

# molecule-cell-gen

## Installation:

- Follow the instructions in [TopoModelX](#)

## GDSS Overview:

### Graph Diffusion Process

A graph  $G$  with  $N$  nodes is defined by its node features  $\mathbf{X} \in \mathbb{R}^{N \times F}$  and the weighted adjacency matrix  $\mathbf{A} \in \mathbb{R}^{N \times N}$  as  $\mathbf{G} = (\mathbf{X}_0, \mathbf{A}_0) \in \mathbb{R}^{N \times F} \times \mathbb{R}^{N \times N} := \mathcal{G}$ , where  $F$  is the dimension of the node features. The diffusion process  $\{\mathbf{G}_t = (\mathbf{X}_t, \mathbf{A}_t)\}_{t \in [0, T]}$  transforms the node features and the adjacency matrices to a simple noise distribution. We can model the diffusion process by the following Ito SDE

$$d\mathbf{G}_t = \mathbf{f}_t(\mathbf{G}_t)dt + \mathbf{g}_t(\mathbf{G}_t)d\mathbf{w}, \quad \mathbf{G}_0 \sim p_{\text{data}}, \quad (1)$$

where  $\mathbf{f}_t(\cdot) : \mathcal{G} \rightarrow \mathcal{G}^1$  is the linear drift coefficient,  $\mathbf{g}_t(\cdot) : \mathcal{G} \rightarrow \mathcal{G} \times \mathcal{G}$  is the diffusion coefficient, and  $\mathbf{w}$  is the standard Wiener process. The reverse of the diffusion process in time is also a diffusion process described by the SDE:

$$d\mathbf{G}_t = [\mathbf{f}_t(\mathbf{G}_t) - g_t^2 \nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t)]d\bar{t} + g_t d\bar{\mathbf{w}} \quad (2)$$

where  $p_t$  denotes the marginal distribution under the forward diffusion process at time  $t$ ,  $\bar{\mathbf{w}}$  is a reverse-time standard Wiener process, and  $d\bar{t}$  is an infinitesimal negative time step.

Solving **Eq. (2)** requires the estimation of  $\nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t) \in \mathbb{R}^{N \times F} \times \mathbb{R}^{N \times N}$  which is expensive to compute. So, the paper proposes a new reverse-time diffusion process.

$$\begin{cases} d\mathbf{X}_t = [\mathbf{f}_{1,t}(\mathbf{X}_t) - g_{1,t}^2 \nabla_{\mathbf{X}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)]d\bar{t} + g_{1,t} d\bar{\mathbf{w}}_1 \\ d\mathbf{A}_t = [\mathbf{f}_{2,t}(\mathbf{A}_t) - g_{2,t}^2 \nabla_{\mathbf{A}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)]d\bar{t} + g_{2,t} d\bar{\mathbf{w}}_2 \end{cases} \quad (3)$$

where  $\mathbf{f}_{1,t}$  and  $\mathbf{f}_{2,t}$  are linear drift coefficients satisfying  $\mathbf{f}_t(\mathbf{X}, \mathbf{A}) = (\mathbf{f}_{1,t}(\mathbf{X}), \mathbf{f}_{2,t}(\mathbf{A}))$ ,  $g_{1,t}$  and  $g_{2,t}$  are scalar diffusion coefficients, and  $\bar{\mathbf{w}}_1, \bar{\mathbf{w}}_2$  are reverse-time standard Wiener Process.

## Estimating the Partial Score Functions

### Training Objectives

Note that the diffusion processes of  $\mathbf{X}_0, \mathbf{A}_0$  in the system are dependent on each other, related by the gradients of the joint log-density  $\nabla_{\mathbf{X}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)$  and  $\nabla_{\mathbf{A}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)$ , which are the **partial score functions**. The partial score functions can be estimated by training the time-dependent score-based models  $\mathbf{s}_{\theta,t}$  and  $\mathbf{s}_{\phi,t}$ , so that  $\mathbf{s}_{\theta,t}(\mathbf{G}_t) \approx \nabla_{\mathbf{X}_t} \log p_t(\mathbf{G}_t)$  and  $\mathbf{s}_{\phi,t}$ , so that  $\mathbf{s}_{\phi,t}(\mathbf{G}_t) \approx \nabla_{\mathbf{A}_t} \log p_t(\mathbf{G}_t)$ .

The paper proposes a novel training objective

$$\begin{aligned} & \min_{\theta} \mathbb{E}_t \left\{ \lambda_1(t) \mathbb{E}_{\mathbf{G}_0} \mathbb{E}_{\mathbf{G}_t | \mathbf{G}_0} \| \mathbf{s}_{\theta,t}(\mathbf{G}_t) - \nabla_{\mathbf{X}_t} \log p_{0t}(\mathbf{X}_t | \mathbf{X}_0) \|_2^2 \right\} \\ & \min_{\phi} \mathbb{E}_t \left\{ \lambda_2(t) \mathbb{E}_{\mathbf{G}_0} \mathbb{E}_{\mathbf{G}_t | \mathbf{G}_0} \| \mathbf{s}_{\phi,t}(\mathbf{G}_t) - \nabla_{\mathbf{A}_t} \log p_{0t}(\mathbf{A}_t | \mathbf{A}_0) \|_2^2 \right\} \end{aligned} \quad (4)$$

where  $\lambda_1(t)$  and  $\lambda_2(t)$  are positive weighting functions and  $t$  is uniformly sampled from  $[0, T]$ . The expectations in Eq. (4) can be efficiently computed using the Monte Carlo estimate with the samples  $(t, \mathbf{G}_0, \mathbf{G}_t)$ .

Then, the paper proposes new architectures for the time-dependent score-based models that can capture the dependencies of  $\mathbf{X}_t$  and  $\mathbf{A}_t$  through time, based on GNNs. The score-based model  $\mathbf{s}_{\phi,t}$  used to estimate  $\nabla_{\mathbf{A}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)$  is defined by

#### Permutation-Equivariant Score-based Model

$$\mathbf{s}_{\phi,t}(\mathbf{G}_t) = \text{MLP} \left( \left[ \{ \text{GMH}(\mathbf{H}_i, \mathbf{A}_t^p) \}_{i=0, p=1}^{K, P} \right] \right), \quad (5)$$

where  $\mathbf{A}_t^p$  are the higher-order adjacency matrices,  $\mathbf{H}_{i+1} = \text{GNN}(\mathbf{H}_i, \mathbf{A}_t)$  with  $\mathbf{H}_0 = \mathbf{X}_t$  given,  $[\cdot]$  denotes the concatenation operation, GMH denotes the graph multi-head attention block, and  $K$  denotes the number of GMH layers. And the score-based model  $\mathbf{s}_{\theta,t}$  to estimate  $\nabla_{\mathbf{X}_t} \log p_t(\mathbf{X}_t, \mathbf{A}_t)$  is defined as

$$\mathbf{s}_{\theta,t}(\mathbf{G}_t) = \text{MLP} \left( \left[ \{ \mathbf{H}_i \}_{i=0}^L \right] \right), \quad (6)$$

where  $\mathbf{H}_{i+1} = \text{GNN}(\mathbf{H}_i, \mathbf{A}_t)$  with  $\mathbf{H}_0 = \mathbf{X}_t$  given and  $L$  denotes the number of GNN layers. Since the message-passing operations of GNNs and the attention function used in GMH are permutation-equivariant, the score-based models are also equivariant, and the log-likelihood implicitly defined by the models is guaranteed to be permutation-invariant.

#### Solving the system of Reverse-time SDEs

Then, we simulate the system of reverse o solve the system of reverse-time SDEs

## Main tasks (order by priority)

**Intitution:** Improve the GDSS model by considering molecules as cellular complexes, adding more atomic features and ring features. We want to know if topological deep learning can improve traditional graph representation learning.

- Update code to lift molecular graph  $G : (\mathbf{X}_0, \mathbf{A}_0)$  to cellular complexes implemented by TopomodelX
  - Find the right place to conduct cellular lifting of the *GDSS* code repo.
  - Once we have a NetworkX graph, we can easily convert it to a CC object
  - Extract rings from the graph and add it to the CC object
    - Current Solution: use the graph-tool package to extract rings and add them to CC objects
    - Difficulty: the graph-tool package is not compatible with TMX now; this is a potential problem during training
- Design two new time-dependent score-based models based on cell attention layers to estimate  $\nabla_{\mathbf{X}_t} \log p_{0t}(\mathbf{G}_t)$  and  $\nabla_{\mathbf{A}_t} \log p_{0t}(\mathbf{G}_t)$ 
  - The cell attention network takes  $\mathbf{X}_0, \mathbf{X}_1, \mathbf{A}_0, \mathbf{A}_{\uparrow,1}$  and  $\mathbf{A}_{\downarrow,1}$ , which can be extracted from the CC object converted from the molecular graph
  - The *CAN* implemented in TopoModelX is designed for classification; we need to slightly update the model structure, including adding an MLP, to estimate a matrix with the same shapes  $\mathbf{X}_0$  and  $\mathbf{A}_0$

- Validate our new cell attention layer with Prof. Derr
  - Implement the new score functions in *ScoreNetwork\_A.py* and *ScoreNetwork\_X.py*
- Adding more atomic features to  $X_0$  before doing the diffusion process
  - Currently, the atom feature is just an array including the atomic numbers, e.g C:6
  - Discuss with Lance to select the feature candidates derived from SMILES
  - Embed the feature extraction method in *smile\_to\_graph.py*
- Since rank-2 cells (rings) are added, we should consider assigning attributes to them with the help of biological knowledge. (Optional)
  - With this new feature, we have to update the score function again. Thus, we should consider this as a later task