

# Graph Self-supervised Learning and Pre-training

# Yukuo Cen

GNN Center, Zhipu Al KEG, Tsinghua University Advisors: Yuxiao Dong, Jie Tang

Course Link: https://cogdl.ai/gnn2022/

CogDL is publicly available at <a href="https://github.com/THUDM/cogdl">https://github.com/THUDM/cogdl</a>



### **Real-world Graphs**



# **Pre-training and Fine-tuning**



Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. In KDD'20.

GNN pre-training problem.

# PROBLEM

# **The GNN Pre-Training Problem**

- Problem:
  - Learn a function *f* that maps a vertex to a low-dimensional vector
  - Structural similarity: map vertices with similar local network topologies close in the vector space
  - Transferability: compatible with vertices and graphs from various sources, even unseen during training time.



Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. In KDD'20.

Graph Contrastive Coding

# **Graph Contrastive Coding (GCC)**



Hypothesis: Graph structural patterns are universal and transferable across networks.

Jiezhong Qiu, Qibin Chen, Yuxiao Dong, Jing Zhang, Hongxia Yang, Ming Ding, Kuansan Wang, and Jie Tang. GCC: Graph Contrastive Coding for Graph Neural Network Pre-Training. In KDD'20.

# **GCC Pre-training**

- **Pre-training Task: Instance** Discrimination
- InfoNCE Loss: output instance representations that are capable of capturing the similarities between instances

$$\mathcal{L} = -\log \frac{\exp\left(\boldsymbol{q}^{\top}\boldsymbol{k}_{+}/\tau\right)}{\sum_{i=0}^{K}\exp\left(\boldsymbol{q}^{\top}\boldsymbol{k}_{i}/\tau\right)}$$

- query instance  $x^q$
- query q (embedding of  $x^q$ ), i.e.,  $q = f(x^q)$
- dictionary of keys  $\{k_0, k_1, \cdots, k_K\}$
- key  $\boldsymbol{k} = f(x^k)$
- Contrastive learning for graphs?
  - Q1: How to define instances in graphs?
  - Q2: How to define (dis) similar instance pairs?
  - Q3: What are the proper encoders?

# **GCC Pre-training**

- **Q1:** How to define **instances** in graphs?
- Q2: How to define (dis) similar instance?
- Q3: What are the proper encoders?



# GCC Pre-training: Learning Algorithms

- Optimizing Contrastive Loss
  - Encoded query q
  - K + 1 encoded keys { $k_0, \cdots, k_K$ }

$$\mathcal{L} = -\log \frac{\exp\left(\boldsymbol{q}^{\top}\boldsymbol{k}_{+}/\tau\right)}{\sum_{i=0}^{K}\exp\left(\boldsymbol{q}^{\top}\boldsymbol{k}_{i}/\tau\right)}$$



#### figure credit:

Momentum Contrast for Unsupervised Visual Representation Learning arxiv.org/abs/1911.05722

## **GCC** Fine-tuning



**Fine-Tuning** 

Fine-tuning



# GCC Fine-tuning: Full v.s. Freezing

**Full fine-tuning** 

**Freezing fine-tuning** 



# **EXPERIMENTS**

# GCC Pre-Training / Fine-tuning

• Six real-world information networks for pre-training.

Table 1: Datasets for pr	e-training, sorted b	by number of vertices.
--------------------------	----------------------	------------------------

Dataset	Academia	DBLP (SNAP)	DBLP (NetRep)	IMDB	Facebook	LiveJournal
V	137,969	317,080	540,486	896,305	3,097,165	4,843,953
E	739,384	2,099,732	30,491,458	7,564,894	47,334,788	85,691,368

- Fine-tuning Tasks:
  - Node classification
  - Graph classification
  - Top-k Similarity search



Subgraph Instance Discrimination

GCC: Graph Contrastive Coding

**Fine-Tuning** 

# **Task 1: Node Classification**

- Setup
  - US-Airport
  - AMiner academic graph



Datasets	US-Airport	H-index
V	1,190 13 599	5,000 44 020
ProNE	62.3	69.1
GraphWave	60.2	70.3
Struc2vec	66.2	> 1 Day
GCC (E2E, Heeze) GCC (MoCo, freeze)	65.6	78.3
GCC (rand, full)	64.2	76.9
GCC (E2E, full)	68.3	80.5
GCC (MoCo, full)	67.2	80.6

## **Task 2: Graph Classification**

• Setup

#### – COLLAB, RDT-B, RDT-M, & IMDB-B, IMDB-M



Datasets	IMDB-B	IMDB-M	COLLAB	RDT-B	RDT-M
# graphs	1,000	1,500	5,000	2,000	5,000
# classes	2	3	3	2	5
Avg. # nodes	19.8	13.0	74.5	429.6	508.5
DGK	67.0	44.6	73.1	78.0	41.3
graph2vec	71.1	50.4	-	75.8	47.9
InfoGraph	73.0	49.7	-	82.5	53.5
GCC (E2E, freeze)	71.7	49.3	74.7	87.5	52.6
GCC (MoCo, freeze)	72.0	49.4	78.9	89.8	53.7
DGCNN	70.0	47.8	73.7	_	_
GIN	75.6	51.5	80.2	89.4	54.5
GCC (rand, full)	75.6	50.9	79.4	87.8	52.1
GCC (E2E, full)	70.8	48.5	79.0	86.4	47.4
GCC (MoCo, full)	73.8	50.3	81.1	87.6	53.0

# **Task 3: Top-k Similarity Search**

• Setup

- AMiner academic graph



	KDD-	ICDM	SIGIR	-CIKM	SIGMOD-ICDE		
V	2,867	2,607	2,851	3,548	2,616	2,559	
E	7,637	4,774	6,354	7,076	8,304	6,668	
# groud truth		697		874		898	
k	20	40	20	40	20	40	
Random	0.0198	0.0566	0.0223	0.0447	0.0221	0.0521	
RolX	0.0779	0.1288	0.0548	0.0984	0.0776	0.1309	
Panther++	0.0892	0.1558	0.0782	0.1185	0.0921	0.1320	
GraphWave	0.0846	0.1693	0.0549	0.0995	0.0947	0.1470	
GCC (E2E)	0.1047	0.1564	0.0549	0.1247	0.0835	0.1336	
GCC (MoCo)	0.0904	0.1521	0.0652	0.1178	0.0846	0.1425	

# Conclusion



- Study the pre-training of GNN with the goal of characterizing and transferring structural representations in social and information networks.
- Present Graph Contrastive Coding, which is a graph-based contrastive learning framework to pre-train GNN.
- The pre-trained GNN achieves competitive performance to its supervised trained-from-scratch counterparts in 3 graph learning tasks on 10 graph datasets.



# GraphMAE: Self-Supervised Masked Graph Autoencoders

### Self-supervised Learning

- Self-supervised learning enables the model to learn informative representations from **unlabeled data**
- Various SSL methods on graph have been developed



### Self-supervised learning on graph

- Contrastive SSL has been the dominant approach in recent years.
  - Especially in classification tasks.
  - Generative methods fail to achieve comparable results

	Dataset	Cora	CiteSeer	PubMed
Come constant d	GCN	81.5	70.3	79.0
Supervised	GAT	83.0±0.7	$72.5 \pm 0.7$	$79.0 {\pm} 0.3$
	GAE	71.5±0.4	65.8±0.4	$72.1 {\pm} 0.5$
metho	GPT-GNN	80.1±1.0	68.4±1.6	$76.3 {\pm} 0.8$
Contrastive	GATE	83.2±0.6	$71.8 \pm 0.8$	$80.9 \pm 0.3$
Com	DGI	82.3±0.6	$71.8 \pm 0.7$	$76.8 {\pm} 0.6$
	MVGRL	83.5±0.4	73.3±0.5	$80.1 {\pm} 0.7$
Self-supervised	GRACE <sup>1</sup>	81.9±0.4	$71.2 \pm 0.5$	$80.6 {\pm} 0.4$
	BGRL <sup>1</sup>	82.7±0.6	$71.1 \pm 0.8$	79.6±0.5
	InfoGCL	83.5±0.3	73.5±0.4	79.1±0.2
	CCA-SSG <sup>1</sup>	$84.0 \pm 0.4$	73.1±0.3	$81.0 \pm 0.4$

### Self-supervised learning on graph

- Contrastive learning relies on complicated and elaborate designs,
- Contrastive SSL could fail if lacking any one component.
  - Negative sampling design
    - In-batch negatives (GRACE, GCA, GraphCL)
    - Dynamic queues as negatives (GCC,)
    - Shuffle node features as negatives (DGI, MVGRL)
  - Architecture design
    - Asymmetric encoder, Projection head (BGRL, SimGRACE)
    - Feature de-correlation (CCA-SSG,)
  - Data augmentation design
    - Node dropping, Edge perturbation, Subgraph Sampling (GraphCL, CCA-SSG, BGRL)
    - Graph Diffusion (MVGRL, ), Random-walk (GCC, ), Infomax Augmentation (Info-GCL)
    - ...

Generative SSL can naturally avoid these issues

#### Generative SSL on graphs

- Generative SSL can naturally avoid the issue of relying on complicated strategies.
- But previous generative SSL often fail to catch up with the performance of contrastive methods.

What makes graph generative learning lag behind contrastive learning ?

Generative SSL has been gaining increasing significance

- MAE[He, 2021] leads the revolution and opens the new era of generative SSL
  - Get rid of all constraints of contrastive SSL



How do we unleash the potential of generative SSL on graphs?

# Problem

Graph Autoencoders

#### Graph AutoEncoder

- G = (V, A, X)-  $A \in \{0, 1\}^{N \times N}$ : adjacency matrix,
  - $X \in \mathbb{R}^{N \times d}$ : node features
- Encoding:  $H = f_E(A, X)$ ,
- Decoding:  $G' = f_D(A, H)$
- Reconstruction objectives:
   graph structure (link)
  - Node features



#### Summary of Graph AutoEncoder

4. Error functio					3. Deco	ding stra	tegy
	1.]	Reco	onstru	ction 7	Farget		
				2. R	econstru	ction me	thod
	↓ ·		<b>I</b>	<b>I</b>			
Methods	Feat.	٨F	No	Mask	GNN	Re-mask	Space
Methods	Loss	ΛĽ	Struc.	Feat.	Decoder	Dec.	Space
VGAE [20]	n/a	$\checkmark$	-	-	-	-	$O(N^2)$
ARVGA [26]	n/a	$\checkmark$	-	-	-	-	$O(N^2)$
MGAE [42]	MSE	$\checkmark$	-	$\checkmark$	-	-	O(N)
GALA [27]	MSE	$\checkmark$	$\checkmark$	-	$\checkmark$	-	O(N)
GATE [31]	MSE	$\checkmark$	-	-	$\checkmark$	-	O(N)
AttrMask [16]	CE	$\checkmark$	$\checkmark$	$\checkmark$	-	-	O(N)
GPT-GNN [17]	MSE	-	-	$\checkmark$	-	-	O(N)
AGE [3]	n/a	$\checkmark$	-	-	-	-	$O(N^2)$
NodeProp [18]	MSE	$\checkmark$	$\checkmark$	$\checkmark$	-	-	O(N)
GraphMAE	SCE	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	O(N)

(a) Technical comparison between generative SSL methods.

#### Critical components

- 1. Link reconstruction may be over-emphasized.
- 2. Reconstruction without corruption may not be robust



(b) The effect of GraphMAE designs on the performance on Cora dataset.

#### Critical components

- 3. Linear/MLP is a less expressive decoding strategy
- 4. MSE may not be a good criterion for feature reconstruction in graph



(b) The effect of GraphMAE designs on the performance on Cora dataset.

#### Generative SSL for graph?

- 1. What to reconstruct ?
- 2. How to avoid trivial solutions ?
- 3. How to design the decoding ?
- 4. What error function to use ?



(b) The effect of GraphMAE designs on the performance on Cora dataset.

# GraphMAE Framework

Self-Supervised Masked Graph Autoencoders

#### GraphMAE Method



- Masked feature reconstruction
- GNN as decoder with re-mask decoding
- Scaled cosine error as the Criterion

#### Masked feature reconstruction

- Feature construction as the learning objective
- Masked feature reconstruction
  - 1. Sample a subset of nodes

$$\widetilde{V} \subset V$$

2. Replace node feature with [MASK]



2. How to avoid trivial solutions ?



#### GNN as decoder with re-mask decoding

- Use GNN as the decoder
  - A more expressive decoder helps reconstruct low informative features
- Re-mask node features before decoder
  - Re-mask the "masked" nodes

3. How to design the decoding ?

34



#### Scaled cosine error as the criterion

• MSE fails, especially for continuous features

- Sensitivity & low selectivity

$$L_{MSE} = \frac{1}{|\tilde{V}|} \sum_{v_i \in \tilde{V}} (x_i - z_i)^2$$

- Scaled cosine error as the criterion.
  - Cosine error & Scaled coefficient

$$\mathcal{L}_{\text{SCE}} = \frac{1}{|\widetilde{\mathcal{V}}|} \sum_{v_i \in \widetilde{\mathcal{V}}} (1 - \frac{\boldsymbol{x}_i^T \boldsymbol{z}_i}{\|\boldsymbol{x}_i\| \cdot \|\boldsymbol{z}_i\|})^{\gamma}, \ \gamma \ge 1,$$

4. What error function to use ?



#### Overall Improvement of GraphMAE



(b) The effect of GraphMAE designs on the performance on Cora dataset. Node classification

# Experiments

- Unsupervised representation learning for *Node classification*
- Unsupervised representation learning for *Graph* classification
- *Transfer learning* on molecular property prediction

#### Node classification

**Table 1: Experiment results in unsupervised representation learning for** <u>**node classification</u>**. We report Micro-F1(%) score for PPI and accuracy(%) for the other datasets.</u>

	Dataset	Cora	CiteSeer	PubMed	Ogbn-arxiv	PPI	Reddit
Supervised	GCN	81.5	70.3	79.0	71.74±0.29	75.7±0.1	95.3±0.1
Supervised	GAT	83.0±0.7	$72.5 \pm 0.7$	$79.0 {\pm} 0.3$	$72.10 {\pm} 0.13$	$97.30 {\pm} 0.20$	96.0±0.1
	GAE	71.5±0.4	$65.8 {\pm} 0.4$	$72.1 {\pm} 0.5$	-	-	-
methoe	GPT-GNN	80.1±1.0	$68.4 \pm 1.6$	$76.3 \pm 0.8$	-	-	-
antrastive .	GATE	83.2±0.6	$71.8 \pm 0.8$	$80.9 \pm 0.3$	-	-	-
Com	DGI	82.3±0.6	$71.8 \pm 0.7$	$76.8 \pm 0.6$	$70.34 \pm 0.16$	$63.80 {\pm} 0.20$	$94.0 \pm 0.10$
	MVGRL	83.5±0.4	$73.3 \pm 0.5$	$80.1 {\pm} 0.7$	-	-	-
Self-supervised	GRACE <sup>1</sup>	81.9±0.4	$71.2 \pm 0.5$	$80.6 \pm 0.4$	$71.51 \pm 0.11$	$69.71 \pm 0.17$	$94.72 {\pm} 0.04$
	BGRL <sup>1</sup>	82.7±0.6	$71.1 \pm 0.8$	$79.6 \pm 0.5$	$71.64 \pm 0.12$	$73.63 \pm 0.16$	$94.22 \pm 0.03$
	InfoGCL	83.5±0.3	73.5±0.4	$79.1 \pm 0.2$	-	-	-
	CCA-SSG <sup>1</sup>	$\underline{84.0\pm0.4}$	$73.1 \pm 0.3$	$\underline{81.0\pm0.4}$	$71.24 {\pm} 0.20$	$73.34 {\pm} 0.17$	$95.07 \pm 0.02$
	GraphMAE	84.2±0.4	73.4±0.4	81.1±0.4	71.75±0.17	74.50±0.29	96.01±0.08

#### Graph classification

#### Table 2: Experiment results in unsupervised representation learning for graph classification. We report accuracy(%) for all datasets.

	Dataset	IMDB-B	IMDB-M	PROTEINS	COLLAB	MUTAG	REDDIT-B	NCI1
Supervised	GIN	75.1±5.1	$52.3 \pm 2.8$	$76.2 \pm 2.8$	80.2±1.9	89.4±5.6	92.4±2.5	82.7±1.7
Supervised	DiffPool	72.6±3.9	-	$75.1 \pm 3.5$	$78.9 \pm 2.3$	$85.0 \pm 10.3$	92.1±2.6	-
Crearly Vormala	WL	72.30±3.44	46.95±0.46	$72.92 \pm 0.56$	-	80.72±3.00	$68.82 \pm 0.41$	80.31±0.46
Graph Kernels	DGK	66.96±0.56	$44.55 \pm 0.52$	$73.30 {\pm} 0.82$	-	$87.44 \pm 2.72$	$78.04 \pm 0.39$	$80.31 {\pm} 0.46$
d	graph2vec	71.10±0.54	$50.44 {\pm} 0.87$	$73.30{\pm}2.05$	-	83.15±9.25	75.78±1.03	$73.22 \pm 1.81$
ive method	Infograph	73.03±0.87	$49.69 \pm 0.53$	$74.44 \pm 0.31$	$70.65 \pm 1.13$	$89.01 \pm 1.13$	$82.50 \pm 1.42$	$76.20 \pm 1.06$
Contrastive	GraphCL	71.14±0.44	$48.58 \pm 0.67$	$74.39 {\pm} 0.45$	$71.36 \pm 1.15$	$86.80 \pm 1.34$	<u>89.53±0.84</u>	$77.87 \pm 0.41$
Colf annowigod	JOAO	70.21±3.08	$49.20 \pm 0.77$	$74.55 \pm 0.41$	$69.50 {\pm} 0.36$	$87.35 \pm 1.02$	$85.29 \pm 1.35$	$78.07 \pm 0.47$
Sell-supervised	GCC	72.0	49.4	-	78.9	-	89.8	-
	MVGRL	$74.20 \pm 0.70$	$51.20 \pm 0.50$	-	-	$89.70 \pm 1.10$	$84.50 \pm 0.60$	-
	InfoGCL	$75.10 \pm 0.90$	$\underline{51.40{\pm}0.80}$	-	$80.00 \pm 1.30$	$91.20{\pm}1.30$	-	$\underline{80.20{\pm}0.60}$
	GraphMAE	75.52±0.66	51.63±0.52	75.30±0.39	80.32±0.46	88.19±1.26	88.01±0.19	80.40±0.30

#### Transfer learning

 Table 3: Experiment results in transfer learning on molecular property prediction benchmarks. The model is first pre-trained on ZINC15 and then finetuned on the following datasets. We report ROC-AUC(%) scores.

_	ethos									
Contra	stive me	BBBP	Tox21	ToxCast	SIDER	ClinTox	MUV	HIV	BACE	Avg.
	No-pretrain	65.5±1.8	74.3±0.5	63.3±1.5	57.2±0.7	58.2±2.8	71.7±2.3	75.4±1.5	70.0±2.5	67.0
_	ContextPred	64.3±2.8	<u>75.7±0.7</u>	63.9±0.6	60.9±0.6	65.9±3.8	75.8±1.7	77.3±1.0	79.6±1.2	70.4
	AttrMasking	64.3±2.8	76.7±0.4	64.2±0.5	<u>61.0±0.7</u>	$71.8 \pm 4.1$	74.7±1.4	$77.2 \pm 1.1$	79.3±1.6	71.1
	Infomax	68.8 ±0.8	75.3 ±0.5	$62.7 \pm 0.4$	$58.4 \pm 0.8$	69.9±3.0	75.3 ±2.5	$76.0 \pm 0.7$	75.9 ±1.6	70.3
	GraphCL	69.7±0.7	73.9±0.7	62.4±0.6	60.5±0.9	76.0±2.7	69.8±2.7	78.5±1.2	75.4±1.4	70.8
	JOAO	70.2±1.0	75.0±0.3	62.9±0.5	60.0±0.8	<u>81.3±2.5</u>	71.7±1.4	76.7±1.2	77.3±0.5	71.9
	GraphLoG	72.5±0.8	<u>75.7±0.5</u>	63.5±0.7	61.2±1.1	76.7±3.3	<u>76.0±1.1</u>	77.8±0.8	83.5±1.2	<u>73.4</u>
-	GraphMAE	72.0±0.6	75.5±0.6	<u>64.1±0.3</u>	60.3±1.1	82.3±1.2	76.3±2.4	77.2±1.0	<u>83.1±0.9</u>	73.8

rads

#### Ablation study

• Effect of decoder type, objective function and mask strategy

Table 4: Ablation studies of decoder type, re-mask and reconstruction criterion on node- and graph-level benchmarks.

	Dataset		Node-Leve	el	Graph	-Level
	Dutabet	Cora	PubMed	Arxiv	MUTAG	IMDB-B
	GraphMAE	84.2	81.1	71.75	88.19	75.52
MP.	w/o mask	79.7	77.9	70.97	82.58	74.42
CO	w/o re-mask	82.7	80.0	71.61	86.29	74.42
_	w/ MSE	79.1	73.1	67.44	86.30	74.04
	MLP	82.2	80.4	71.54	87.16	73.94
ode	GCN	81.3	79.1	71.59	87.78	74.54
)ecc	GIN	81.8	80.2	71.41	88.19	75.52
П	GAT	84.2	81.1	71.75	86.27	74.04

## Summary

- Explore generative self-supervised learning in graphs
- Identify the common issues in current graph autoencoders.
- Present a simple and improved masked autoencoder— GraphMAE
- The experimental results show that generative SSL can have great potential



# Summary

- GCC: Graph Contrastive Coding for GNN Pre-Training
  - Motivation: universal structural patterns across networks?
  - Pre-training on the graph structure via contrastive learning
  - Fine-tuning on different downstream graph tasks

- GraphMAE: Self-Supervised Masked Graph Autoencoders
  - Motivation: why graph generative learning lag behind graph contrastive learning?
  - Target generative self-supervised learning on graphs
  - Based on a simple architecture: masked auto-encoder

# Homework 7: Graph SSL & Pre-training

- Comment on graph SSL & pre-training:
  - Due by 4<sup>th</sup> Sept.
  - No coding
  - Write your comments (Reading other papers if possible)
  - Post them to https://discuss.cogdl.ai/t/topic/96
  - Discuss with others
  - Send your comments and discussions (via screenshots) to our email
- Reminder: Homework 1~6 & course proposal



# Thank you!

#### **Collaborators:**

Zhenyu Hou, Yuxiao Dong, Jie Tang, et al. (THU) Qingfei Zhao, Xinije Zhang, Peng Zhang, et al. (Zhipu Al) Hongxiao Yang, Chang Zhou, et al. (Alibaba)

Yang Yang (ZJU)

Yukuo Cen, KEG, Tsinghua U. Online Discussion Forum https://github.com/THUDM/cogdl https://discuss.cogdl.ai/