# Unlocking the Potential of Black-box Pre-trained GNNs for Graph Few-shot Learning

Qiannan Zhang[1], Shichao Pei[2], Yuan Fang[3], Xiangliang Zhang[4]

[1]Cornell University
[2]University of Massachusetts Boston
[3]Singapore Management University
[4]University of Notre Dame

## Motivation

- Graph few-shot learning deals with label scarcity on graphs.
- Pre-trained GNNs have promoted the advancement of graph few-shot learning by providing rich graph knowledge via large-scale label-free training.
- For economic and security considerations, pre-trained GNNs might be accessed only via Model-as-a-Service (Black-box).
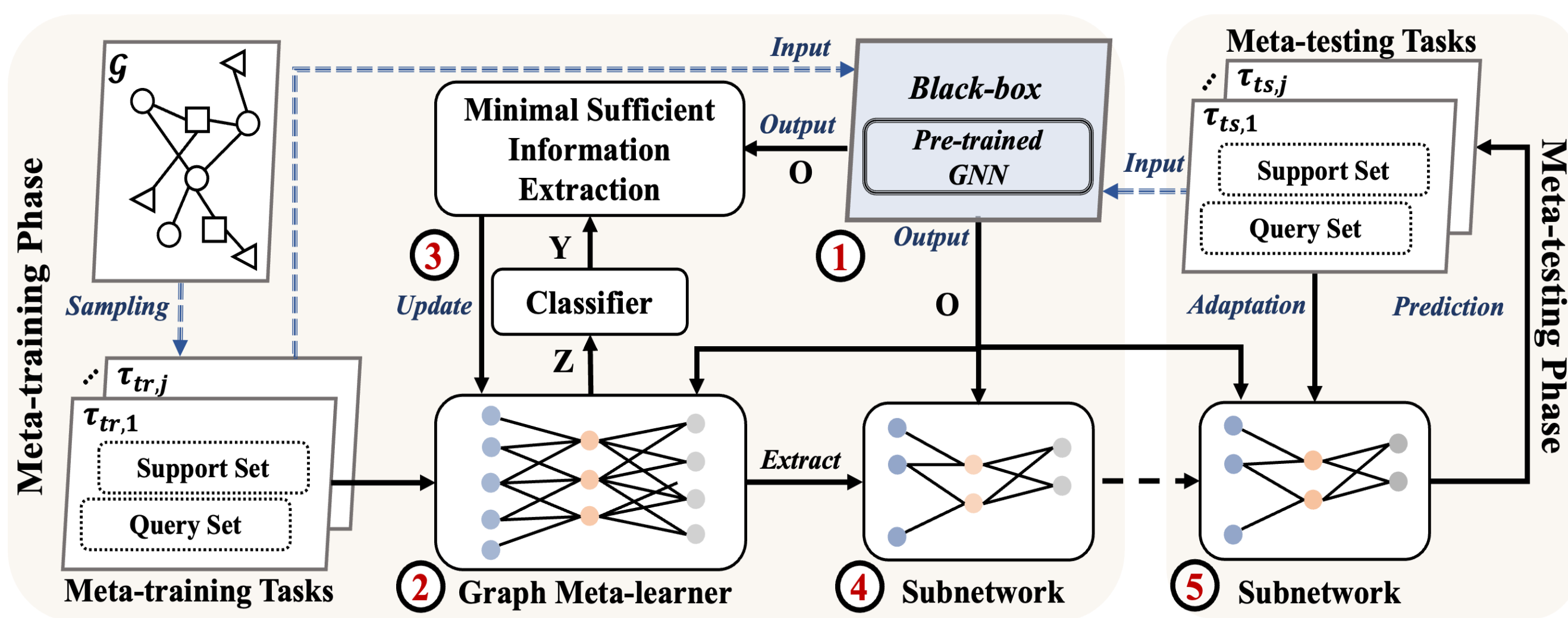
### Challenges

- When applying pre-trained GNNs in few-shot learning:
  - ❖ *Task Gap:* Pre-training objectives may include task-irrelevant information for downstream few-shot tasks.
  - ❖ *Black-box Setting:* Conventional fine-tuning isn't applicable without access to internal gradients.
  - ❖ *Overfitting:* Limited meta-training tasks can lead to memorization and poor generalization

### Key Idea

- Leverage a lightweight graph meta-learner, enabling:
  - ✓ Capturing just task-relevant knowledge needed for downstream node classification
  - ✓ Leveraging outputs from black-box pre-trained GNNs, no access to the gradients or parameters
  - ✓ Fast adapting to novel few-shot tasks with a compact subnetwork.

## Graph Meta-learning with Black-box Pre-trained GNNs (Meta-BP)



## Experiments

- Evaluated on four real-world graph datasets
- For node classification under various N-way K-shot settings

| Methods | Cora | Computers | Cora-full | |
|---|---|---|---|---|
| | 2-way | 3-way | 5-way | 10-way |
| | | 1-shot | | |
| GCN | $55.21_{(\pm5.64)}$ | $37.33_{(\pm3.91)}$ | $43.75_{(\pm2.92)}$ | $31.26_{(\pm3.29)}$ |
| GraphSage | $58.33_{(\pm5.22)}$ | $39.98_{(\pm5.17)}$ | $44.26_{(\pm2.64)}$ | $32.54_{(\pm3.53)}$ |
| GMI | $60.25_{(\pm4.32)}$ | $62.28_{(\pm4.92)}$ | $56.81_{(\pm1.42)}$ | $40.98_{(\pm1.72)}$ |
| DGI | $61.56_{(\pm4.46)}$ | $64.52_{(\pm5.03)}$ | $56.52_{(\pm1.53)}$ | $40.36_{(\pm1.76)}$ |
| PN | $52.60_{(\pm5.23)}$ | $47.63_{(\pm5.23)}$ | $48.25_{(\pm1.80)}$ | $35.65_{(\pm1.98)}$ |
| MAML | $53.66_{(\pm4.92)}$ | $66.05_{(\pm5.16)}$ | $58.38_{(\pm2.01)}$ | $38.72_{(\pm2.19)}$ |
| Meta-SGC | $57.72_{(\pm5.99)}$ | $67.40_{(\pm6.79)}$ | $61.34_{(\pm4.53)}$ | $41.29_{(\pm4.06)}$ |
| GPN | $57.22_{(\pm4.20)}$ | $63.78_{(\pm5.27)}$ | $54.36_{(\pm2.29)}$ | $43.27_{(\pm1.92)}$ |
| G-Meta | $62.24_{(\pm4.93)}$ | $67.22_{(\pm5.61)}$ | $55.21_{(\pm2.15)}$ | $46.23_{(\pm1.79)}$ |
| TLP | $61.11_{(\pm3.14)}$ | $65.05_{(\pm5.40)}$ | $61.28_{(\pm2.41)}$ | $48.12_{(\pm1.53)}$ |
| TENT | $61.25_{(\pm5.15)}$ | $66.24_{(\pm5.24)}$ | $62.42_{(\pm2.16)}$ | $47.95_{(\pm1.88)}$ |
| TEG | $\underline{63.14}_{(\pm4.43)}$ | $\underline{67.58}_{(\pm5.11)}$ | $\underline{63.73}_{(\pm2.08)}$ | $\underline{48.36}_{(\pm2.03)}$ |
| Meta-BP | $\mathbf{66.38}_{(\pm5.01)}$ | $\mathbf{69.35}_{(\pm4.28)}$ | $\mathbf{66.05}_{(\pm1.46)}$ | $\mathbf{51.41}_{(\pm1.91)}$ |

### ② Graph Meta-Learner (GML)

- Receives output representations (denoted as **O**) from the black-box pre-trained GNN as the initial embeddings of GML
- Pre-computes *neighbor abstractions* with graph topology

$$\mathbf{h}_{\mathcal{N}_v} = \frac{1}{|\mathcal{N}_v|}\sum_{u\in\mathcal{N}_v}\mathbf{h}_u$$

- Fuses node features and neighbor abstractions by linear transformation

$$\mathbf{z}_v = \sigma([\mathbf{h}_v||\mathbf{h}_{\mathcal{N}_v}]\mathbf{W})$$

### ④ Meta-Learner Pruning

- Extracts a *sparse subnetwork* from GML based on *capacity ratio c*
  - deals with overparameterization and overfitting the insufficient meta-training tasks
  - inspired by the *Lottery Ticket Hypothesis* to learn a binary mask of GML parameters

$$\mathcal{L}_S = \min_{\mathbf{m}}\frac{1}{B}\sum_{j=1}^{B}\{\frac{1}{b}\sum_{i=1}^{b}(\mathcal{L}_{CE}(f_C(\mathrm{GML}(\mathbf{o}_i,\mathcal{G};\phi_j\odot\mathbf{m})),\hat{y}_i) - C_j)\}.$$

$$C_j = \mathcal{L}_{CE}(f_C(\mathrm{GML}(\mathbf{o}_i,\mathcal{G};\phi_j)),\hat{y}_i)$$

$$\text{subject to}\quad |\mathbf{m}^*| \le c\cdot|\phi|$$

### ③ Extracting Minimal Sufficient Information

- Outputs include both task-relevant and task-irrelevant information.
- *Information bottleneck:* learning node representations **Z** to be:
  - a *minimal* knowledge (simplest mapping) from the pre-trained GNN
  - yet *sufficient* to adapt well to meta-tasks

$$\mathcal{L}_I = \min_{\mathbf{Z}\sim\mathrm{GML}(\cdot)} I(\mathbf{O};\mathbf{Z}) - \beta I(\mathbf{Z};\mathbf{Y})$$
$$= \min_{\mathbf{Z}\sim\mathrm{GML}(\cdot)} I(\mathbf{O};\mathbf{Z}) - \beta H(\mathbf{Y}) + \beta H(\mathbf{Y}|\mathbf{Z})$$

constant classification loss

- *Neural estimation:* using neural networks to efficiently estimate I(**O**, **Z**) alongside meta-optimization

$$I(\mathbf{O};\mathbf{Z}) \approx \frac{1}{b}\sum_{i=1}^{b}T_\theta(\mathbf{o}_i,\mathbf{z}_i) - \log(\frac{1}{b}\sum_{i=1}^{b}e^{T_\theta(\mathbf{o}_i,\overline{\mathbf{z}}_i)})$$

meta-task samples  shuffled order

### Optimization

- $\mathcal{L}_I$: Minimal sufficient *information extraction* during meta-optimization
- $\mathcal{L}_S$: *Subnetwork* extraction during meta-optimization
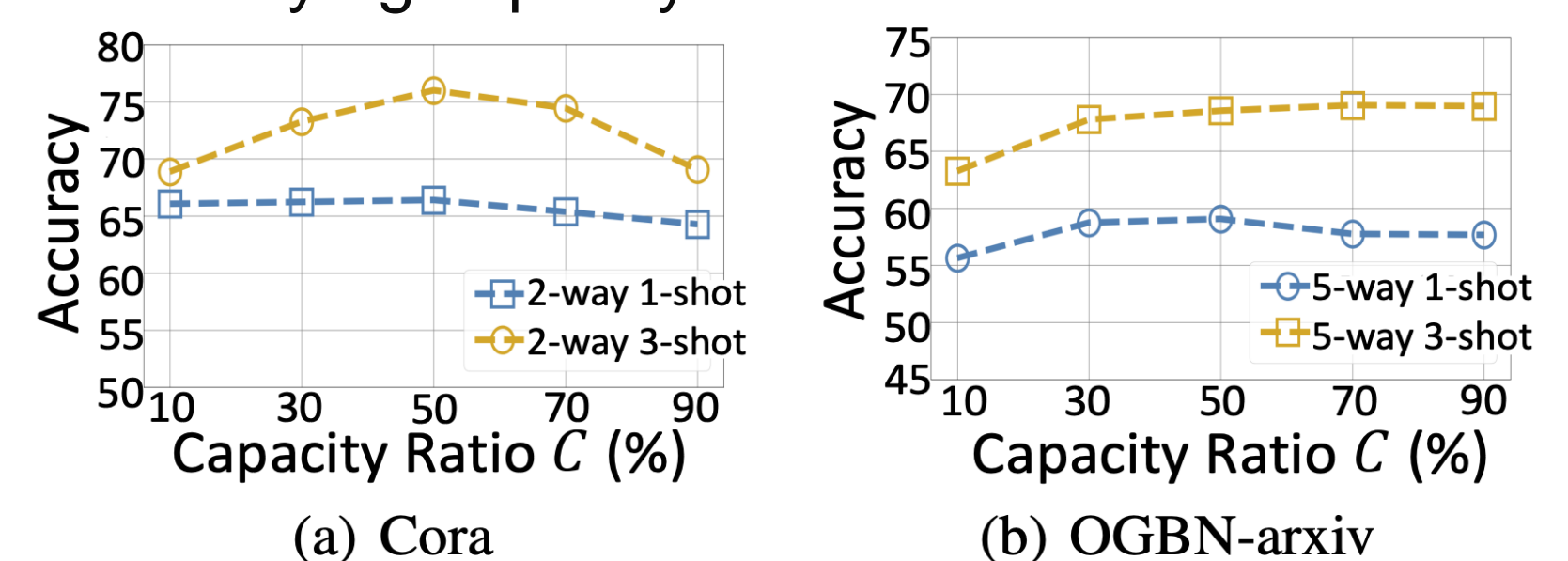- Total loss:

$$\mathcal{L}_{Meta} = \mathcal{L}_I + \alpha\mathcal{L}_S$$

*Parameter sensitivity*: training Meta-BP with varying capacity ratios



(a) Cora  (b) OGBN-arxiv

*Pretrained model impact:* performance varies with different pre-trained backbones

| Methods | Cora-full | | | |
|---|---|---|---|---|
| | 5-way 1-shot | 5-way 3-shot | 10-way 1-shot | 10-way 3-shot |
| Meta-BP-GMI | $66.13_{(\pm1.48)}$ | $75.58_{(\pm1.23)}$ | $51.88_{(\pm1.44)}$ | $62.28_{(\pm1.86)}$ |
| Meta-BP-BGRL | $65.92_{(\pm1.62)}$ | $73.64_{(\pm1.42)}$ | $51.36_{(\pm1.65)}$ | $59.47_{(\pm1.94)}$ |
| Meta-BP-DGI | $66.05_{(\pm1.46)}$ | $72.98_{(\pm1.86)}$ | $51.41_{(\pm1.91)}$ | $57.79_{(\pm2.16)}$ |

*Parameter efficiency*: the number of parameters shows the extracted final subnetwork for adaptation is much smaller

| Methods | Meta-PreGNN | Meta-BP |
|---|---|---|
| **Computers** | **3-way** | |
| **#Total** params | 263.16 K | 313.56 K |
| **#Trainable Params** in meta-training | 263.16 K | 50.40 k |
| **#Params of GML** | - | 8.24 K |
| **#Tunable params** in meta-testing | 263.16 K | 2.47 K |
| **Ratio of tunable params** w.r.t. Meta-PreGNN | 100% | 0.94% |