

# Self-Supervised Graph Transformer on Large-Scale Molecular Data

Yu Rong\*, Yatao Bian\*, Tingyang Xu,  
Weiyang Xie, Ying Wei, Wenbing Huang,  
Junzhou Huang

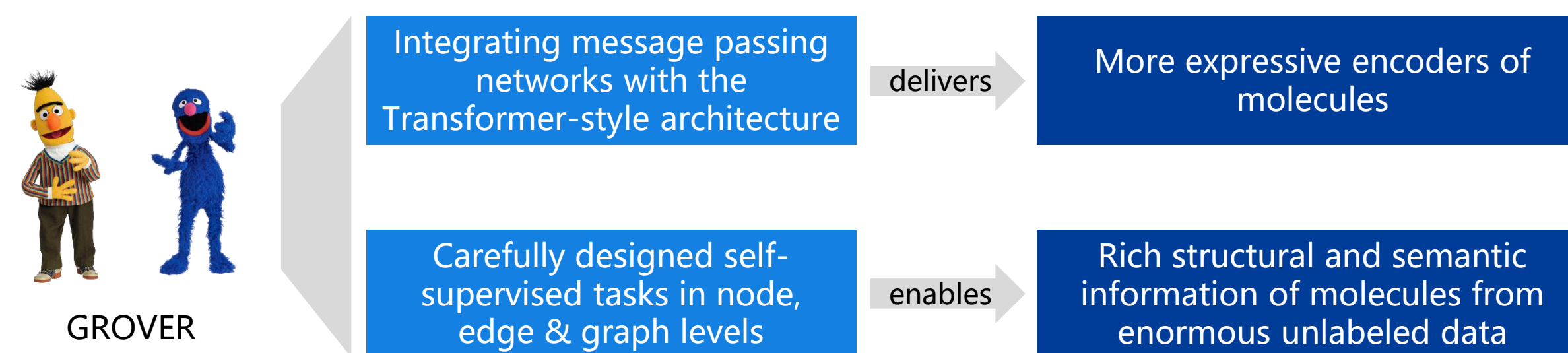
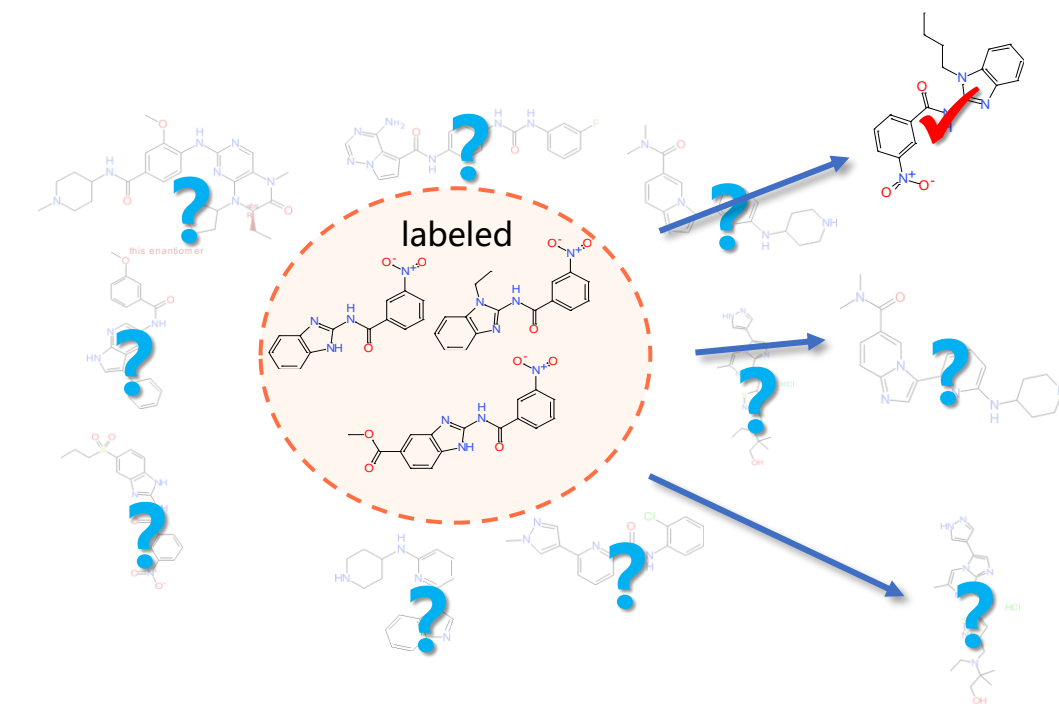


## Background & Contributions

GNNs are widely adopted for molecular tasks.

### Challenges:

- ❑ Insufficient labeled molecules for supervised training
- ❑ Poor generalization capabilities to newly-synthesized molecules

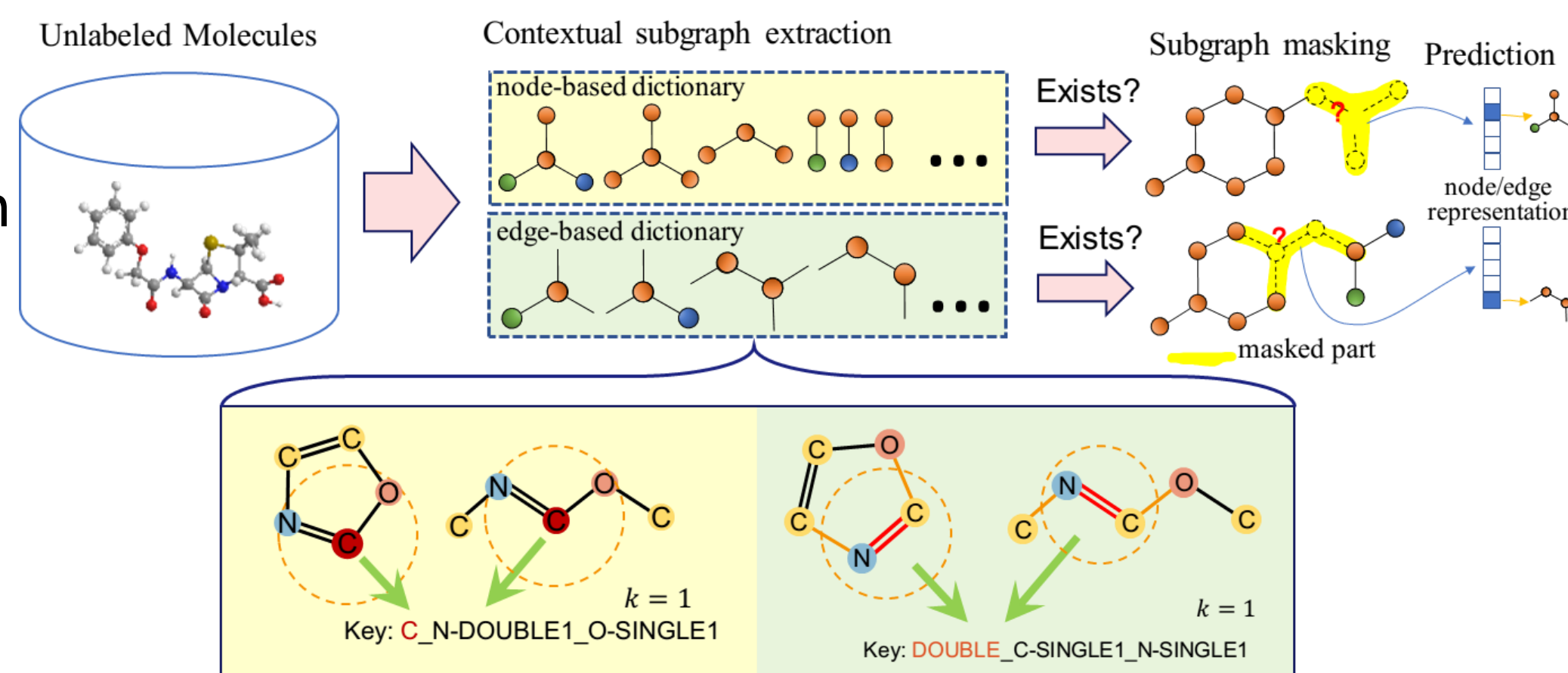


## Proposed: Self-supervised Tasks

Good self-supervision task shall have *reliable* and *cheap* prediction target

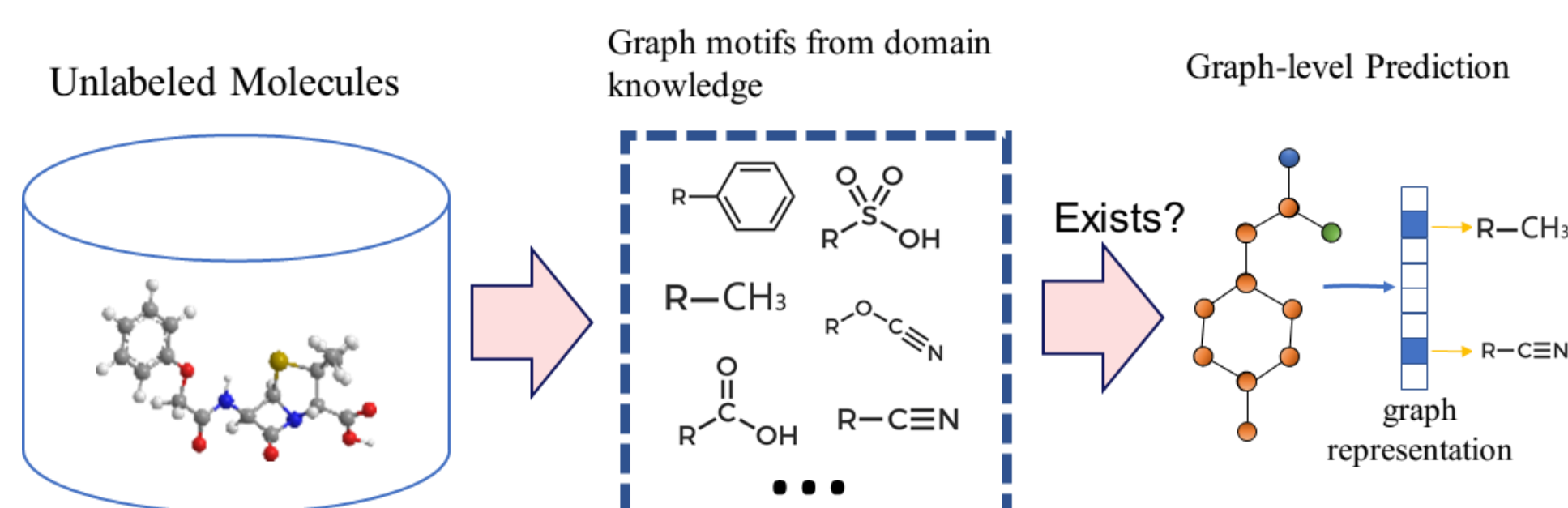
### Node/edge level task: contextual property prediction

Target reflects contextual property:  
recurrent statistical properties of  
local subgraph → A multi-class  
classification problem



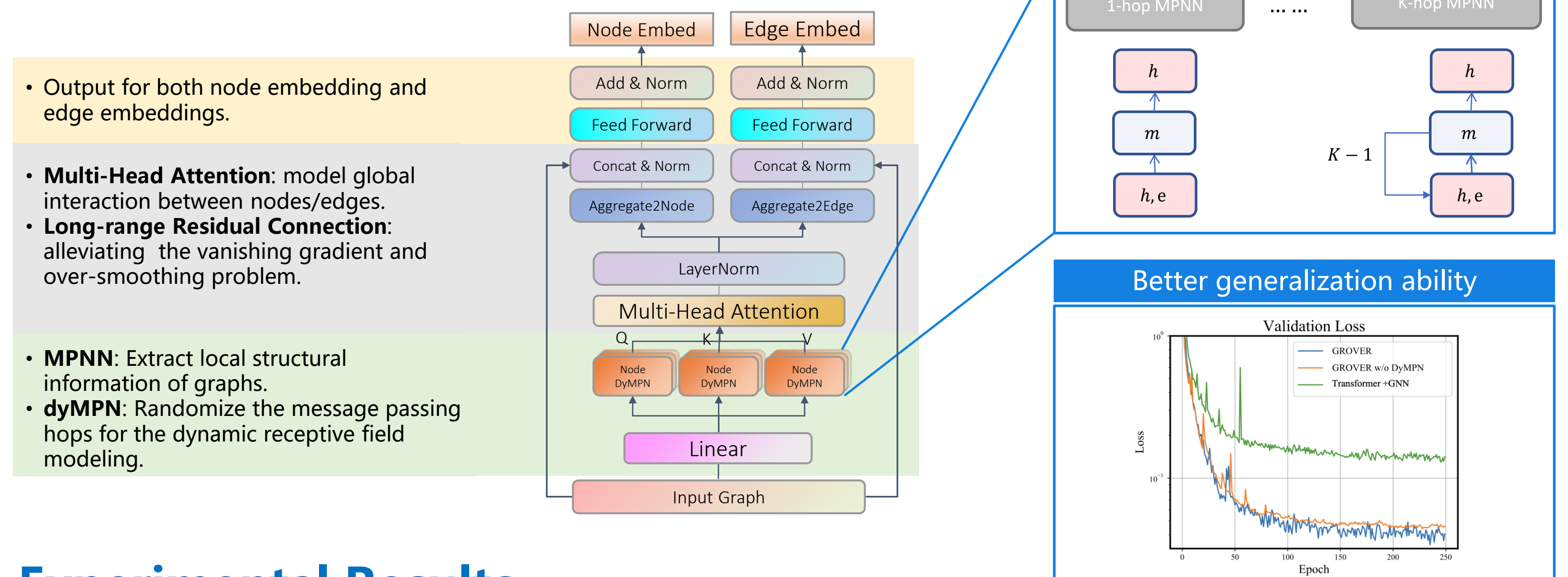
### Graph level task: motif prediction

Motifs: recurrent sub-graphs,  
such as functional groups →  
A multi-label classification  
problem



## Proposed: GTransformer Architecture

A transformer based neural network with tailored GNNs  
as the self-attention building blocks.



## Experimental Results

GROVER<sub>base</sub>: 48M parameters

GROVER<sub>large</sub>: 100M parameters

Pre-training GROVER  
models on 10M  
unlabelled molecules

Verifying on down-  
stream tasks with  
fine-tuning

Significant improvement over  
SOTA models on 11 challenging  
benchmarks

Performance comparison for classification problems (regression results see appendix)

Dataset # Molecules	Classification (Higher is better)					
	BBBP 2039	SIDER 1427	ClinTox 1478	BACE 1513	Tox21 7831	ToxCast 8575
TF_Robust [40]	0.860 <sub>(0.087)</sub>	0.607 <sub>(0.033)</sub>	0.765 <sub>(0.085)</sub>	0.824 <sub>(0.022)</sub>	0.698 <sub>(0.012)</sub>	0.585 <sub>(0.031)</sub>
GraphConv [24]	0.877 <sub>(0.036)</sub>	0.593 <sub>(0.035)</sub>	0.845 <sub>(0.051)</sub>	0.854 <sub>(0.011)</sub>	0.772 <sub>(0.041)</sub>	0.650 <sub>(0.025)</sub>
Weave [23]	0.837 <sub>(0.065)</sub>	0.543 <sub>(0.034)</sub>	0.823 <sub>(0.023)</sub>	0.791 <sub>(0.008)</sub>	0.741 <sub>(0.044)</sub>	0.678 <sub>(0.024)</sub>
SchNet [45]	0.847 <sub>(0.024)</sub>	0.545 <sub>(0.038)</sub>	0.717 <sub>(0.042)</sub>	0.750 <sub>(0.033)</sub>	0.767 <sub>(0.025)</sub>	0.679 <sub>(0.021)</sub>
MPNN [13]	0.913 <sub>(0.041)</sub>	0.595 <sub>(0.030)</sub>	0.879 <sub>(0.054)</sub>	0.815 <sub>(0.044)</sub>	0.808 <sub>(0.024)</sub>	0.691 <sub>(0.013)</sub>
DMPNN [63]	0.919 <sub>(0.030)</sub>	0.632 <sub>(0.023)</sub>	0.897 <sub>(0.040)</sub>	0.852 <sub>(0.053)</sub>	0.826 <sub>(0.023)</sub>	0.718 <sub>(0.011)</sub>
MGCN [30]	0.850 <sub>(0.064)</sub>	0.552 <sub>(0.018)</sub>	0.634 <sub>(0.042)</sub>	0.734 <sub>(0.030)</sub>	0.707 <sub>(0.016)</sub>	0.663 <sub>(0.009)</sub>
AttentiveFP [61]	0.908 <sub>(0.050)</sub>	0.605 <sub>(0.060)</sub>	0.933 <sub>(0.020)</sub>	0.863 <sub>(0.015)</sub>	0.807 <sub>(0.020)</sub>	0.579 <sub>(0.001)</sub>
N-GRAM [29]	0.912 <sub>(0.013)</sub>	0.632 <sub>(0.005)</sub>	0.855 <sub>(0.037)</sub>	0.876 <sub>(0.035)</sub>	0.769 <sub>(0.027)</sub>	0.714 <sub>(0.019)</sub>
HU. et.al[18]	0.915 <sub>(0.040)</sub>	0.614 <sub>(0.006)</sub>	0.762 <sub>(0.058)</sub>	0.851 <sub>(0.027)</sub>	0.811 <sub>(0.015)</sub>	0.723 <sub>(0.010)</sub>
GROVER <sub>base</sub>	0.936 <sub>(0.008)</sub>	0.656 <sub>(0.006)</sub>	0.925 <sub>(0.013)</sub>	0.878 <sub>(0.016)</sub>	0.819 <sub>(0.020)</sub>	0.723 <sub>(0.010)</sub>
GROVER <sub>large</sub>	<b>0.940</b> <sub>(0.019)</sub>	<b>0.658</b> <sub>(0.023)</sub>	<b>0.944</b> <sub>(0.021)</sub>	<b>0.894</b> <sub>(0.028)</sub>	<b>0.831</b> <sub>(0.025)</sub>	<b>0.737</b> <sub>(0.010)</sub>

250 Nvidia V100  
GPUs used

Pre-trained  
methods

❑ GROVER models achieve the  
best performance

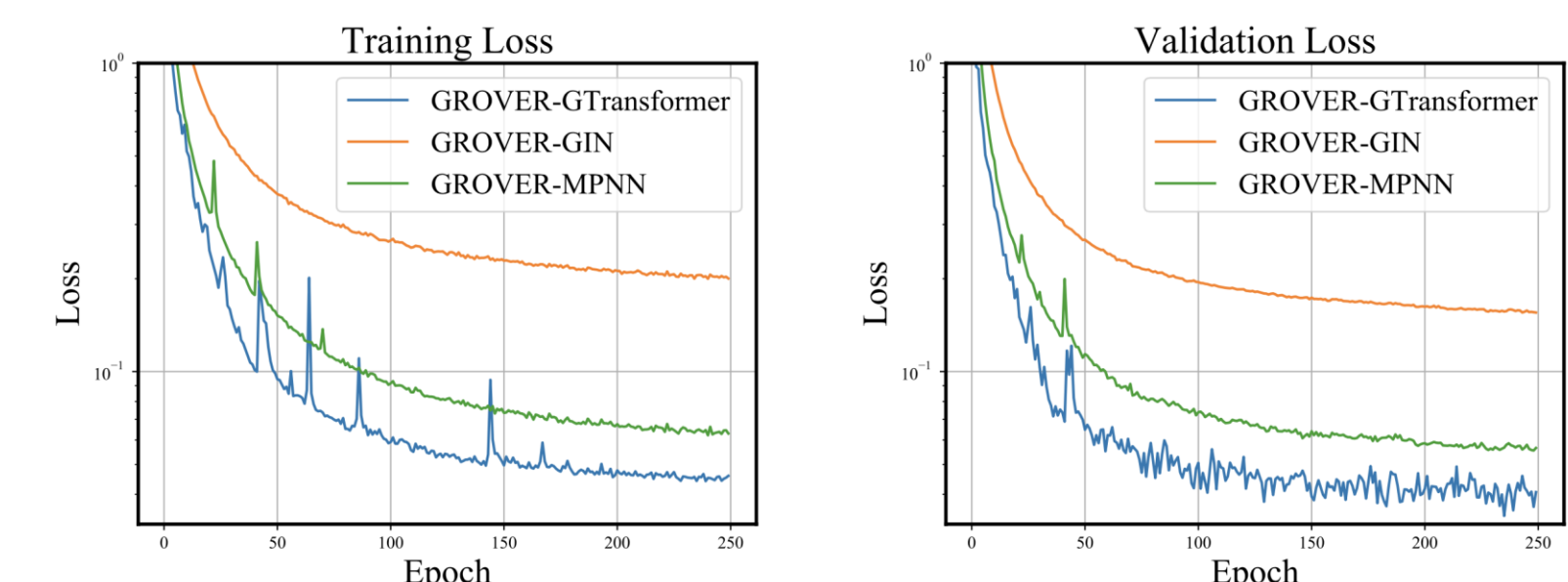
❑ Large model (GROVER<sub>large</sub>)  
enjoys high expressive  
power

❑ Larger improvement  
achieved for dataset with  
less label info.

### Self-supervised pre-training is useful

	GROVER	No Pretrain	Abs. Imp.
BBBP (2039)	<b>0.940</b>	0.911	+0.029
SIDER (1427)	<b>0.658</b>	0.624	+0.034
ClinTox (1478)	<b>0.944</b>	0.884	+0.060
BACE (1513)	<b>0.894</b>	0.858	+0.036
Tox21 (7831)	<b>0.831</b>	0.803	+0.028
ToxCast (8575)	<b>0.737</b>	0.721	+0.016
Average	<b>0.834</b>	0.803	+0.038

### GTransformer backbone matters



### References

Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*, 2019