

# Group Meeting - November 13, 2020

Paper review & Research progress

Truong Son Hy \*

\*Department of Computer Science  
The University of Chicago

Ryerson Physical Lab



- 1 If you want to be successful, you must respect one rule: **Never lie to yourself.**
- 2 All human **unhappiness** comes from **not facing reality** squarely, exactly as it is.
- 3 Ego never accepts the truth.

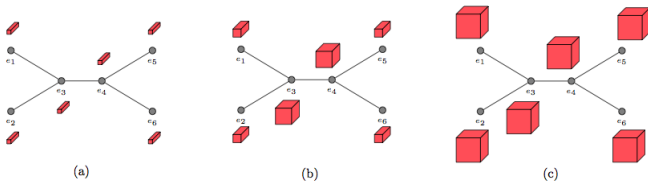


- ① **Self-Supervised Graph Transformer on Large-Scale Molecular Data** (NeurIPS 2020)  
**GROVER: Self-supervised Message Passing Transformer on Large-scale Molecular Data** (preprint)  
<https://arxiv.org/abs/2007.02835>
- ② **A Flexible Generative Framework for Graph-based Semi-supervised Learning** (NeurIPS 2019),  
<https://arxiv.org/abs/1905.10769>



# Brainstorm (1)

**Problem:** How input vertex/atom featurization affects the adjacency reconstruction. Let consider the Ethylene ( $C_2H_4$ ) molecular graph for example example.



- The input atom features of 4 hydrogen atoms are **identical** (the same). The input atom features of 2 carbon atoms are also the same.
- Because of the symmetry, all the tensors (let's consider CCNs for example) associating with hydrogen/carbon atoms are the same.



# Brainstorm (2)

## Problem

There is an edge between e1 (hydrogen) and e3 (carbon). But there is no edge between e5 (hydrogen) and e3 (carbon). Therefore, there is **NO** function can classify edge/non-edge for pairs (e1, e3) and (e5, e3).

## Experiment

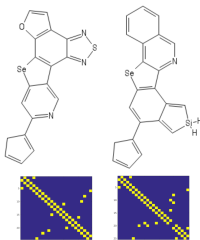
- Use Support Vector Machine (SVM) or simple logistic regression
- To classify edge/non-edge.
- The data are pairs of input atomic features (2 atoms).
- We want to know what is the minimum amount of atomic features needed for each atom to accurately predict edge/non-edge.
- SVM would give us a **reasonable baseline** for the auto-encoding task.
- **Dictionary Weisfeiler-Lehman** to enrich the input features.

# Brainstorm (3)

Markov Random Field (MRF):



- 3/5 have aromatic rings. 2/5 have Benzene. I am thinking of **why?**
- Moving to **Harvard Clean Energy Project (HCEP)** dataset that has a lot of rings?



## **Self-Supervised Graph Transformer on Large-Scale Molecular Data** (NeurIPS 2020)

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

<https://arxiv.org/abs/2007.02835>



# Proposals

## Problem of GNNs

- 1 Insufficient labeled molecules for supervised training.
- 2 Poor generalization capabilities to new-synthesized molecules.

## Proposals

GROVER = **G**raph representation from self-supervised message passing transformer.

- Self-supervised tasks in node, edge and graph-level.
- Transformer-style architecture.
- Pre-train GROVER with 100M parameters on 10M unlabelled molecules → the biggest GNN and the largest training dataset.





# Multi-head Attention

The multi-head attention mechanism is the main building block of various Transformer - style models:

- Parallel running.
- Stacks several scaled dot-product attention layers together. One scaled dot-product attention layer takes a set of queries, keys, values  $(\mathbf{q}, \mathbf{k}, \mathbf{v})$  as inputs:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d}) \mathbf{V}$$

- Suppose we arrange  $k$  attention layers into the multi-head attention:

$$\text{Multi-Head}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Concat}(\text{head}_1, \dots, \text{head}_k) \mathbf{W}^O$$

$$\text{head}_i = \text{Attention}(\mathbf{Q}\mathbf{W}_i^Q, \mathbf{K}\mathbf{W}_i^K, \mathbf{V}\mathbf{W}_i^V)$$

where  $\mathbf{W}_i^Q$ ,  $\mathbf{W}_i^K$  and  $\mathbf{W}_i^V$  are the projection matrices of head  $i$ .



# Graph Neural Networks (GNNs)

Suppose there are  $L$  **iterations**, and iteration  $\ell$  contains  $K_\ell$  **hops**:

$$\mathbf{m}_v^{(\ell,k)} = \text{AGGREGATE}(\{(\mathbf{h}_v^{(\ell,k-1)}, \mathbf{h}_u^{(\ell,k-1)}, \mathbf{e}_{uv}) \mid u \in \mathcal{N}_v\})$$

$$\mathbf{h}_v^{(\ell,k)} = \sigma(\mathbf{W}^{(\ell)} \mathbf{m}_v^{(\ell,k)} + \mathbf{b}^{(\ell)})$$

READOUT operation is applied to get the graph - level representation:

$$\mathbf{h}_G = \text{READOUT}(\{\mathbf{h}_v^{(0,K_0)}, \dots, \mathbf{h}_v^{(L,K_L)} \mid v \in \mathcal{V}\})$$

## Dynamic Message Passing Network (dyMPN):

- The number of hops is closely related to the size of the receptor field.
- Instead of pre-specifying the number of hops, we develop a **randomized strategy** for choosing the number of hops.



# GNN Transformer (GTransformer)

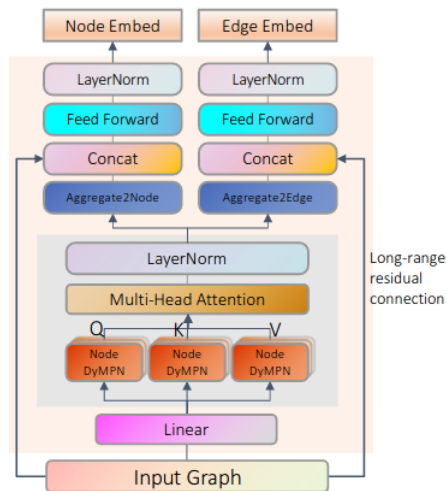


Figure 1: Overview of GTransformer.



# Self-supervised Tasks

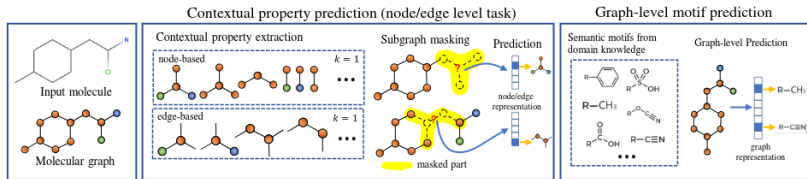


Figure 2: Overview of the designed self-supervised tasks of GROVER.

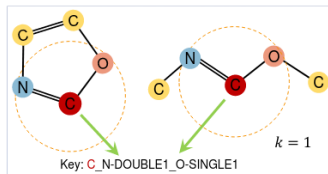


Figure 3: Illustration of contextual properties.



# Benchmark datasets

<http://moleculenet.ai/datasets-1>

## Dataset Details

Category	Dataset	Data Type	Task Type	# Tasks	# Compounds	Rec - Split <sup>a</sup>	Rec - Metric <sup>b</sup>
Quantum Mechanics	QM7	SMILES, 3D coordinates	Regression	1	7160	Stratified	MAE
	QM7b	3D coordinates	Regression	14	7210	Random	MAE
	QM8	SMILES, 3D coordinates	Regression	12	21786	Random	MAE
	QM9	SMILES, 3D coordinates	Regression	12	133885	Random	MAE
Physical Chemistry	ESOL	SMILES	Regression	1	1128	Random	RMSE
	FreeSolv	SMILES	Regression	1	642	Random	RMSE
	Lipophilicity	SMILES	Regression	1	4200	Random	RMSE
Biophysics	PCBA	SMILES	Classification	128	437929	Random	PRC-AUC
	MUV	SMILES	Classification	17	93087	Random	PRC-AUC
	HIV	SMILES	Classification	1	41127	Scaffold	ROC-AUC
	PDBbind	SMILES, 3D coordinates	Regression	1	11908	Time	RMSE
	BACE	SMILES	Classification	1	1513	Scaffold	ROC-AUC
Physiology	BBBP	SMILES	Classification	1	2039	Scaffold	ROC-AUC
	Tox21	SMILES	Classification	12	7831	Random	ROC-AUC
	ToxCast	SMILES	Classification	617	8575	Random	ROC-AUC
	SIDER	SMILES	Classification	27	1427	Random	ROC-AUC
	ClinTox	SMILES	Classification	2	1478	Random	ROC-AUC



# Pre-training performance

Table 1: The performance comparison. The numbers in brackets are the standard deviation. The methods in green are pre-trained methods.

Classification (Higher is better)						
Dataset # Molecules	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>	-
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>	0.714 <sub>(0.019)</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>	0.940 <sub>(0.019)</sub>	0.658 <sub>(0.023)</sub>	0.944 <sub>(0.021)</sub>	0.894 <sub>(0.028)</sub>	0.831 <sub>(0.025)</sub>	0.737 <sub>(0.010)</sub>
Regression (Lower is better)						
Dataset # Molecules	FreeSolv 642	ESOL 1128	Lipo 4200	QM7 6830	QM8 21786	
TF_Robust [40]	4.122 <sub>(0.085)</sub>	1.722 <sub>(0.038)</sub>	0.909 <sub>(0.060)</sub>	120.6 <sub>(9.6)</sub>	0.024 <sub>(0.001)</sub>	
GraphConv [24]	2.900 <sub>(0.135)</sub>	1.068 <sub>(0.050)</sub>	0.712 <sub>(0.049)</sub>	118.9 <sub>(20.2)</sub>	0.021 <sub>(0.001)</sub>	
Weave [23]	2.398 <sub>(0.250)</sub>	1.158 <sub>(0.055)</sub>	0.813 <sub>(0.042)</sub>	94.7 <sub>(2.7)</sub>	0.022 <sub>(0.001)</sub>	
SchNet [45]	3.215 <sub>(0.755)</sub>	1.045 <sub>(0.064)</sub>	0.909 <sub>(0.098)</sub>	74.2 <sub>(6.0)</sub>	0.020 <sub>(0.002)</sub>	
MPNN [13]	2.185 <sub>(0.952)</sub>	1.167 <sub>(0.430)</sub>	0.672 <sub>(0.051)</sub>	113.0 <sub>(17.2)</sub>	0.015 <sub>(0.002)</sub>	
DMPNN [63]	2.177 <sub>(0.914)</sub>	0.980 <sub>(0.258)</sub>	0.653 <sub>(0.046)</sub>	105.8 <sub>(13.2)</sub>	0.0143 <sub>(0.002)</sub>	
MGCN [30]	3.349 <sub>(0.097)</sub>	1.266 <sub>(0.147)</sub>	1.113 <sub>(0.041)</sub>	77.6 <sub>(4.7)</sub>	0.022 <sub>(0.002)</sub>	
AttentiveFP [61]	2.030 <sub>(0.420)</sub>	0.853 <sub>(0.060)</sub>	0.650 <sub>(0.030)</sub>	126.7 <sub>(4.0)</sub>	0.0282 <sub>(0.001)</sub>	
N-GRAM [29]	2.512 <sub>(0.190)</sub>	1.100 <sub>(0.160)</sub>	0.876 <sub>(0.033)</sub>	125.6 <sub>(1.5)</sub>	0.0320 <sub>(0.003)</sub>	
GROVER <sub>base</sub>	1.592 <sub>(0.072)</sub>	0.888 <sub>(0.116)</sub>	0.563 <sub>(0.030)</sub>	72.5 <sub>(5.9)</sub>	0.0172 <sub>(0.002)</sub>	
GROVER <sub>large</sub>	1.544 <sub>(0.397)</sub>	0.831 <sub>(0.120)</sub>	0.560 <sub>(0.035)</sub>	72.6 <sub>(3.8)</sub>	0.0125 <sub>(0.002)</sub>	



# Comparing to its variants

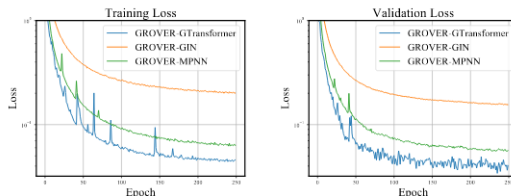


Figure 4: The training and validation losses on different backbones.

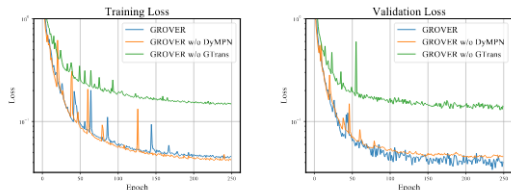


Figure 5: The training and validation loss of GROVER and its variants.



# Comparing to without pre-training

	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

Table 2: Comparison between GROVER with and without pre-training.





## **A Flexible Generative Framework for Graph-based Semi-supervised Learning (NeurIPS 2019)**

Jiaqi Ma, Weijing Tang, Ji Zhu, Qiaozhu Mei

<https://arxiv.org/abs/1905.10769>



# Proposals

## Proposals

- 1 Generative framework for graph-based SSL.
- 2 Node features, labels, and graph structure are modeled in a joint distribution.
- 3 Employ the variational inference techniques to approximate the Bayesian posterior.

## Note

The datasets in the experiment are too small. The work should be scalable for a much larger network like IMAGENET.



# Graph-based regularization for semi-supervised learning

Suppose there are  $n$  data samples in total and  $m$  of them are labeled, the graph-based regularization methods generally conduct semi-supervised learning by optimizing the objective:

$$\sum_{i=1}^m \mathcal{L}_i + \eta \sum_{i,j=1}^n w_{i,j} \mathcal{R}(\mathbf{f}_i, \mathbf{f}_j)$$

where:

- $\mathcal{L}_i$  is the supervised loss function of sample  $i$ .
- $\mathcal{R}(\cdot, \cdot)$  is a regularization function.
- $w_{ij}$  is a graph-based coefficient.
- $\mathbf{f}_i, \mathbf{f}_j$  are the outcome predictions or feature representations.
- $\eta$  is a hyper-parameter trade-off the supervised loss and the graph-based regularization.



# A flexible generative framework for graph-based SSL (1)

Denote  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\mathbf{Y} \in \mathbb{R}^{n \times \ell}$  as feature and outcome vectors. Partition the outcome matrix as:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_{obs} \\ \mathbf{Y}_{miss} \end{bmatrix}$$

The goal of graph-based SSL is to infer  $\mathbf{Y}_{miss}$  based on  $(\mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})$ .

- **Generation process.** The generation process can be illustrated by the following factorization of the joint distribution:

$$p(\mathbf{X}, \mathbf{Y}, \mathbf{G}) = p(\mathbf{G}|\mathbf{X}, \mathbf{Y})p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})$$

where the first and second terms are parameterized by  $\theta$ :  
 $p_{\theta}(\mathbf{G}|\mathbf{X}, \mathbf{Y})$  and  $p_{\theta}(\mathbf{Y}|\mathbf{X})$ .



# A flexible generative framework for graph-based SSL (2)

- **Model inference.** To infer the missing outcomes  $\mathbf{Y}_{miss}$ , we would need the posterior distribution  $p_{\theta}(\mathbf{Y}_{miss}|\mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})$  that is usually intractable.  
→ Introduce the recognition model  $q_{\phi}(\mathbf{Y}_{miss}|\mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})$  (variational principles)
- **Model learning.**

$$\log p(\mathbf{Y}_{obs}, \mathbf{G}|\mathbf{X}) \geq -\mathcal{L}_{ELBO}(\theta, \phi, \mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})$$

where the lower bound is defined as:

$$\mathbb{E}_{q_{\phi}(\mathbf{Y}_{miss}|\mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})} \left[ \log p_{\theta}(\mathbf{Y}_{miss}, \mathbf{Y}_{obs}, \mathbf{G}|\mathbf{X}) - \log q_{\phi}(\mathbf{Y}_{miss}|\mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G}) \right]$$

Optimization:

$$\hat{\theta}, \hat{\phi} = \arg \min_{\theta, \phi} \mathcal{L}_{ELBO}(\theta, \phi; \mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G})$$



# A flexible generative framework for graph-based SSL (3)

- **Latent space models (LSM).** Based on the conditional independence assumption of edges:

$$p_{\theta}(\mathbf{G}|\mathbf{X}, \mathbf{Y}) = \prod_{i,j} p_{\theta}(e_{ij}|\mathbf{X}, \mathbf{Y})$$

where:

$$p_{\theta}(e_{ij}|\mathbf{X}, \mathbf{Y}) = p_{\theta}(e_{ij}|\mathbf{x}_i, \mathbf{y}_i, \mathbf{x}_j, \mathbf{y}_j)$$

is modeled by a logistic regression.

- **Supervised loss.** Additional term with weight hyper-parameter  $\eta$  for the approximate posterior model:

$$\mathcal{L}(\theta, \phi) = \mathcal{L}_{ELBO}(\theta, \phi; \mathbf{X}, \mathbf{Y}_{obs}, \mathbf{G}) - \eta \cdot \log q_{\phi}(\mathbf{Y}_{obs}|\mathbf{X}, \mathbf{G})$$



# Standard benchmark

Table 1: Summary of benchmark datasets.

Dataset	# Classes	# Nodes	# Edges	Avg. 2-Neighborhood Size
Cora	7	2,708	5,278	35.8
Pubmed	3	19,717	44,324	59.1
Citeseer	6	3,327	4,552	14.1

Table 2: Classification accuracy under the standard benchmark setting. The upper block lists the discriminative baselines. The lower block lists the proposed variants of G<sup>3</sup>NN. The **bold** marker denotes the best performance on each dataset. The underline marker denotes that the generative model outperforms its discriminative counterpart, e.g., LSM-GCN outperforms GCN; and the asterisk (\*) marker denotes the difference is statistically significant by a t-test at significance level 0.05. The ( $\pm$ ) error bar denotes the standard deviation of the test performance of 10 independent trials.

	Cora	Pubmed	Citeseer
MLP	0.583 $\pm$ 0.009	0.734 $\pm$ 0.002	0.569 $\pm$ 0.008
GCN	0.815 $\pm$ 0.002	<b>0.794</b> $\pm$ 0.004	0.718 $\pm$ 0.003
GAT	0.825 $\pm$ 0.005	0.785 $\pm$ 0.004	0.715 $\pm$ 0.007
LSM_GCN	<u>0.825</u> $\pm$ 0.002*	0.779 $\pm$ 0.004	<u>0.744</u> $\pm$ 0.003*
LSM_GAT	<b><u>0.829</u></b> $\pm$ 0.003	0.776 $\pm$ 0.007	<u>0.731</u> $\pm$ 0.005*
SBM_GCN	<u>0.822</u> $\pm$ 0.002*	0.784 $\pm$ 0.006	<b><u>0.745</u></b> $\pm$ 0.004*
SBM_GAT	<u>0.829</u> $\pm$ 0.003	0.774 $\pm$ 0.004	<u>0.740</u> $\pm$ 0.003*



# Missing-edges setting

Setting: Remove all the edges of the test nodes from the graph.

Table 3: Classification accuracy under the missing-edge setting. The **bold** marker, the underline marker, the asterisk (\*) marker, and the ( $\pm$ ) error bar share the same definitions in Table 2

	Cora	Pubmed	Citeseer
MLP	$0.583 \pm 0.009$	$0.734 \pm 0.002$	$0.569 \pm 0.008$
GCN	$0.665 \pm 0.007$	$0.746 \pm 0.004$	$0.652 \pm 0.005$
GAT	$0.682 \pm 0.004$	$0.744 \pm 0.006$	$0.642 \pm 0.004$
LSM_GCN	<u><math>0.711 \pm 0.005^*</math></u>	<u><b><math>0.766 \pm 0.006^*</math></b></u>	<u><math>0.704 \pm 0.002^*</math></u>
LSM_GAT	<u><math>0.710 \pm 0.007^*</math></u>	<u><math>0.766 \pm 0.004^*</math></u>	<u><math>0.691 \pm 0.005^*</math></u>
SBM_GCN	<u><b><math>0.718 \pm 0.004^*</math></b></u>	<u><math>0.762 \pm 0.005^*</math></u>	<u><b><math>0.716 \pm 0.004^*</math></b></u>
SBM_GAT	<u><math>0.716 \pm 0.007^*</math></u>	<u><math>0.761 \pm 0.005^*</math></u>	<u><math>0.709 \pm 0.008^*</math></u>





# Reduced-labels setting

Setting: Drop half of the training labels for each class compared to the standard benchmark setting.

Table 4: Classification accuracy under the reduced-label setting. The **bold** marker, the underline marker, the asterisk (\*) marker, and the ( $\pm$ ) error bar share the same definitions in Table 2

	Cora	Pubmed	Citeseer
MLP	$0.498 \pm 0.004$	$0.674 \pm 0.005$	$0.493 \pm 0.010$
GCN	$0.750 \pm 0.003$	<b><math>0.724 \pm 0.005</math></b>	$0.666 \pm 0.003$
GAT	$0.771 \pm 0.004$	$0.711 \pm 0.006$	$0.675 \pm 0.005$
LSM_GCN	<u><math>0.777 \pm 0.002^*</math></u>	$0.709 \pm 0.003$	<u><math>0.691 \pm 0.005^*</math></u>
LSM_GAT	<u><math>0.792 \pm 0.004^*</math></u>	$0.699 \pm 0.003$	<u><math>0.691 \pm 0.004^*</math></u>
SBM_GCN	<u><math>0.780 \pm 0.002^*</math></u>	$0.710 \pm 0.004$	<u><b><math>0.703 \pm 0.006^*</math></b></u>
SBM_GAT	<u><b><math>0.796 \pm 0.008^*</math></b></u>	$0.699 \pm 0.003$	<u><math>0.698 \pm 0.003^*</math></u>

