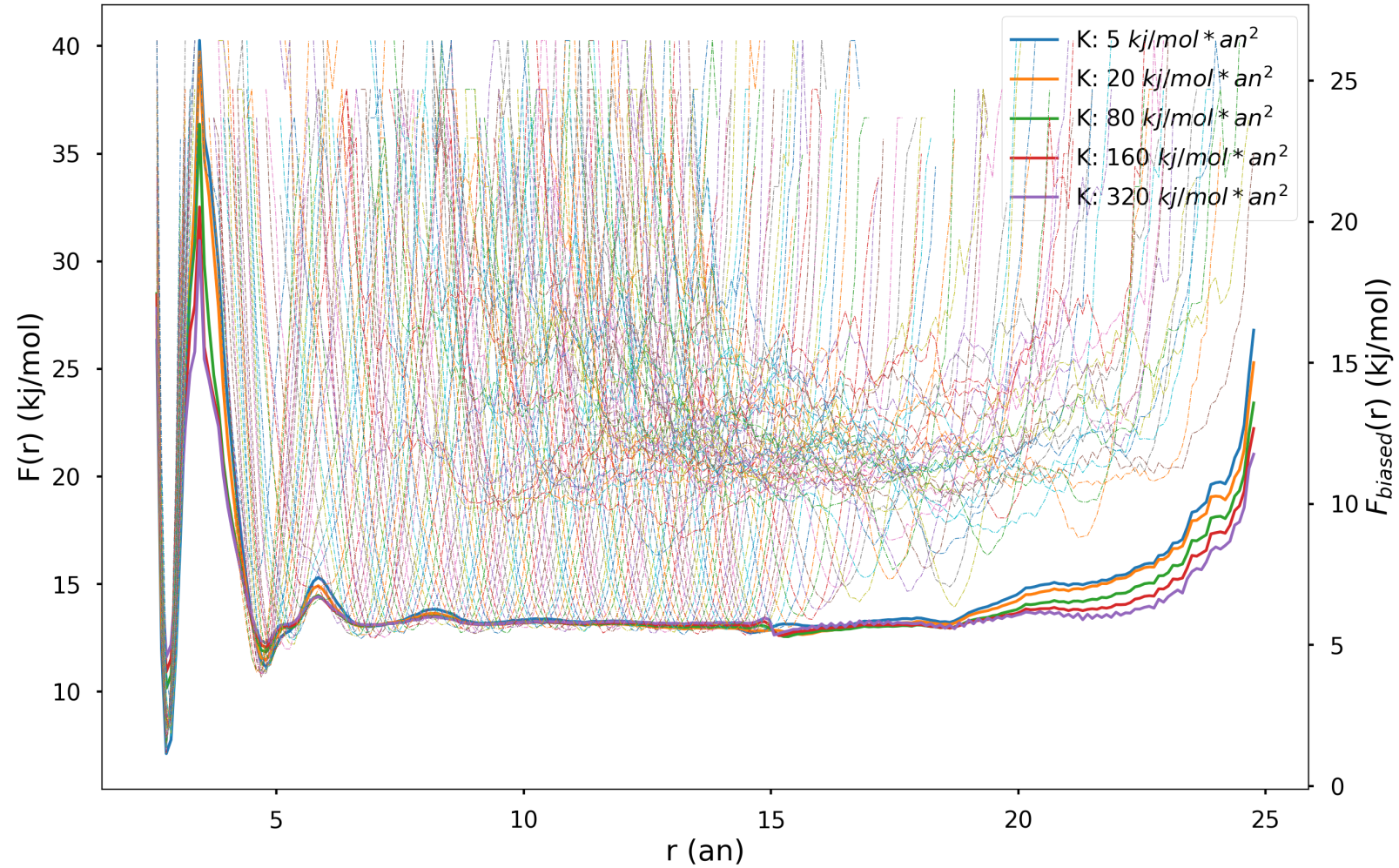
analyzer.loadData()
analyzer.iterativeSolver(num_steps=500)
```

Validation:

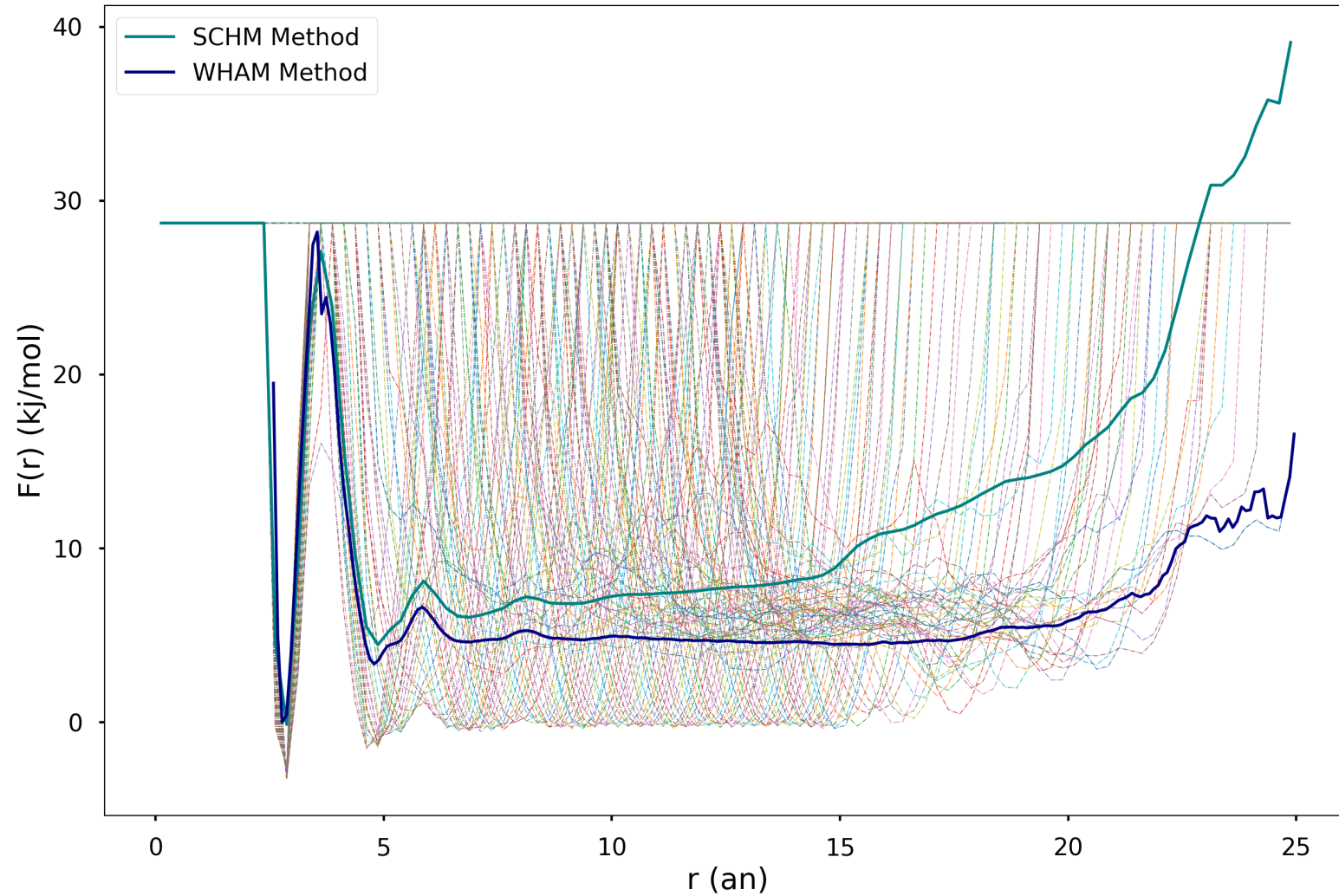
- Quickly convergence



- Sensitive to the constant K of harmonic oscillator bias potential

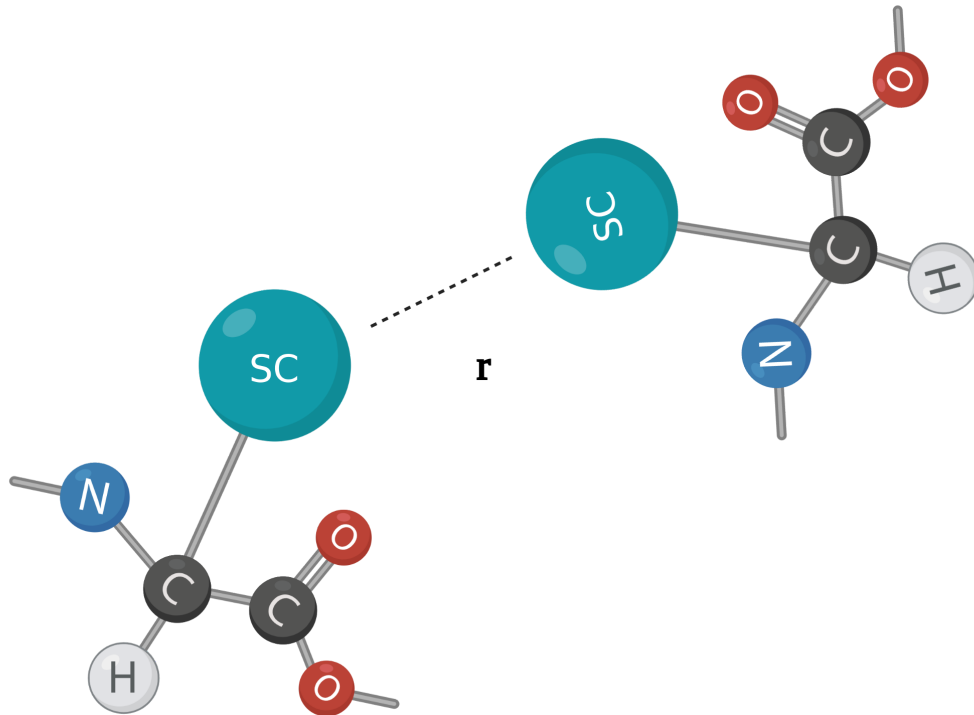


- Better performance than **Self-Consistent Histogram Method (SCHM)** in sparse sampling part



Free energy matrix

Model details



- Two amino acid
- number of atoms: around 5000
- NaCl concentration: 0.15mol/L
- time step: 1 fs
- thermostat: Langevin
- barostat: MonteCarlo 1bar
- integrator: Langevin Integrator
- forcefield: 'amber14/tip3pfb.xml'

Collective variable details

- **Defination:**

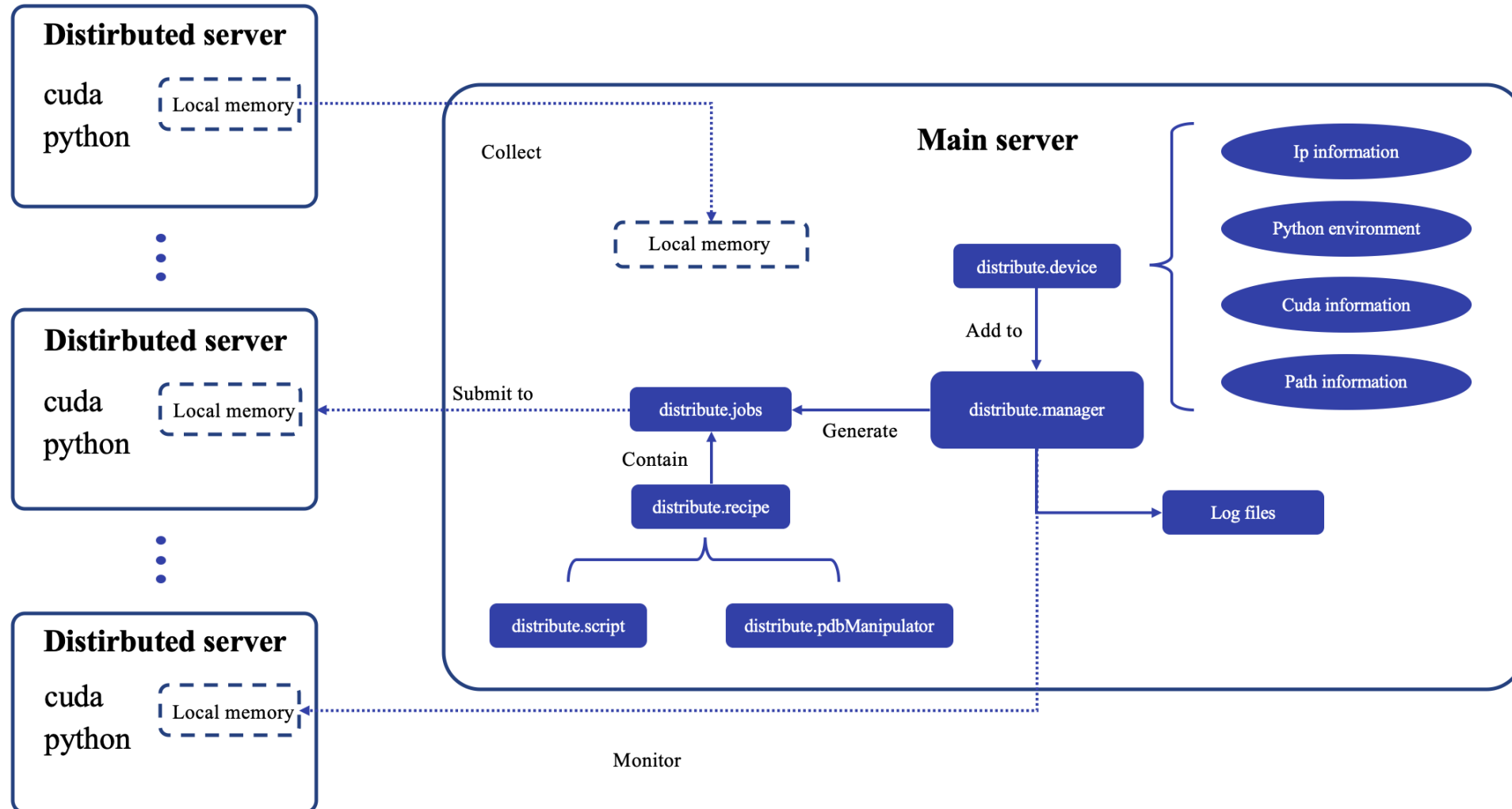
The distance between the centroid center of two amino acid

$$\Delta U = \frac{1}{2}K(\mathbf{c}_{dist} - \mathbf{c}_0)^2, \quad \mathbf{c}_{dist} = \frac{1}{M_{tot}} \sum_{i=1}^N m_i \mathbf{q}_i$$

$$F_{x_i} = \frac{\partial \Delta U}{\partial c_{dist,x}} \frac{\partial c_{dist,x}}{\partial x_i} = \frac{m_i}{M_{tot}} K c_{dist,x}$$

- **Range:** $r \in [1, 20] \text{ \AA}$
- **Grid:** 200 trajectories: $150 \in [1, 15), 50 \in [10, 20)$

Pdff-distribute



Result

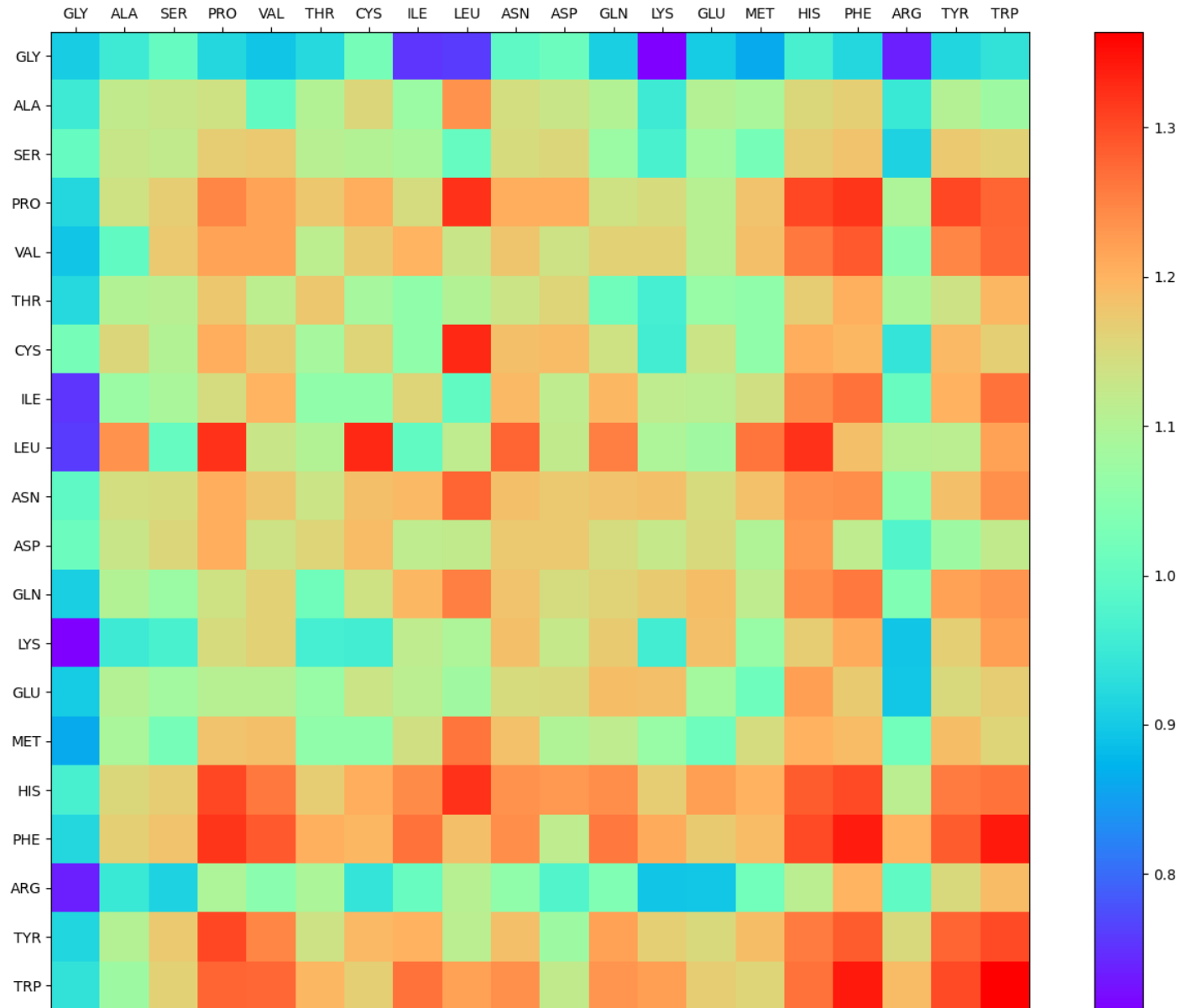
	ALA	ARG	ASN	ASP	CYS	GLN	GLU	GLY	HIS
ALA	1.119	0.949	1.142	1.129	1.156	1.103	1.103	0.954	1.156
ARG	0.949	0.996	1.058	0.98	0.942	1.039	0.897	0.735	1.111
ASN	1.142	1.058	1.186	1.174	1.185	1.18	1.148	0.995	1.201
ASP	1.129	0.98	1.174	1.172	1.192	1.145	1.151	1.013	1.201
CYS	1.156	0.942	1.185	1.192	1.157	1.137	1.134	1.026	1.201
GLN	1.103	1.039	1.18	1.145	1.137	1.16	1.188	0.908	1.201
GLU	1.103	0.897	1.148	1.151	1.134	1.188	1.082	0.901	1.201

氨基酸性质表

	缩写	名称	结构	分子量	等电点	溶解度 水 (0,20°C)g/L	分类
G	Gly	甘氨酸 Glycine		75.07	6.06	141.8 225.2	脂肪族类
A	Ala	丙氨酸 Alanine		89.09	6.11	127.3 157.8	脂肪族类
V	Val	缬氨酸 Valine		117.15	6	59.6 68.1	脂肪族类
L	Leu	亮氨酸 Leucine		131.17	6.01	22.70 23.74	脂肪族类
I	Ile	异亮氨酸 Isoleucine		131.17	6.05	37.91 40.25	脂肪族类
F	Phe	苯丙氨酸 Phenylalanine		165.19	5.49	19.83 27.35	芳香族类
W	Trp	色氨酸 Tryptophan		204.23	5.89	8.23 10.57	芳香族类
Y	Tyr	酪氨酸 Tyrosine		181.19	5.64	0.196 0.384	芳香族类
D	Asp	天冬氨酸 Aspartic acid		133.1	2.85	2.62 6.33	酸性氨基酸类
H	His	组氨酸 Histidine		155.16	7.6	41.9 (25°C)	碱性氨基酸类
N	Asn	天冬酰胺 Asparagine		132.12	5.41	8.49 23.51	酰胺类
E	Glu	谷氨酸 Glutamic acid		147.13	3.15	8.55 17.22	酸性氨基酸类
K	Lys	赖氨酸 Lysine		146.19	9.6	400 630	碱性氨基酸类
Q	Gln	谷氨酰胺 Glutamine		146.15	5.65	7.2 (37°C)	酰胺类

M	Met	甲硫氨酸 Methionine		149.21	5.74	18.18 29.95	含硫类
R	Arg	精氨酸 Arginine		174.2	10.76	855.6 (25°C)	碱性氨基酸类
S	Ser	丝氨酸 Serine		105.09	5.68	22.04 42.95	羟基类
T	Thr	苏氨酸 Threonine		119.12	5.6	13.2 (25°C)	羟基类
C	Cys	半胱氨酸 Cysteine		121.16	5.05	0.11 (25°C)	含硫类
P	Pro	脯氨酸 Proline		115.13	6.3	1620 (25°C)	亚氨基酸
U	Sec	硒半胱氨酸 Selenocysteine		168.07			
O	Pyl	吡咯赖氨酸 Pyrrolysine		255.31			

Sort by molecular weight

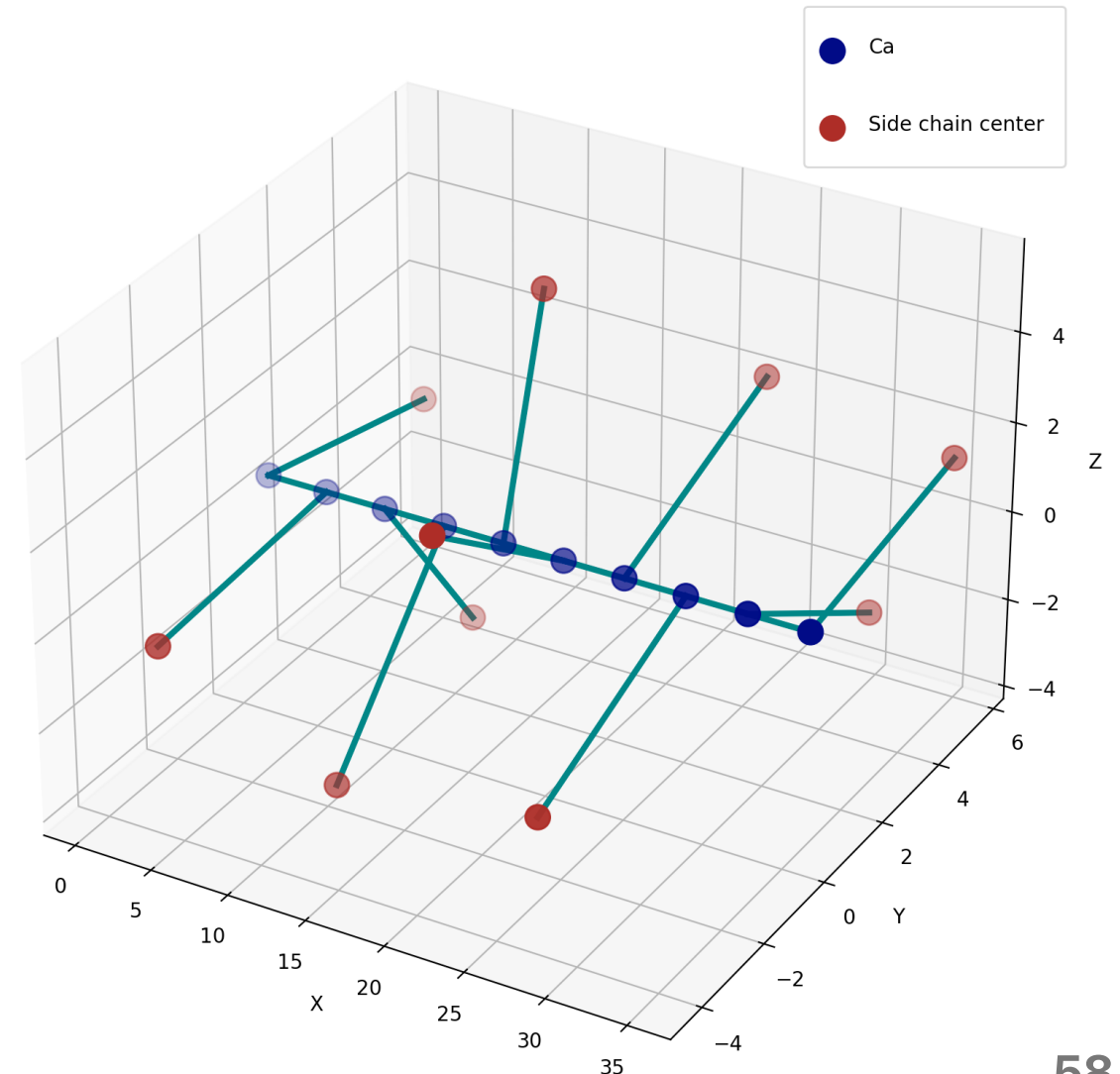


OpenPD

Coarse-Grained Model

The protein is treated as a peptides chain, each of which consists two beads:

- **CA** bead represents the backbone part of the peptide;
- **SC** bead represents the side chain part of the peptide;



Force field expression

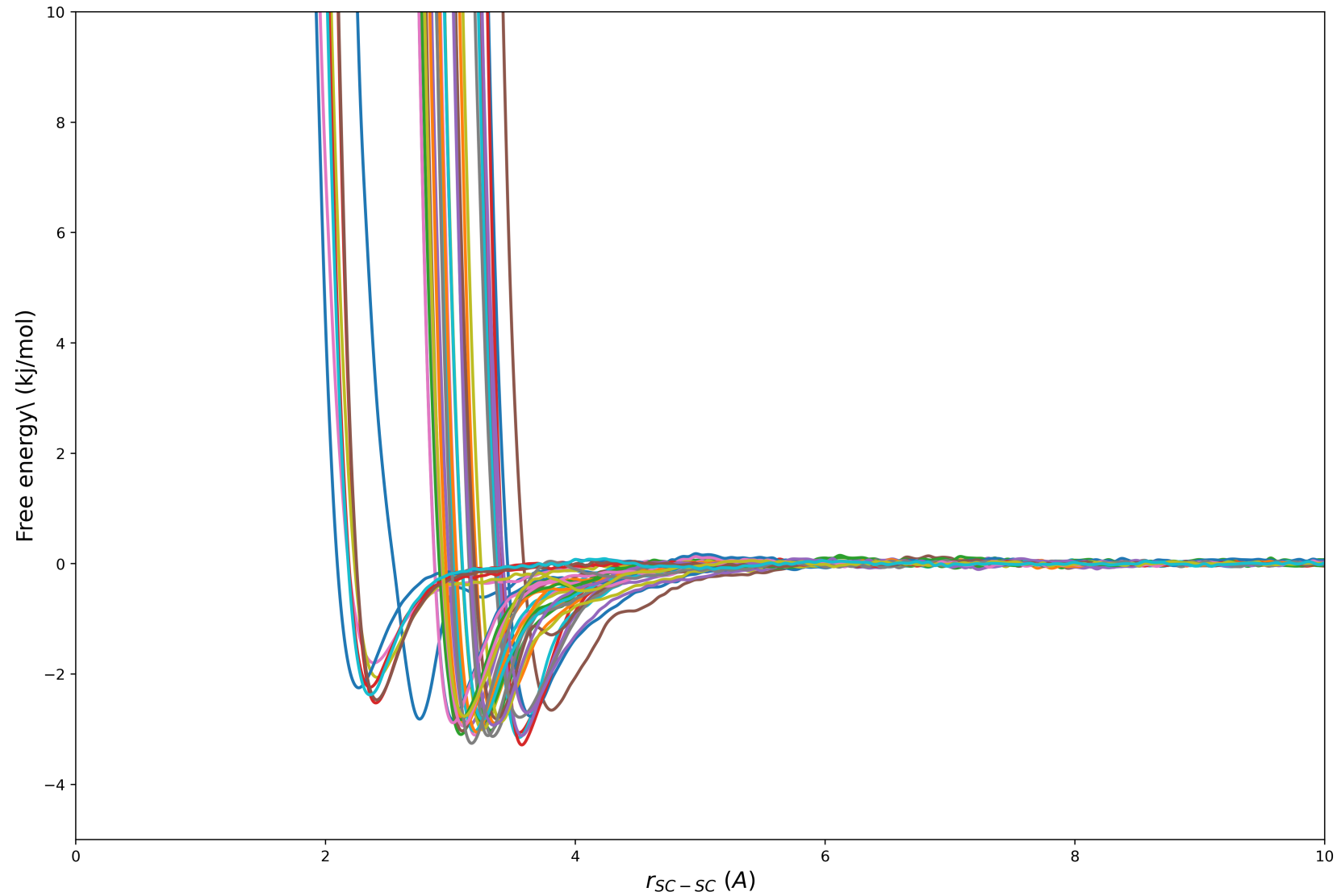
$$\begin{aligned} U_{tot} &= U_{Non-Bonded} + U_{Bond} + U_{Torsion} \\ &= \sum_{i=1}^{N_A} \sum_{j=i+1}^{N_A} U_{Non-Bonded}^{(ij)}(\mathbf{r}_{ij}) + \sum_{i=1}^{N_B} U_{Bond}^{(i)}(l_i) + \sum_{i=1}^{N_T} U_{Torsion}^{(i)}(\theta_i) \end{aligned}$$

- $U_{Non-Bonded}$ describes the interaction between **SC** and **SC** (no neighbor) bead.
- U_{Bond} describes the bond that constraint the connected bead.
- $U_{Torsion}$ describes the torsion angle preference between neighbor **SC** bead.

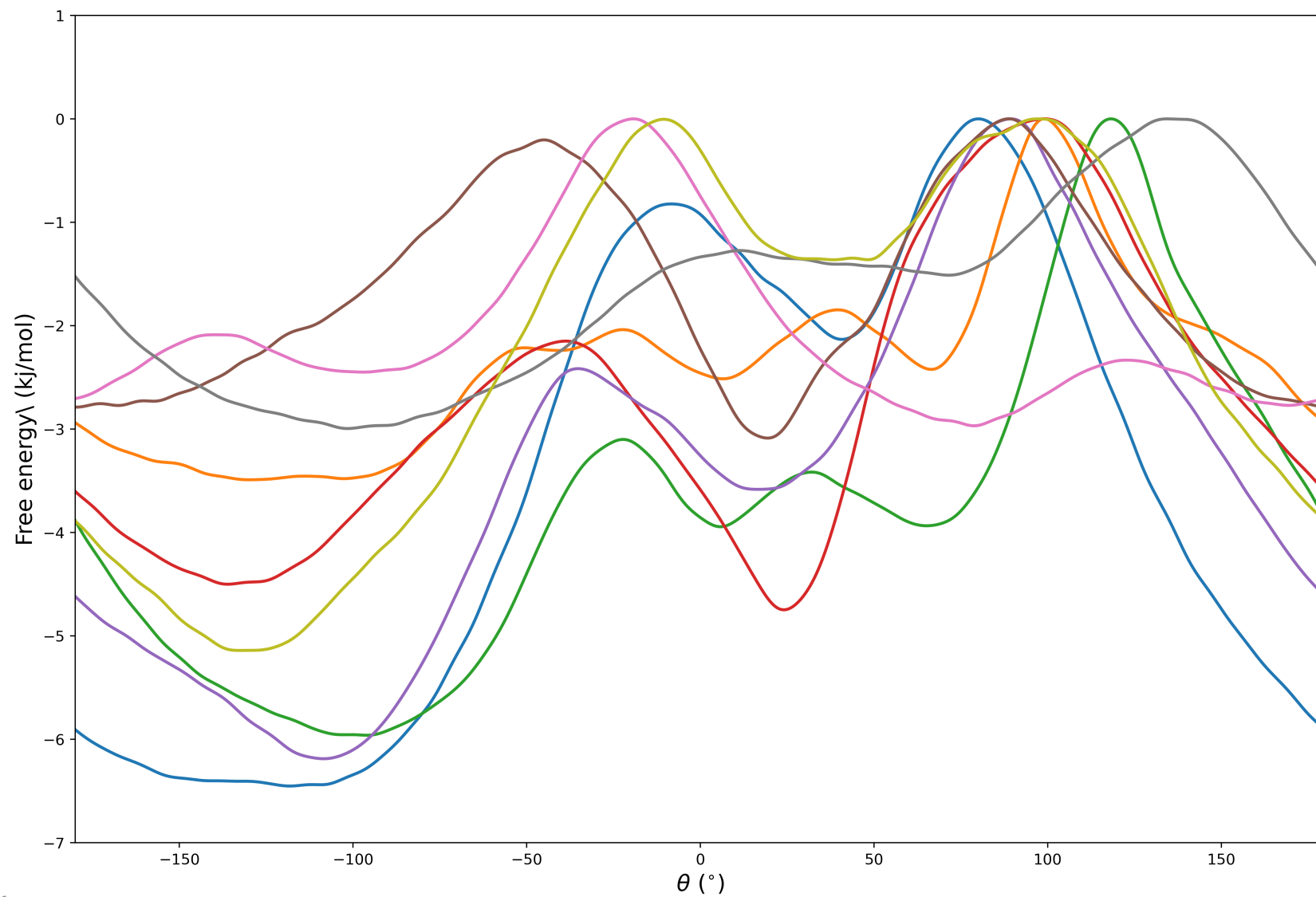
Computational detail

- Non bonded interaction:
 - CV: the distance between centroid of two peptides' side chain;
 - Sampling algorithm: **Umbrella sampling**; 500ns;
- Bond interaction:
 - CV: the distance between two target beads;
 - Sampling algorithm: **Umbrella sampling**; 500ns;
- Torsion interaction:
 - CV: $^2 \angle SC - CA - CA - SC$;
 - Sampling algorithm: **Well-tempered Metadynamics**; 200ns;

Non Bonded Result



Torsion Result



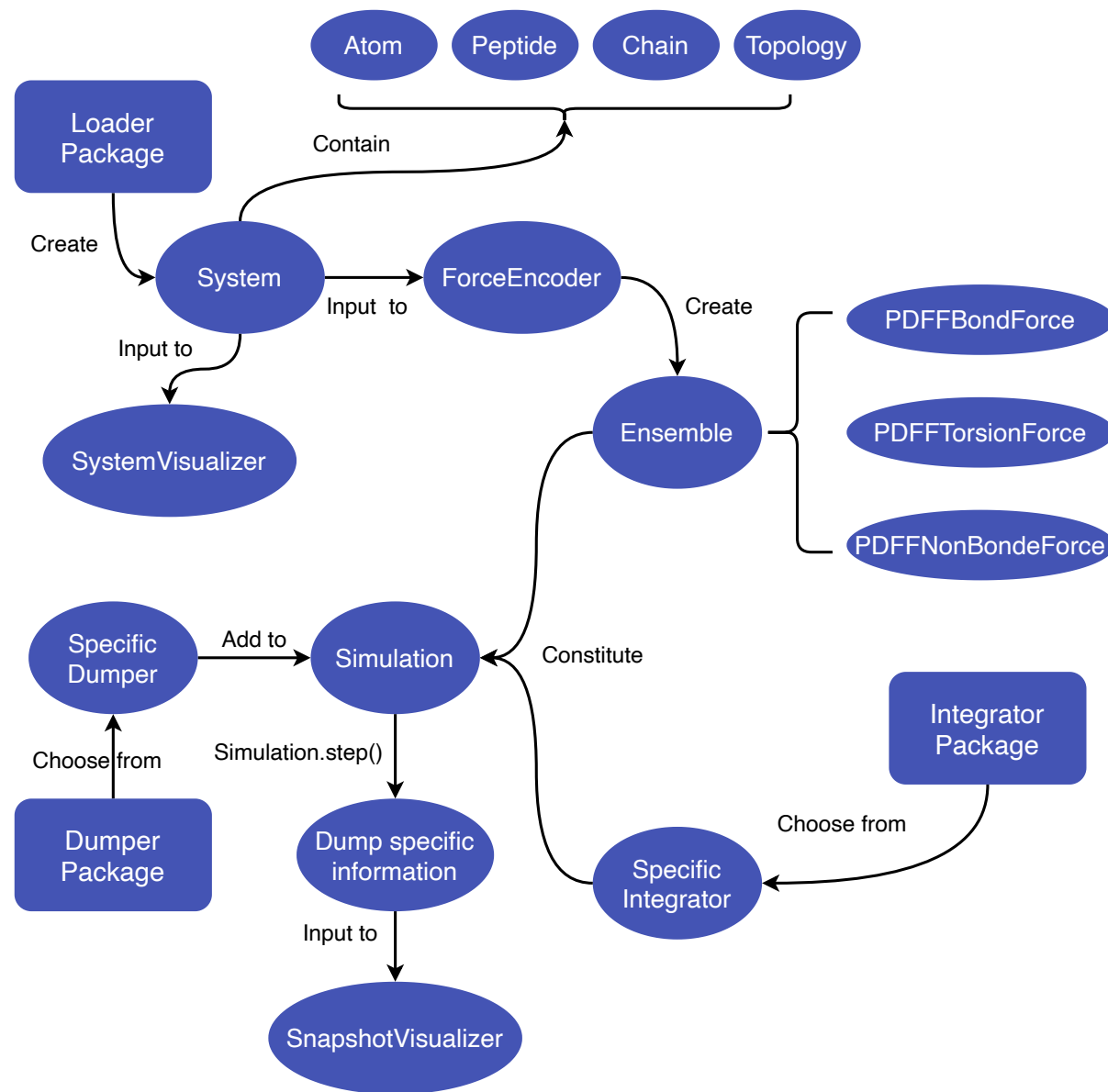
Difference from ordinary coarse grained force field

- PDFF does not contain solvent term;
- PDFF is not fitting from experimental data or other result;
- PDFF is, theoretically, generalized enough for **all** protein, no matter whether exists homologous pair in the database

Overview

OpenPD, standing for **Open Peptide Dynamics**, is a **python** package, distributed freely under the terms of **GPLv3**, for peptide dynamics simulation.

- <https://github.com/openpd-dev/openpd>
- <https://openpd.net>



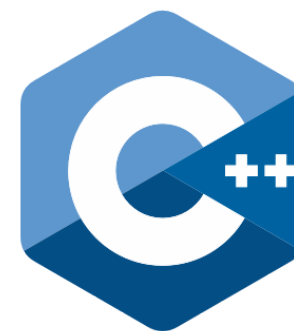
Extensible packages

- **openpd.loader** package:
 - PDBLoader, SequenceLoader
- **openpd.force** package:
 - PDFFNonBondedForce, PDFFBondForce, PDFFTorsionForce
- **openpd.integrator** package:
 - VerletIntegrator, VelocityVerletIntegrator, BrownianIntegrator ...
- **openpd.dumper** package:
 - LogDumper, SnapshotDumper, PDBDumper, XYZDumper

Current



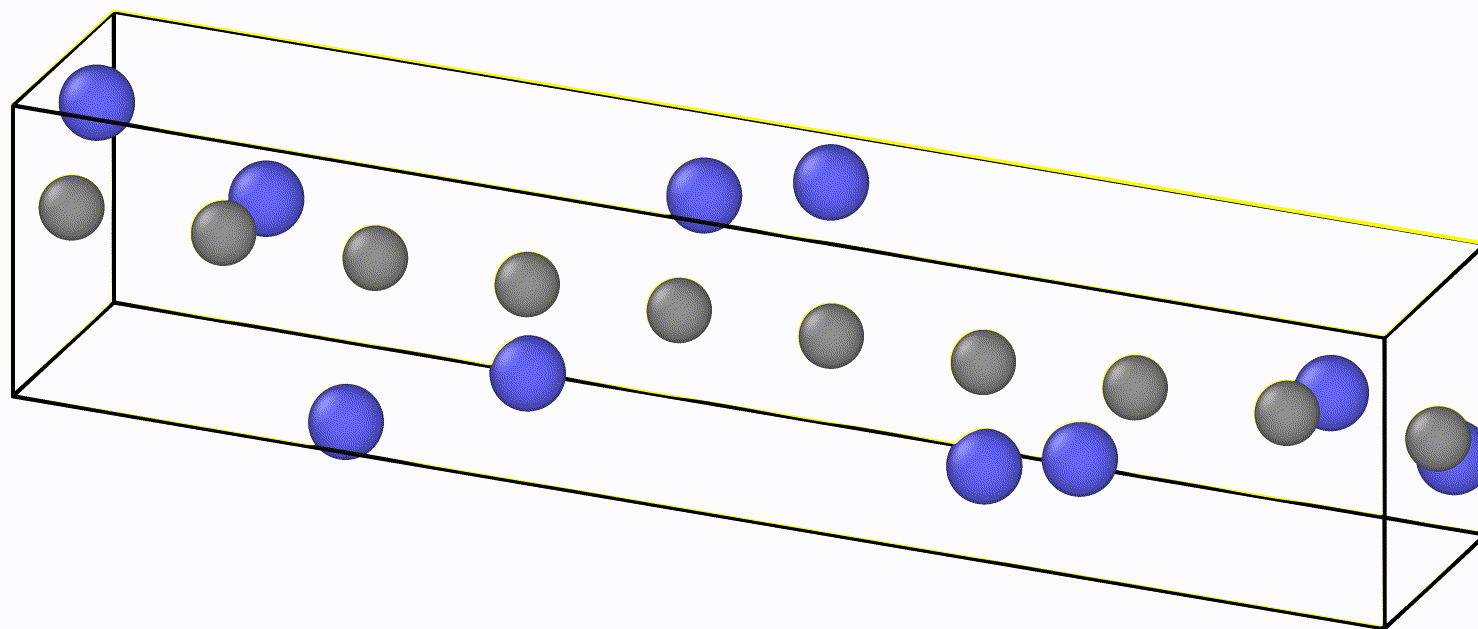
Planned

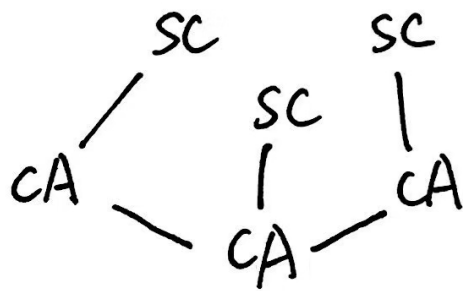
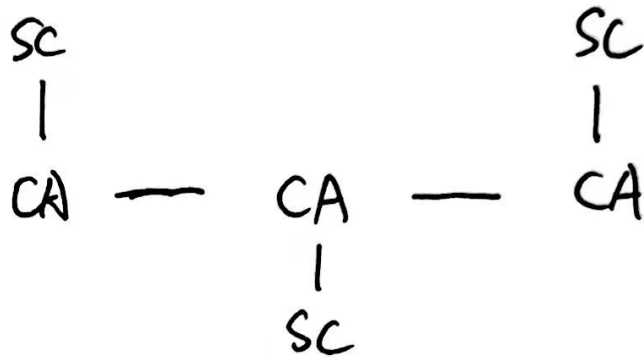


Demo Code

```
1 import openpd as pd
2 from openpd.unit import *
3
4 system = pd.SequenceLoader('data/simulation.json').createSystem()
5 ensemble = pd.ForceEncoder(system).createEnsemble()
6 integrator = pd.VelocityVerletIntegrator(1*femtosecond, 300*kelvin)
7 simulation = pd.Simulation(ensemble, integrator)
8
9 log_dumper = pd.LogDumper(
10     'simulation.log', 50, get_potential_energy=True, get_temperature=True,
11     get_steps=True, get_elapsed_time=True, get_remain_time=True
12 )
13 snapshot_dumper = pd.SnapshotDumper('simulation.pds', 50)
14 xyz_dumper = pd.XYZDumper('simulation.xyz', 50)
15 simulation.addDumpers(log_dumper, snapshot_dumper, xyz_dumper)
16
17 simulation.minimizeEnergy('gd', max_iteration=20)
18 simulation.step(50000)
```

Result





Discussion

Non-Bonded interaction only considered the **SC** bead;

The constraint on **CA** bead is completely **locally**:

- Neglect the interaction of distance between **CA** and non-neighbored **SC**;
- Neglect the interaction of distance between **CA** and non-neighbored **CA**;

Future proposal

- Implicit solvent
 - Machine learning solver for PE, PBE, LPBE
 - Solvent boundary potential for **arbitrary boundary**
- Coarse grained model
 - Refinement of PDFF
- Research **large functional** proteins
 - DNA polymerase

Acknowledge

- Thanks for Prof. **Yunfei Chen**'s guidance with patience
- Thanks the **Big Data Computing Center of Southeast University** for providing the facility support on the numerical calculations in this paper
- Thanks for the sponsorship of **National Natural Science Foundation of China**



Questions & Answers