

Simulate Time-integrated Coarse-grained Molecular Dynamics with Geometric Machine Learning

Xiang Fu^{1†}, Tian Xie^{1†}, Nathan J. Rebello², Bradley D. Olsen² and Tommi Jaakkola^{1*}

¹CSAIL, MIT, Cambridge, 02139, MA, United States.
²Department of Chemical Engineering, MIT, Cambridge, 02139, MA, United States.

*Corresponding author(s). E-mail(s): xiangfu@csail.mit.edu; tommi@csail.mit.edu;
 Contributing authors: txie@csail.mit.edu; nrebello@mit.edu; bdolesen@mit.edu;
[†]Equal Contribution.



Tommi Jaakkola
 MIT
 in csail.mit.edu 的电子邮件经过验证 - 主页
 machine learning natural language processing biomolecular design

引用次数 年份

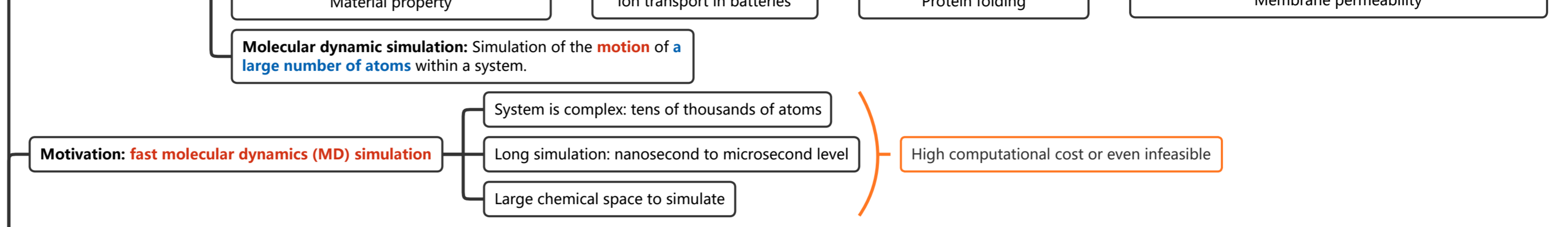
About the last author

标题	引用次数	年份
An introduction to variational methods for graphical models	4008	1999
Exploring generative models in discriminative classifiers	1960	1998
Maximum-margin matrix factorization	1507	2004
Convergence of stochastic iterative dynamic programming algorithms	1114	1993
Weighted low-rank approximations	953	2003



Molecular dynamic
 Material property, Ion transport in batteries, Protein folding, Membrane permeability

Motivation: fast molecular dynamics (MD) simulation
 System is complex: tens of thousands of atoms
 Long simulation: nanosecond to microsecond level
 Large chemical space to simulate
 High computational cost or even infeasible



Ab initio accurate: Reactions, Biopolymers, Molecules, Nucleotides, Proteins, Peptides, Materials, Biological Units

Machine learning force field method
 Ab initio method: 一种基于量子化学的方法。通过求解与系统所有原子相关的薛定谔方程得到精确的原子间相互作用力。进而模拟系统内所有原子的运动。但计算成本高昂。
 Force fields method: 一种基于经验力场的方法。原子之间的相互作用由预先定义的势能函数描述。计算成本低，但不能精确描述复杂的化学过程。
 Machine learning force field: 一种基于机器学习的方法。通过训练机器学习模型，输入原子的坐标和类型，输出原子之间的相互作用力。精度介于前两者之间。

Coarse-graining model
 Reduce system complexity

Enhanced sampling method
 Conformational State
 Modify the potential energy surface to enable faster sampling of transition between metastable states
 Shortage: No dynamics

Machine learning generative modeling
 Shortage: No dynamics

Solutions to the above shortage
 System is complex: tens of thousands of atoms
 Long simulation: nanosecond to microsecond level
 Large chemical space to simulate
 Spatial coarse-graining
 Very large time integration step at nanosecond-level without force computation
 Learn from short MD trajectories and simulate long trajectories for novel unseen systems

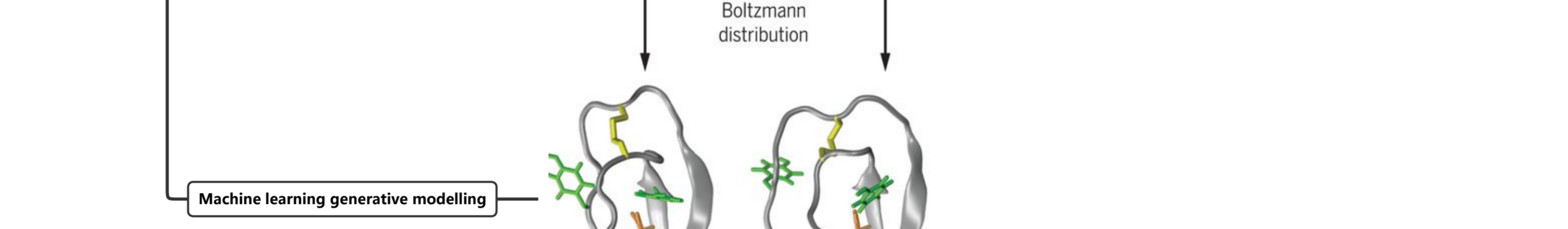
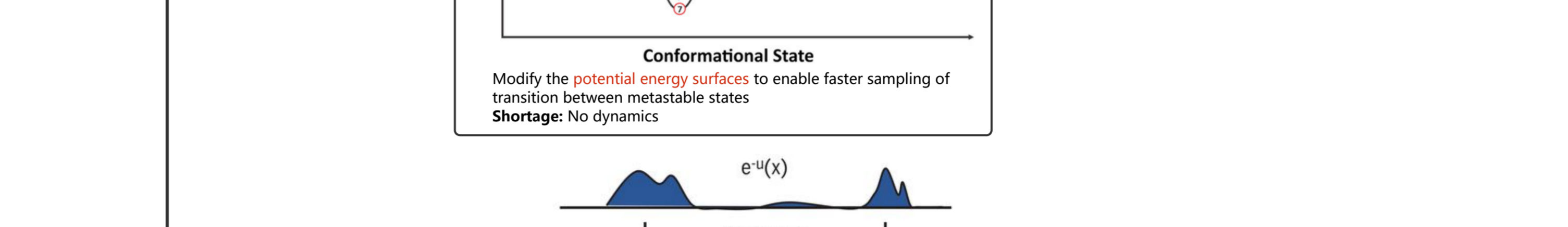


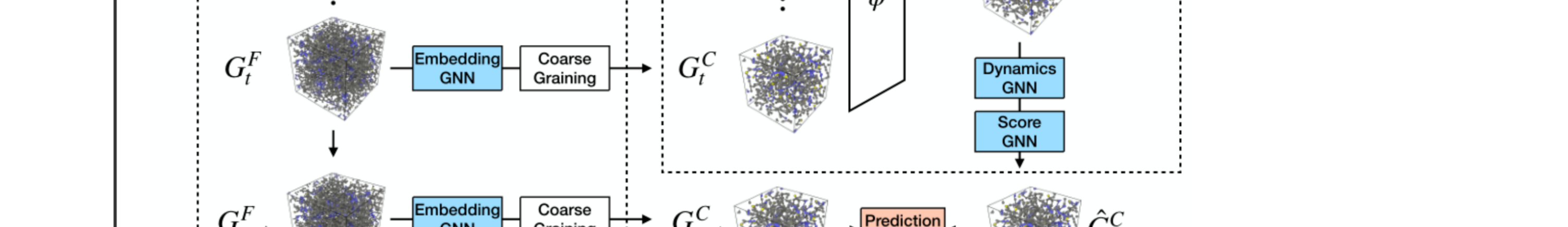
Figure 2 (a) The GNS predicts finite states represented as particles using its learned-dynamics model, A_θ , and a fixed update procedure. (b) The A_θ uses an "encode-process-decode" scheme, which compares dynamics information, Y , from input state, X . (c) The ENCODER constructs latent graphs, G^E , from the input state, X . (d) The PROCESSOR performs M rounds of learned message-passing over the latent graphs, G^E, \dots, G^M . (e) The DECODER extracts dynamics information, Y' , from the final latent graph, G^M .

ICML 20-Learning to simulate complex physics with graph networks

Graph processing with graph neural networks
 Objective: Learn a graph from give X
 $Y = G^M$
 Encoder: $V_i = \sigma^U(X_i)$
 E^i is a 2 layers MLP that is independently applied to each node
 $R_{ij} = \sigma^V(R_{ij})$
 E^i is a 2 layers MLP that is independently applied to each edge
 $G = (G^1, \dots, G^M)$
 Processor: Objective: Generate a sequence of update latent graph
 $G^{m+1} = \text{GNN}^{m+1}(G^m)$
 GNN is a 7 layers MPNN
 Decoder: $G \rightarrow Y'$
 $Y' = \sigma^V(Y)$
 Y' is a 2 layers MLP that is independently applied to each node

Representing MD trajectories as time series of graphs
 Time variant graph: MD simulation trajectory is represented as a time series of fine-level graphs (G^t). Each fine-level graph is represented as $v_{i,j}^t \in V^t$
 atom invariant and variant feature
 $e_{i,j}^t \in E^t$
 chemical bond

Learning CG-bead type embeddings with embedding GNN
 Objective: Learn the node embeddings which contains no positional information
 Input: Time invariant graph that only contains atom type, atom weights, bond types information
 Output: node type embedding e_i^t for all $v_i^t \in V^t$
 $v_i^t = [a_i^t, w_i^t]$
 learnable atom type embedding; vector weight; scale
 $e_i^t = [a_i^t + a_j^t + a_{ij}^t]$
 learnable bond type embedding; vector



Coarse-graining
 Graph clustering CG model: METIS
 Objectives: grouping atoms in the same group into a CG bead
 Fine-level node feature: $v_{i,j}^t = [x_i^t, w_{ij}^t, x_j^t]$
 Create CG-bond: $v_{i,j}^t \in E^t \iff \exists i \in C_m, j \in C_n$
 such that $e_{i,j}^t \in E^t$
 Create CG-bond if there exists a chemical bond between a pair of atoms in group C_m and group C_n
 Create radius cut-off edge
 $v_{i,j}^t = [x_i^t, w_{ij}^t, x_j^t, [e_{i,j}^t - x_i^t, e_{i,j}^t - x_j^t]]$
 Node feature
 Input: Edge feature, $e_{i,j}^t$ indicate whether $e_{i,j}^t$ is a CG-bond or is constructed through radius cut-off

Learning CG MD with Dynamics GNN
 Dynamics GNN: Encoder: MPNN
 Decoder: MLP on Node to learn the acceleration
 $\text{GNN}(\text{node}(G^E, \Delta t, \tau)) = \mathcal{N}(\mu, \sigma^2) = [\mathcal{N}(\mu, \sigma^2), \mathcal{N}(\sigma, \sigma^2)] \in V^t$
 Finally return 3-dimensional Gaussian distribution of acceleration for each coarse-grained node
 Loss: $L_{\text{Dynamics}} = -\log \mathcal{N}(k_i^t | \mu_i, \sigma_i^2)$
 Predict position
 $\hat{x}_{i,t+\Delta t} = x_i + \hat{a}_i \Delta t$
 $\hat{x}_{i,t+\Delta t} = x_i + \hat{a}_i \Delta t$
 Predict the position using Euler integration

Learning to refine CGMD predictions with Score GNN
 Noise Conditional Score Network (NCSN): NeUP19 - Generative Modeling by Estimating Gradients of the Data Distribution
 Goal of NCSN here: Refine the results from Dynamics GNN
 Langevin dynamics: $\dot{x}_i = x_{i-1} + \frac{1}{2} \nabla_x \log p(x_{i-1}) + \sqrt{2} z_i$
 Goal of NCSN: $\mathcal{R}_\theta(x) \approx \nabla_x \log p_{\text{data}}(x)$
 $\mathcal{R}_\theta(x) \approx \frac{1}{2} \nabla_{\text{pos}} \log p(x) - \nabla_x \log p_{\text{data}}(x)$
 $\min_{\theta} \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} \|\mathcal{R}_\theta(x) - \nabla_x \log q_\theta(x)\|_2^2$
 $\mathcal{R}_\theta(x; \sigma) = \frac{\nabla_x \log q_\theta(x) - \nabla_x \log p_{\text{data}}(x)}{\sigma}$
 $\mathcal{L}(\theta; \sigma) \triangleq \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} \|\mathcal{R}_\theta(x; \sigma)\|_2^2$

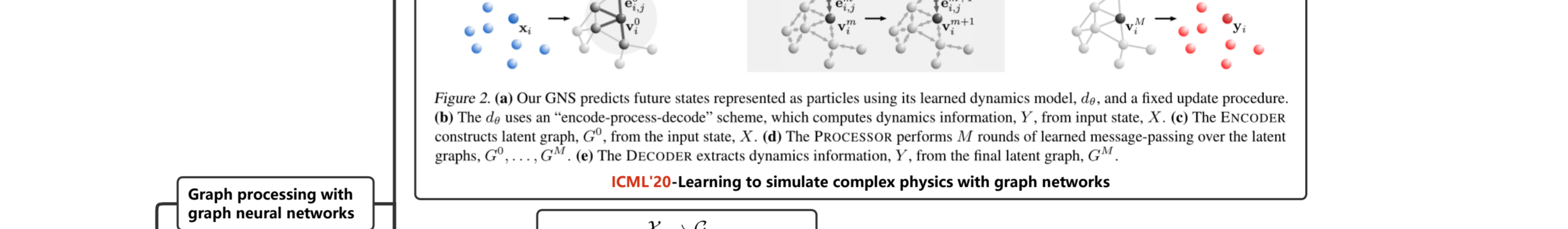


Fig. 2 (a) Class-I polymers are generalizable for training, while class-II polymers are used for testing. The structure variation requires the model to learn generalizable dynamics. (b) Training trajectories are 50k long (black dashed line), while we use 5M long trajectories for evaluation. (c) Short ground truth MD of training trajectory length gives high variance, poor estimation of $\langle R_g^2 \rangle$. (d) Example polymer before and after coarse-graining. The green beads and bonds represent the CG structure. (e, f, g) $\langle R_g^2 \rangle$ estimation performance of the supervised learning baseline using (e) GNN, (f) LSTM, and (g) our learned simulator.

(a) 训练数据(Training data) - 聚合物的 MD 模拟数据
 (b) 测试数据(Testing data) - 聚合物的 MD 模拟数据
 (c) 聚合物的 原子坐标和键长信息。这与聚合物在溶液中的流变性和扩散性有关。
 (d) 对于聚合物的粗粒化建模。
 (e) 粗粒化后的聚合物的 MD 模拟数据。
 (f) 粗粒化后的聚合物的 MD 模拟数据。
 (g) 粗粒化后的聚合物的 MD 模拟数据。

Figure 3 (a) R_g^2 distribution computed from our learned simulation matches the ground truth. The vertical dashed lines associate the mean of the distribution. (b) Two-dimensional free energy surface profiled with PCA, with representative time states with low/high free energy revealed. (c) The autocorrelation function of R_g^2 for the three polymers with the smallest, median, and the largest relaxation time. (d) Predicted performance of our model on the R_g^2 relaxation time. (e) Computational cost comparison of our learned simulation vs. traditional MD. The units are in log scale. (f) Performance of all models in predicting diffusivity of Li-ion, TFSI-ion, and polymer particles.

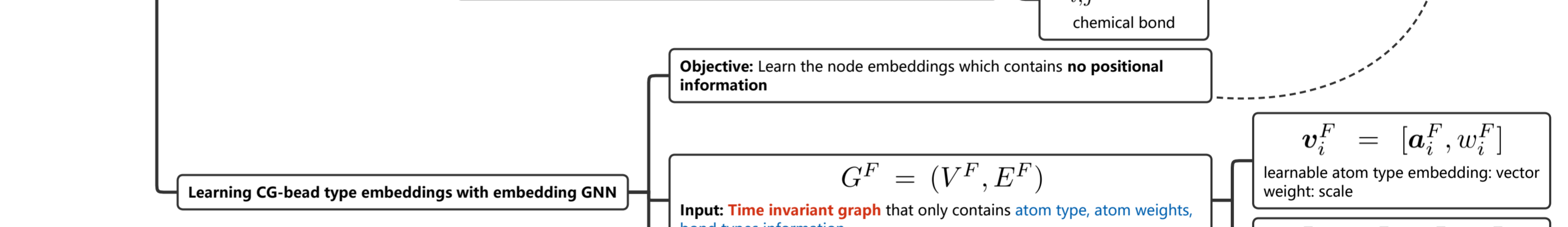


Figure 4 (a) The chemical space for the learned SPEs. (b) Example SPE before and after coarse-graining. (c) Bead collision as a function of simulation time, averaged over the 50 training SPEs for all methods. (d) RDF of Li-ions, for our model with/without the Score GNN refinement, and the ground truth MD simulation, averaged over a 50ks in trajectories of a selected SPE. (e) RDF of Li-ions and TFSI-ions. (f) RDF of Li-ions and polymer particles.

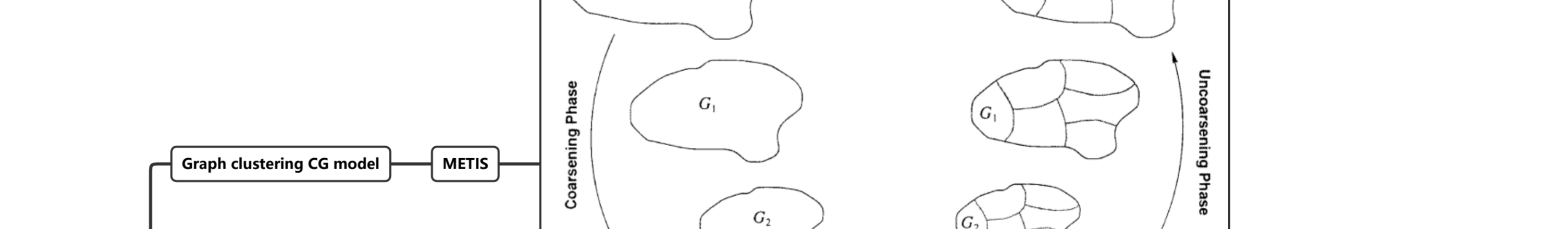


Figure 5 (a) Slow convergence of Li-ion diffusivity. Examining this long-time property with only 5M training trajectories (black dashed line) is very challenging. (b) 5M MD gives a poor estimation of Li-ion diffusivity. (c, d, e) Performance of SL GNN model. Learned simulator without Score GNN refinement, and our full model in Li-ion diffusivity prediction. Only our full model is able to predict long-time property from short training trajectories only. (f) Computational cost of our learned simulator vs. traditional MD. The units are in log scale. (g) Performance of all models in predicting diffusivity of Li-ion, TFSI-ion, and polymer particles.

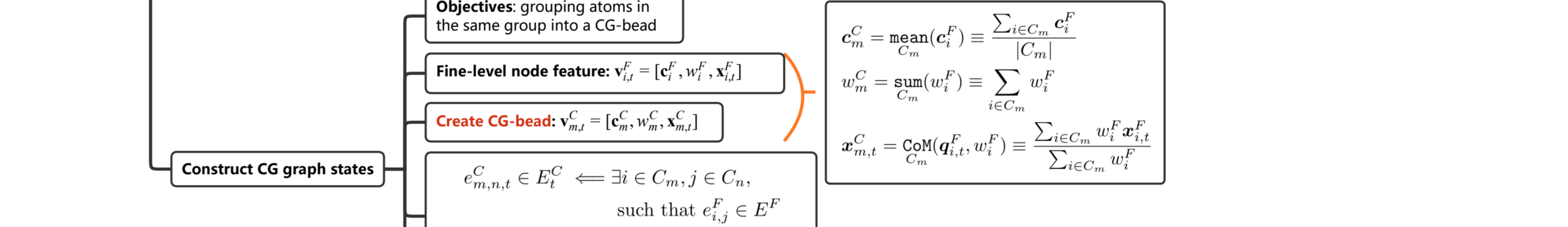


Figure 6 (a) Li-ion, Li-ion. (b) Li-ion, TFSI-ion. (c) Li-ion, Polymer. (d) Li-ion, Li-ion. (e) Li-ion, TFSI-ion. (f) Li-ion, Polymer.

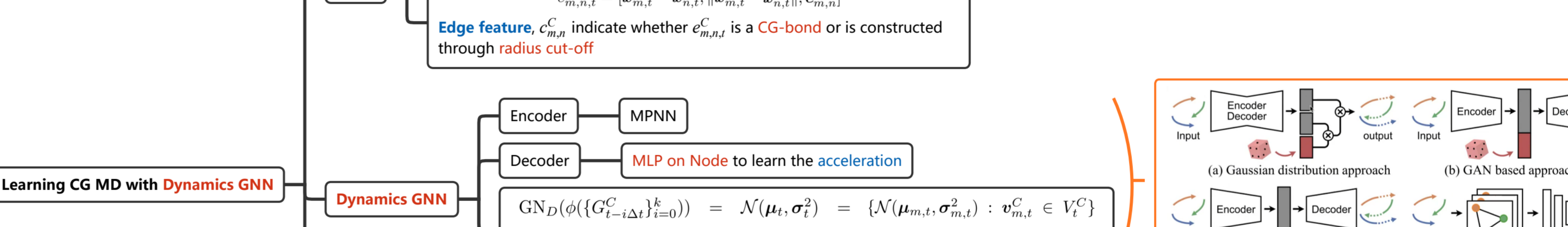


Figure 7 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

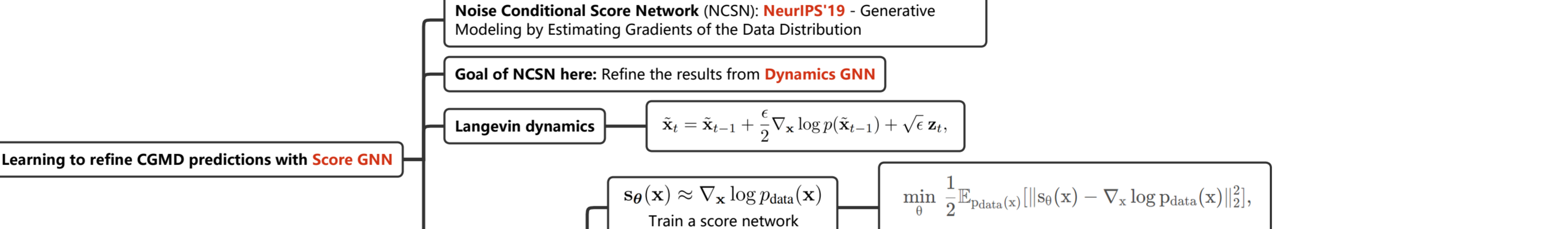


Figure 8 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

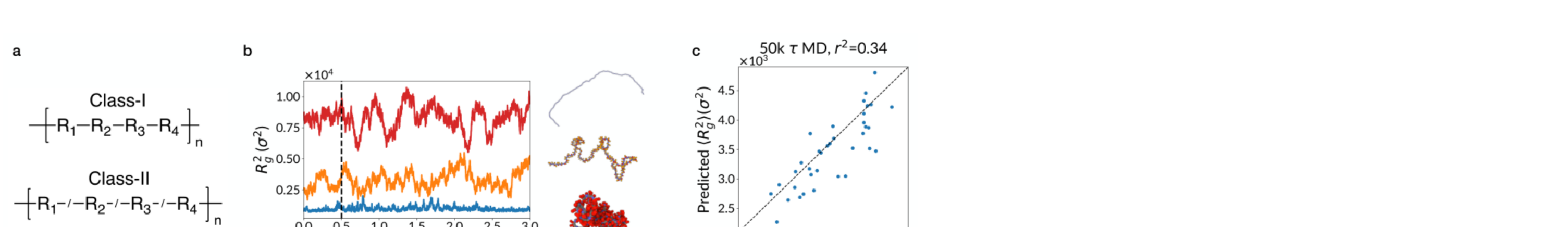


Figure 9 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

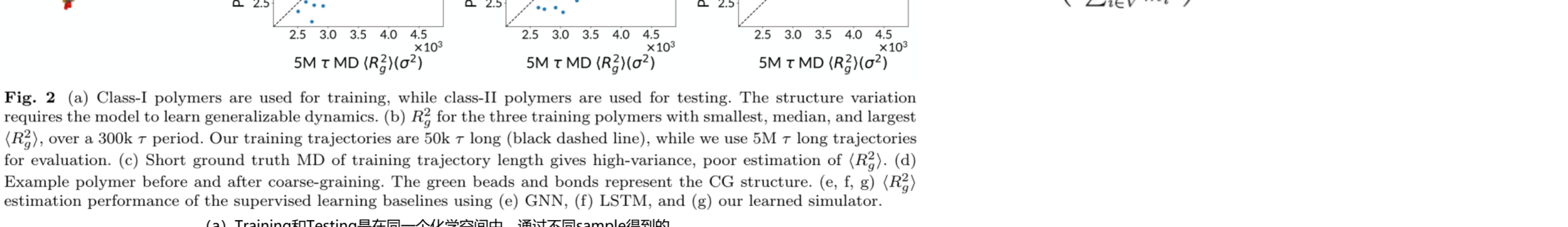


Figure 10 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

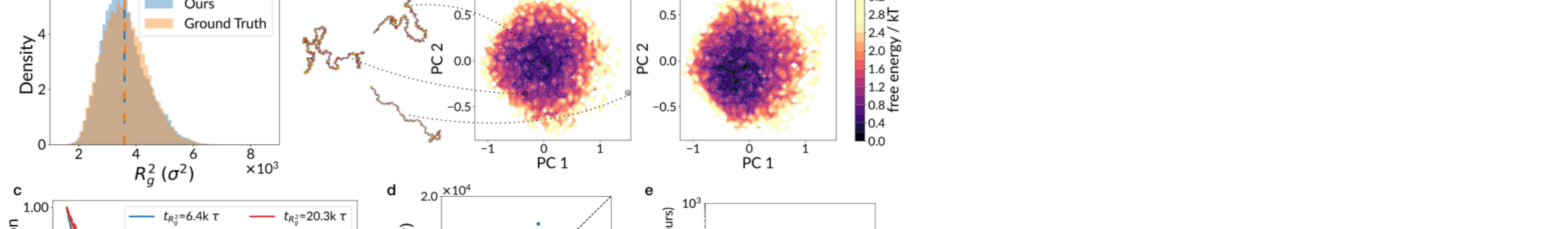


Figure 11 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

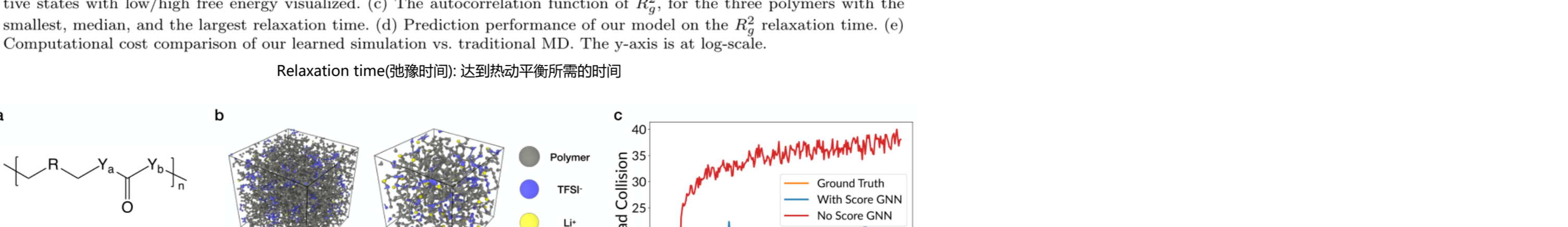


Figure 12 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

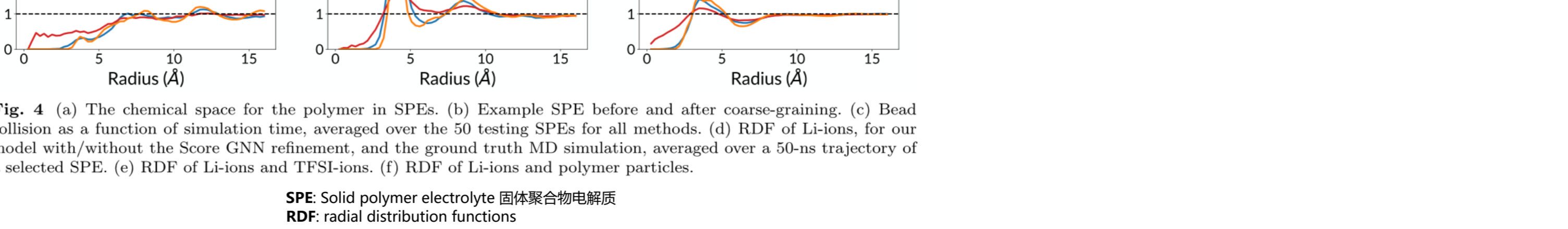


Figure 13 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.

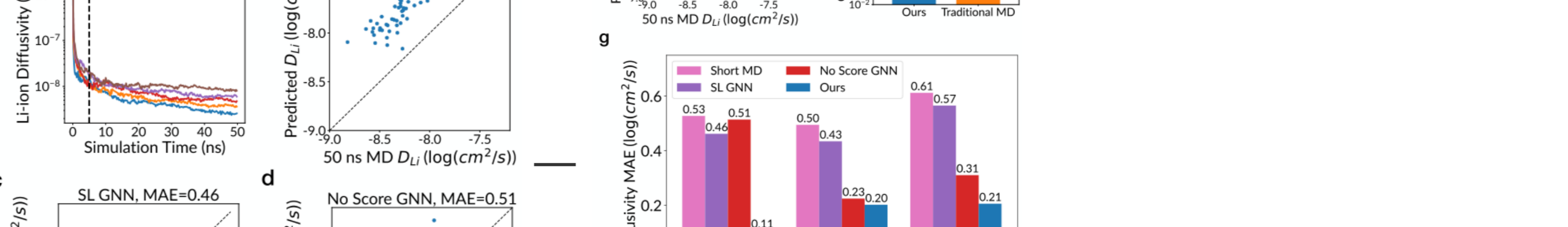


Figure 14 (a) Li-ion Diffusivity. (b) 5M MD, MAE=0.52. (c) SL GNN, MAE=0.46. (d) No Score GNN, MAE=0.55.