

Molecular Pre-trained Models

Wu Hualong

HARBIN INSTITUTE OF TECHNOLOGY, SHEN ZHEN

Problem: 3D geometric information is scarce in downstream tasks due to its expensive cost to obtain.

Goal: Propose a Molecular Pre-trained model exploiting both 2D and 3D information of molecules.

Requirement:

- The model should be pre-trained with the knowledge of 3D geometry and then fine-tuned on downstream tasks with or without 3D information.
- Consider the symmetry of molecules.

Dataset

Datasets:

- Pre-training
 - GEOM, a dataset with over 37 million molecular conformations annotated by energy and statistical weight for over 450,000 molecules.
 - PCQM4Mv2, including 3.37 million molecules with both 2D graphs and 3D geometric structures.(3D available only for training molecules)
- Downstream
 - MoleculeNet:one of the most widely used benchmarks for 2D molecular property prediction
 - 2D, 8 classification tasks
 - PCQM4Mv2 humo-lumo gap (2D)
 - QM9 quantum properties
 - 3D, 12 regression tasks

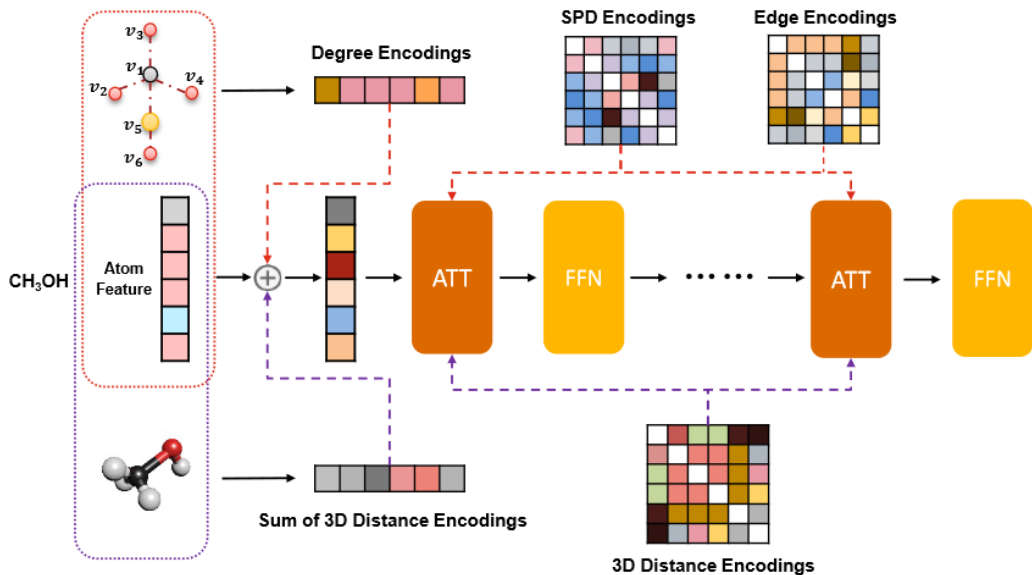
Pre-training Strategies

- AutoEncoding (AE)
- Masked Component Modeling
- Contrastive Learning (CL)
 - Same-Scale Contrast (SSC)
- DeNoising (DN)

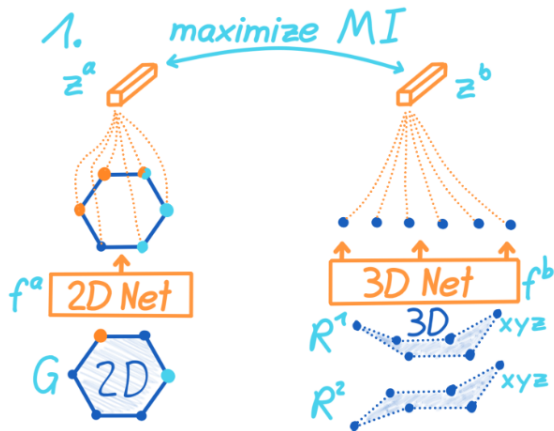
Model	Input	Backbone architecture	Pre-training task	Pre-training database
GraphCL	Graph	5-layer GIN	SSC	ZINC15 (2M)
Graphormer	Graph	Graphormer	Supervised	PCQM4M-LSC (3.8M)
Mole-BERT	Graph	5-layer GIN	MCM + SSC	GEOM (100K)
GraphMVP	Graph+Geometry	5-layer GIN + SchNet	SSC + AE	GEOM (50K)
3D Informax	Graph+Geometry	PNA	SSC	QM9 (50K) + GEOM (140K) + QMugs (620K)
Transformer-M	Graph+Geometry	Transformer-M	Supervised + DN	PCQM4Mv2
Denosing	Geometry	GNS	DN	PCQM4Mv2 (3.4M)

Table 1: A summary of Molecular Pre-trained Models

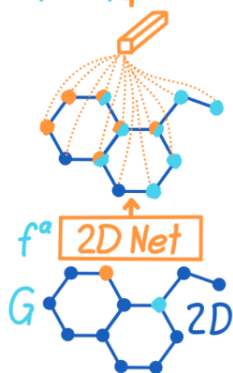
Transformer-M

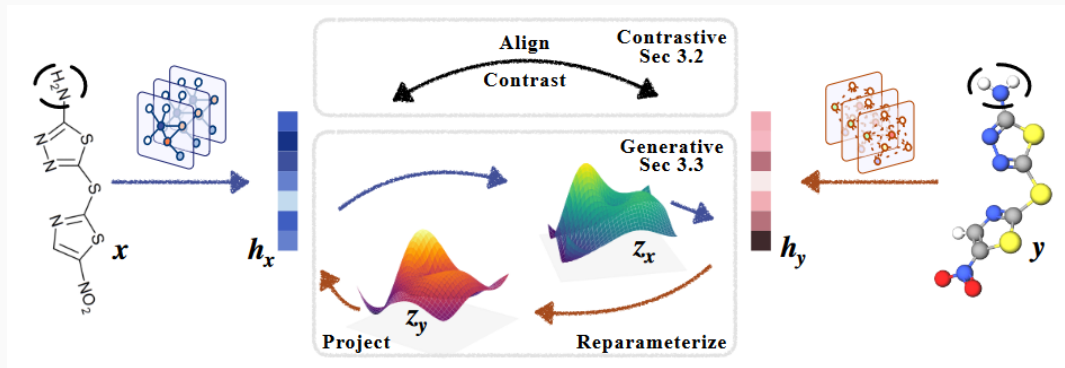


3D Informax

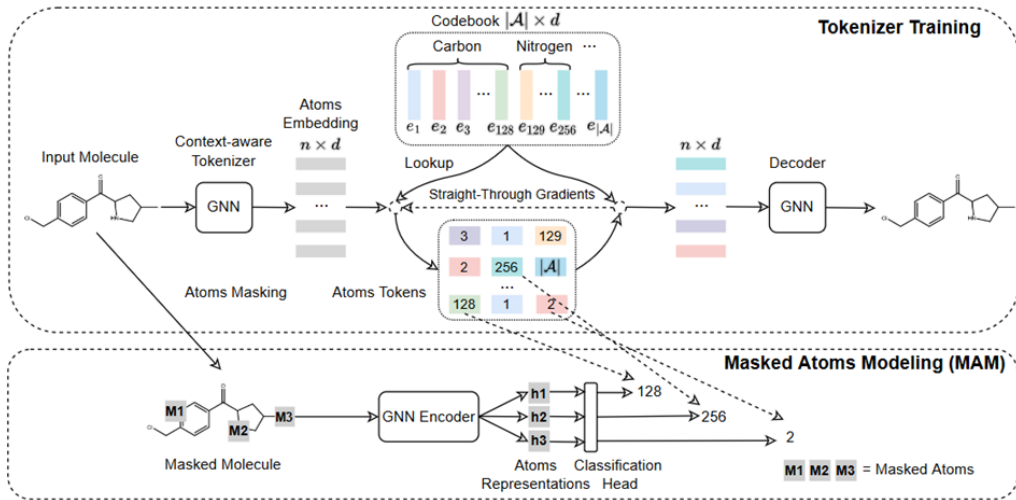


2. *property*

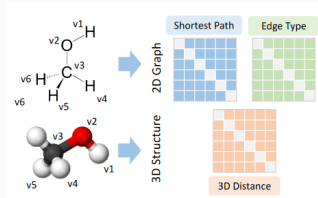




Mole-BERT



Initial idea



We have two channels of input:

$$\begin{aligned} X_{2D}^0 &= X + \psi^{Degree} \\ X_{3D}^0 &= X + \psi^{Sum \text{ of Distance}} \end{aligned} \tag{1}$$

where ψ^{Degree} denote the degree encoding of the atom, $\psi_i^{Sum \text{ of 3D Distance}} = \sum_{j \in [n]} \psi_{(i,j)} W_D^3$

Then in self-attention:

$$A_{2D}(X) = \text{softmax} \left(\frac{X_{2D} W_Q (X_{2D} W_K)^{\top}}{\sqrt{d}} + \underbrace{\phi^{\text{SPD}} + \phi^{\text{Edge}}}_{2D \text{ pair-wise channel}} \right) \quad (2)$$

$$A_{3D}(X) = \text{softmax} \left(\frac{X_{3D} W_Q (X_{3D} W_K)^{\top}}{\sqrt{d}} + \underbrace{\phi^{\text{3D Distance}}}_{3D \text{ pair-wise channel}} \right) \quad (3)$$

Then we have two types of representation z_{3D}, z_{2D} ,

$$\begin{aligned} z_{2D} &= A_{2D} V_{2D} \\ z_{3D} &= A_{3D} V_{3D} \end{aligned} \quad (4)$$

these can be atom-wise or molecular-wise. And we can design pre-training task via these representations.