

SUPPORTING INFORMATION

for

OpenMSCG: A Software Tool for Bottom-up Coarse-graining

Yuxing Peng,¹ Alexander J. Pak,² Aleksander E. P. Durumeric,³ Patrick G. Sahrman,⁴
Sriramvignesh Mani,⁴ Jaehyeok Jin,⁴ Timothy D. Loose,⁴ Jeriann Beiter,⁴ and Gregory A. Voth^{*4, a)}

¹ NVIDIA Corporation, 2788 San Tomas Expressway, Santa Clara, CA 95051, USA

²Department of Chemical and Biological Engineering, Colorado School of Mines, Golden,
Colorado 80401, USA

³Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany

⁴Department of Chemistry, Chicago Center for Theoretical Chemistry, James Franck Institute,
and Institute for Biophysical Dynamics, The University of Chicago, Chicago, Illinois 60637,
USA

^{a)} Author to whom correspondence should be addressed: gavoth@uchicago.edu

Dynamic Programming Algorithm for Essential Dynamics Coarse-Graining (EDCG)

The goal of EDCG is to minimize the residual from mapping n C_α atoms of a protein into N CG sites, which is defined as

$$\chi = \frac{1}{3N} \sum_{I=1}^N \mathbf{C}_I$$

where

$$\mathbf{C}_I = \sum_{i \in I} \sum_{j \geq i \in I} \langle |\Delta \mathbf{r}_i^{ED} - \Delta \mathbf{r}_j^{ED}|^2 \rangle$$

is the loss of total covariance from mapping a group of C_α atoms to the CG site I . In a sequential/linear model, C_α atoms associated with each CG site are assumed contiguous in the protein primary sequence, and a loss function for mapping a group of consecutive C_α atoms $\{a, a+1, a+2 \dots b\}$ into a CG site can be rewritten as

$$\mathbf{C}(a, b) = \sum_{i=a}^{b-1} \sum_{j=i+1}^b \langle |\Delta \mathbf{r}_i^{ED} - \Delta \mathbf{r}_j^{ED}|^2 \rangle$$

Therefore, to calculate the loss function for all possible values of $1 \leq a < b \leq n$, the time complexity is $O(N^4)$. To gain a higher computational efficiency, the loss function can be calculated from the recurrence equation:

$$\mathbf{C}(a, b) = \begin{cases} \mathbf{C}(a, b - 1) + \mathbf{C}(a + 1, b) - \mathbf{C}(a + 1, b - 1) + \langle |\Delta \mathbf{r}_a^{ED} - \Delta \mathbf{r}_b^{ED}|^2 \rangle, & a < b \\ 0, & a = b \end{cases}$$

The pseudo code with a time complexity of $O(N^2)$ can be designed as

```

loop L from 1 to n # number of Cα atoms to be mapped
  loop a from 1 to n-L # starting Cα atom
    b = a + L - 1 # ending Cα atom
    Calculate C(a,b) # calculate the loss
  
```

After obtaining all values of the loss functions, we can then define a sub-residual, $\chi(k, m)$, which is the minimum of total loss by mapping the first k C_α atoms into m CG sites, where $1 \leq k \leq n$, and $1 \leq m \leq N$ and $m \leq k$. A recurrence equation can then be defined as

$$\chi(k, m) = \begin{cases} \min_{1 \leq i \leq k-m+1} \{ \chi(k-i, m-1) + C(k-i+1, k) \}, & m > 1 \\ C(1, k), & m = 1 \end{cases}$$

The pseudo code can be designed as

```

Chi(k,1) = C(1,k)
loop m from 1 to N
  loop k from 1 to n
    Chi(k,m) = min{Chi(k-i,m-1)+C(k-i+1,k)}
  
```

The time complexity is $O(N^3)$ for calculating all sub-loss functions and the final residual $\chi(n, N)$ is the global minimum. This algorithm, known as dynamic programming, can be used to obtain the globally optimized solution.