

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ The Instability of GFM Fine-tuning

- Observations from SOTA GFMs: **huge variances** under few-shot fine-tuning

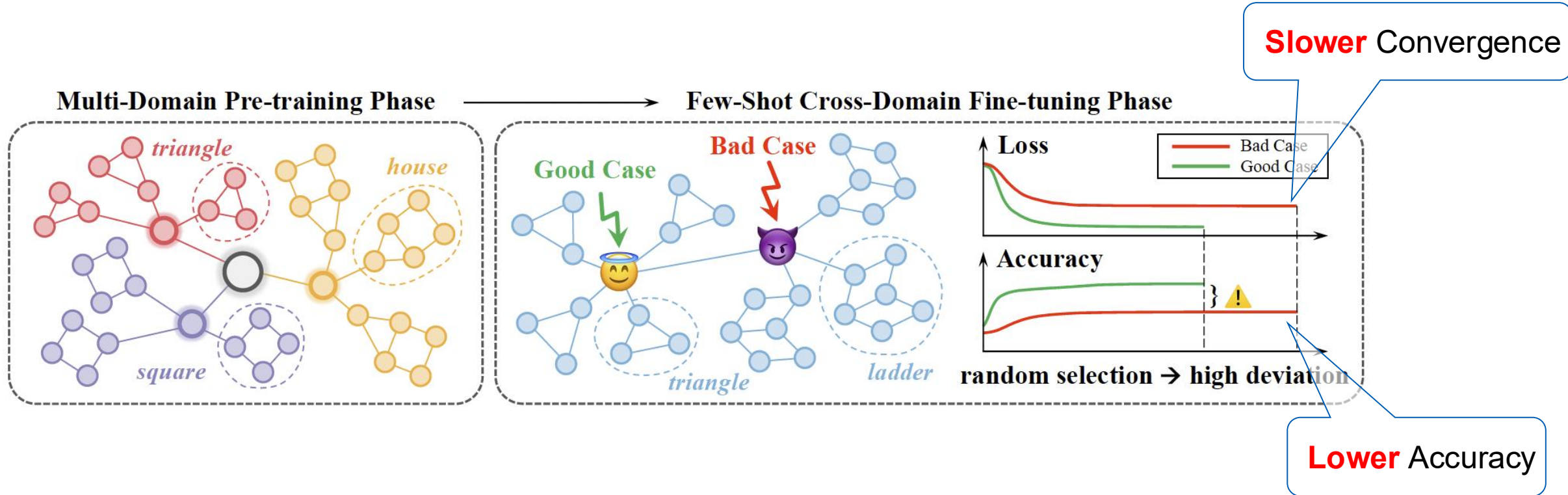
42.26 ± 10.18	42.40 ± 9.26	49.82 ± 8.38	64.82 ± 10.53	49.77 ± 11.00	67.80 ± 6.84	
48.36 ± 11.34	44.28 ± 10.16	54.34 ± 9.76	64.08 ± 9.85	48.29 ± 11.35	68.92 ± 7.32	
47.80 ± 11.88	36.38 ± 9.10	50.25 ± 10.43	58.71 ± 8.69	48.22 ± 8.17	42.70 ± 8.73	33.36 ± 8.11
55.35 ± 13.62	38.75 ± 9.40	48.69 ± 10.16	58.75 ± 11.67	48.72 ± 11.18	43.71 ± 9.54	48.28 ± 9.72
34.23 ± 8.16	39.05 ± 8.82	44.85 ± 6.72	34.02 ± 11.94	22.46 ± 1.96	24.61 ± 3.99	50.79 ± 0.65
39.54 ± 9.02	39.24 ± 8.95	45.39 ± 11.01	33.58 ± 10.38	22.35 ± 3.77	23.68 ± 1.56	50.78 ± 3.05
44.83 ± 7.41	42.18 ± 6.41	46.84 ± 7.31	40.77 ± 5.96	24.30 ± 3.26	28.36 ± 3.65	52.36 ± 0.86
64.56 ± 7.29	61.24 ± 4.82	63.50 ± 5.81	49.56 ± 6.92	23.00 ± 4.39	30.54 ± 2.87	53.58 ± 0.83

- **Why?** Random few-shot sample selection brings fine-tuning uncertainty.

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ The Instability of GFM Fine-tuning

- How does fine-tuning instability happen?

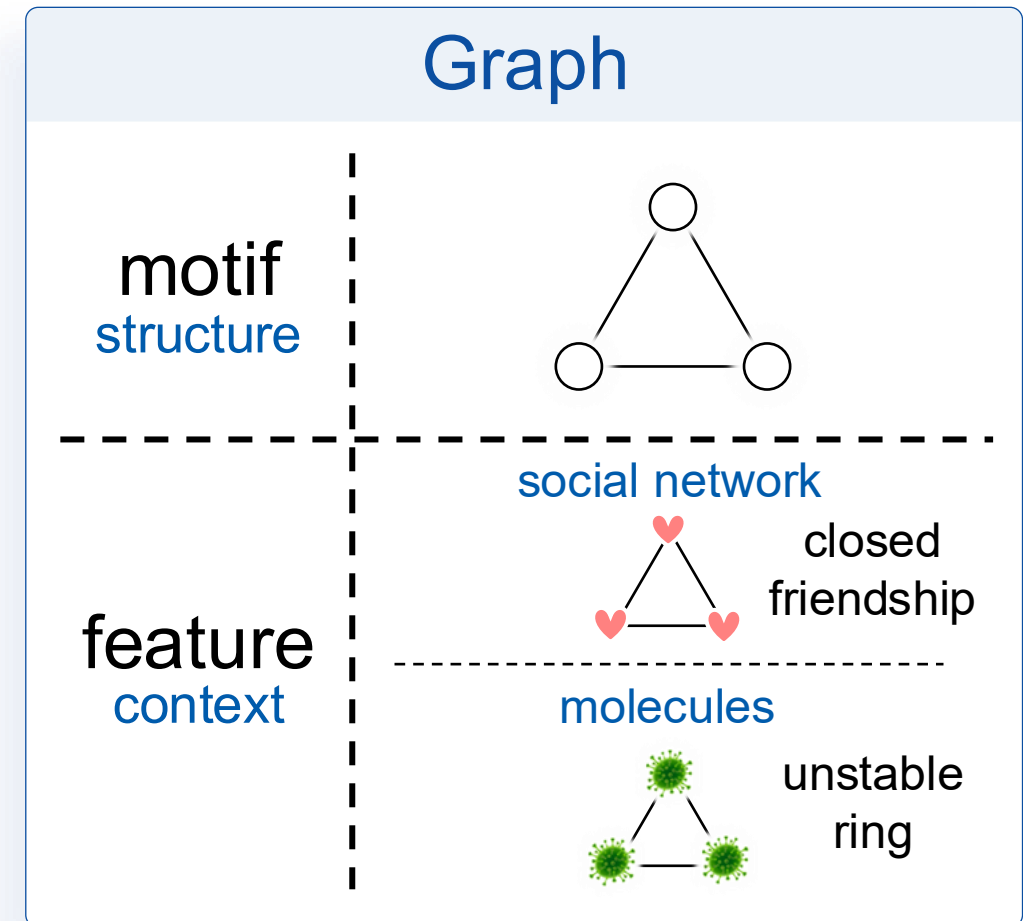
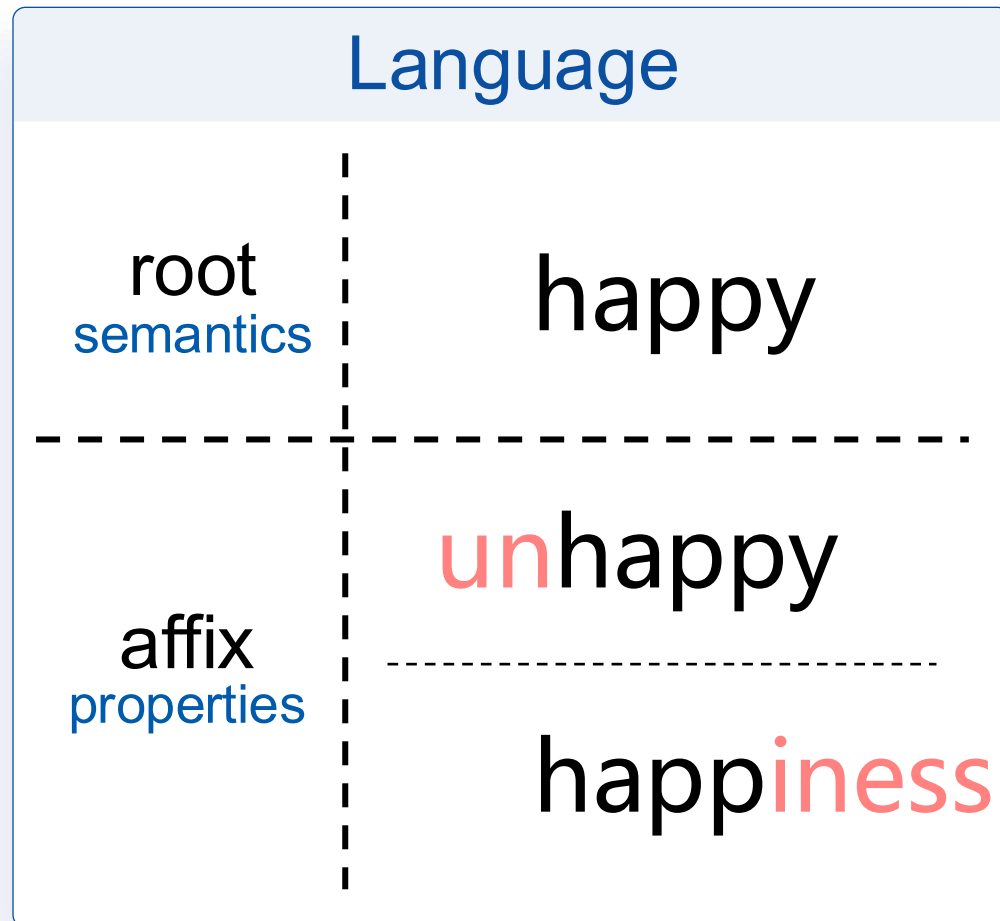


- **Challenge: mismatch** between support samples and pre-trained graph patterns.

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

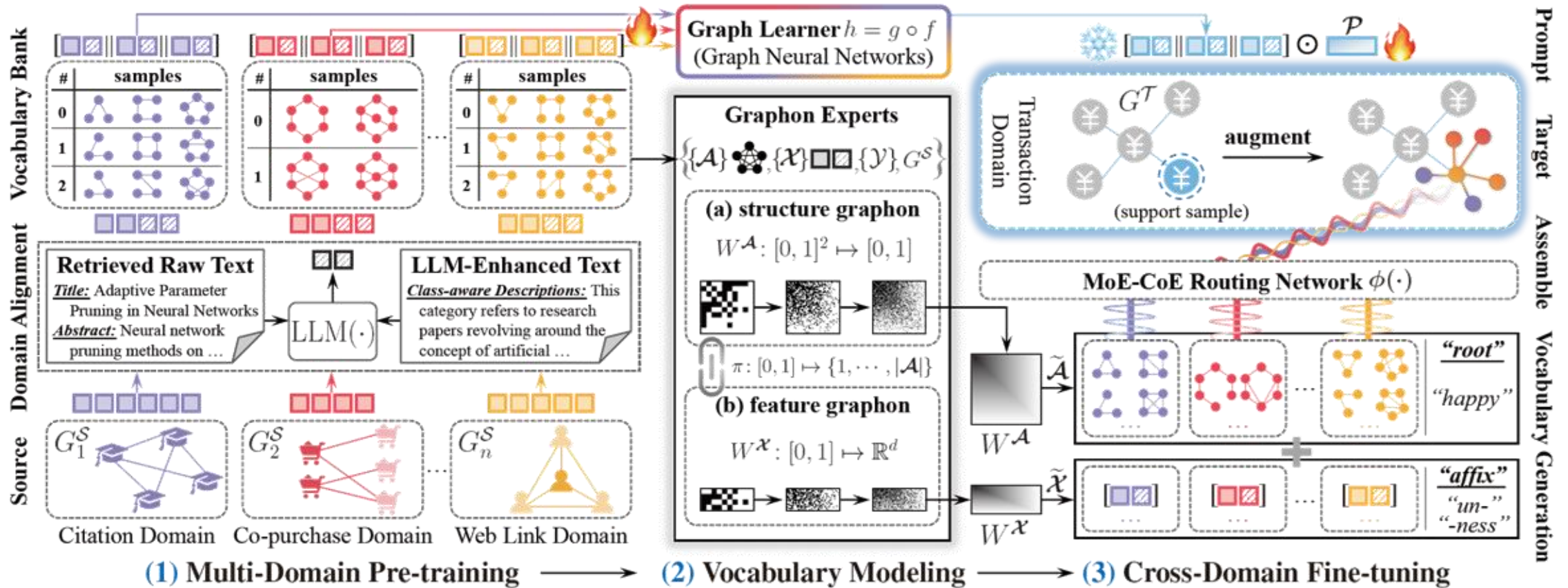
■ Key Insight: Tokenize Graphs with Vocabulary

- Understanding graph data through the lens of language modeling



GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ GRAVER: Pretrain-then-Finetune → Pretrain-then-**Augment**-then-Finetune



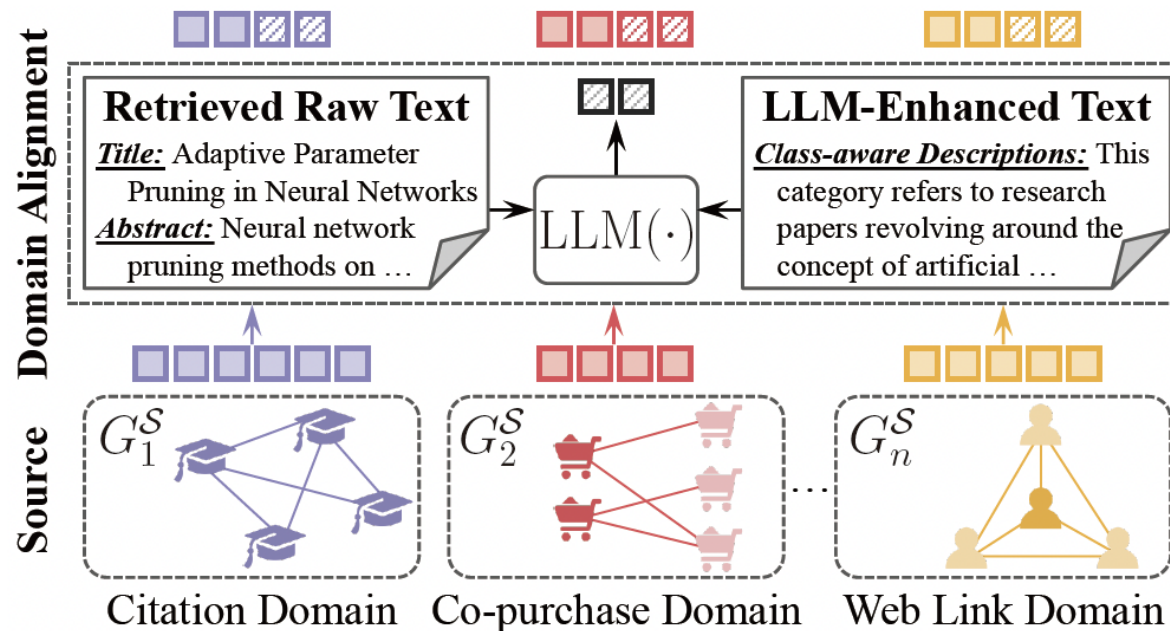
GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 1: Pre-training with Transferable Graph Vocabulary

□ **Goal:** identify subgraph (graph vocabulary) that generalize across tasks and domains.

□ Step 1: Multi-Domain Alignment

$$\hat{\mathbf{X}}_i^S \in \mathbb{R}^{N_i \times d} = \mathbf{W}_i^\top (\mathcal{F}(\mathbf{X}_i^S \parallel \text{LLM}(\mathbf{X}_i^S))), \quad \text{for all } G_i^S \in \{G^S\}$$



GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 1: Pre-training with Transferable Graph Vocabulary

□ **Goal:** identify subgraph (graph vocabulary) that generalize across tasks and domains.

□ Step 2: Factor-aware Ego-Graph Disentanglement

□ multi-channel encoder:

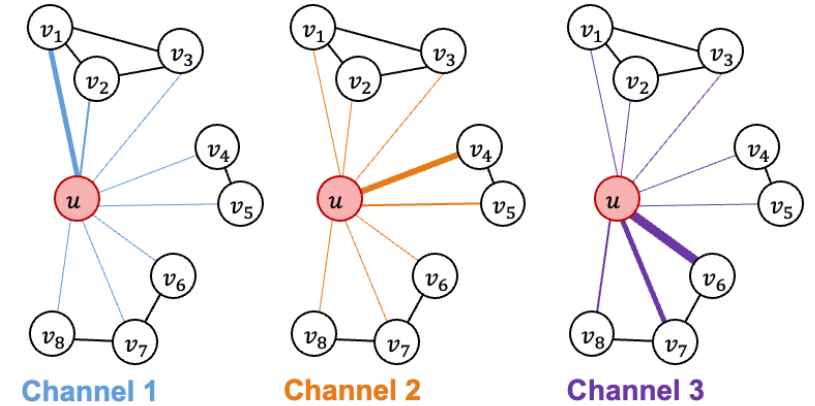
$$f_{\Theta}: \mathbf{g}_u^{\mathcal{S}} \mapsto \{\mathbf{h}_{u,k}^{\mathcal{S}} \in \mathbb{R}^d\}_{k=1}^K$$

□ neighbor attention:

$$\alpha_{v \rightarrow k}^{(t)} \propto \text{Softmax}_k \left(\langle \mathbf{h}_{u,k}^{\mathcal{S}(t)}, \mathbf{h}_{v,k}^{\mathcal{S}(t)} \rangle / \tau \right), \quad \text{s.t.} \quad \sum_{k=1}^K \alpha_{v \rightarrow k}^{(t)} = 1$$

□ updating:

$$\mathbf{h}_u^{\mathcal{S}(T)} = \parallel_{k=1}^K \mathbf{h}_{u,k}^{\mathcal{S}(T)}, \quad \mathbf{h}_{u,k}^{\mathcal{S}(t+1)} := \mathbf{h}_{u,k}^{\mathcal{S}(t)} + \sum_{v \in \mathcal{V}_{u,k}^{(t)}} \left[\alpha_{v \rightarrow k}^{(t)} \mathbf{h}_{v,k}^{\mathcal{S}(t)} \right]$$



GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 1: Pre-training with Transferable Graph Vocabulary

□ **Goal:** identify subgraph (graph vocabulary) that generalize across tasks and domains.

□ Step 3: Semantic-Independence Promotion

□ regularizer: $\mathcal{R}_{\text{MI}}^u = \sum_{i \neq j} I(\mathbf{h}_{u,i}^S; \mathbf{h}_{u,j}^S) \stackrel{\text{Proof C.1}}{\leq} \sum_{i \neq j} \mathbb{E}_{u \in \mathcal{B}} \left[-\log \left(\text{Softmax}_v (\langle \mathbf{h}_{u,i}^S, \mathbf{h}_{v,j}^S \rangle / \tau) \Big|_{v=u} \right) \right]$

□ Why it work?

Proposition 1 (Vocabulary Transferability). Given any nodes u, v , and assume $\|\hat{\mathbf{x}}_u^S - \hat{\mathbf{x}}_v^S\|_2 \leq \epsilon$. The semantic discrepancies $\Delta = \|f(\mathbf{g}_u^S) - f(\mathbf{g}_v^S)\|_2$ is upper bounded by:

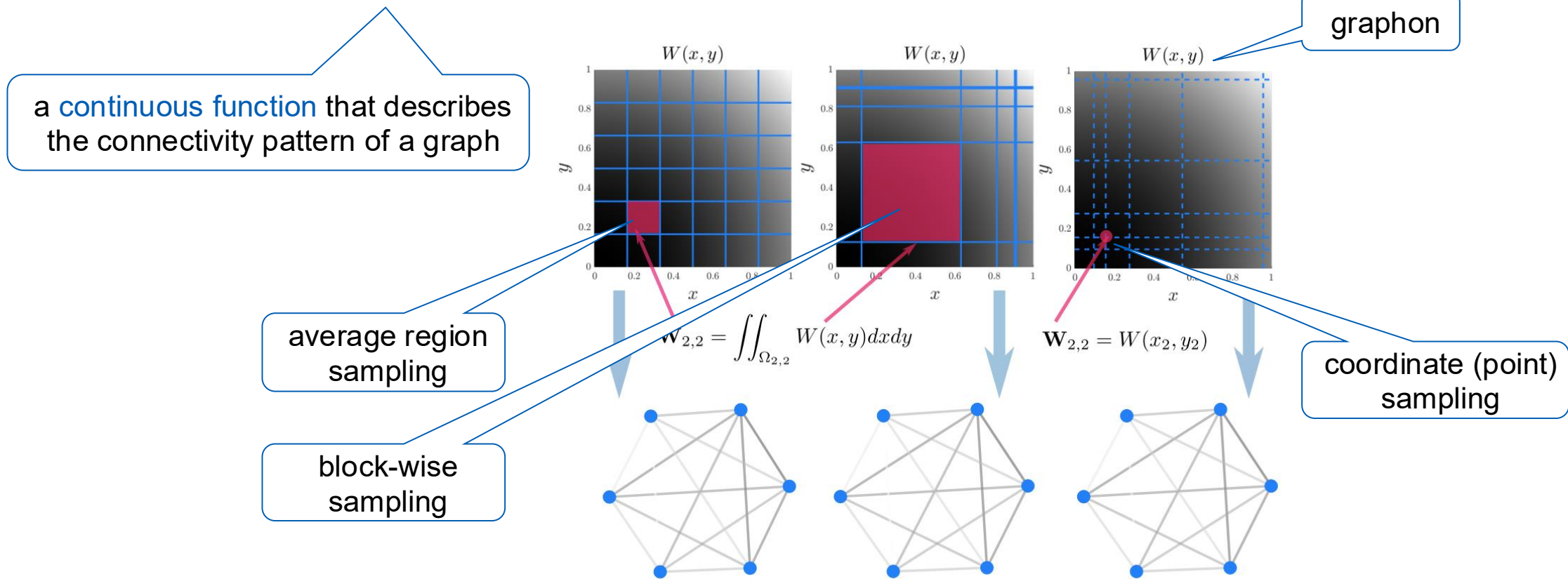
$$\Delta \stackrel{[a] \text{ Proof C.2}}{\leq} \left[\min_{\Pi \in S_K} \sum_{k=1}^K \left\| \mathbf{h}_{u,k}^S - \mathbf{h}_{v,\Pi(k)}^S \right\|_2^2 + \psi(\mathcal{R}_{\text{MI}}^{u,v}) \right]^{\frac{1}{2}} \stackrel{[b] \text{ Proof C.3}}{\leq} \epsilon \sqrt{K} \left(\frac{C_\sigma L_{\mathbf{W}} L_s}{4\rho\tau} \right)^T, \quad (7)$$

Key Takeaway: Graph vocabulary is the **smallest & indivisible** unit that preserves structural semantics.

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

Phase 2: Vocabulary Modeling with Generative Graphon Experts

- Goal: generate transferable, class-consistent graph vocabularies from structure-feature **graphon** distributions (establish a “vocabulary bank”).



[Picture Credit] Graphon Pooling in Graph Neural Networks. *Alejandro Parada-Mayorga, et al.*

Haonan Yuan, Qingyun Sun, et al. GRAVER: Generative Graph Vocabularies for Robust Graph Foundation Models Fine-tuning, NeurIPS 2025.

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 2: Vocabulary Modeling with Generative Graphon Experts

□ **Goal:** generate transferable, class-consistent graph vocabularies from structure-feature **graphon** distributions (establish a “**vocabulary bank**”).

□ Step 1: Structure Token Modeling

□ structure token graphon approximation:

$$W_c^{\mathcal{A}}: [0, 1]^2 \mapsto [0, 1], \quad W_c^{\mathcal{A}}(u, v) = \frac{1}{N_c} \sum_{i=1}^{N_c} \mathcal{A}_i^{(c)}[\pi_i(u), \pi_i(v)]$$
$$\pi_i: [0, 1] \mapsto \{1, \dots, |\mathcal{A}_i|\}$$

□ Step 2: Feature Token Modeling

□ feature token graphon approximation:

$$W_c^{\mathcal{X}}: [0, 1] \mapsto \mathbb{R}^d, \quad W_c^{\mathcal{X}}(u) = \frac{1}{N_c} \sum_{i=1}^{N_c} \mathcal{X}_i^{(c)}[\pi_i(u), :]$$

shared

non-parametric

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 2: Vocabulary Modeling with Generative Graphon Experts

□ **Goal:** generate transferable, class-consistent graph vocabularies from structure-feature **graphon** distributions (establish a “**vocabulary bank**”).

□ Step 3: Conditional Vocabulary Generation

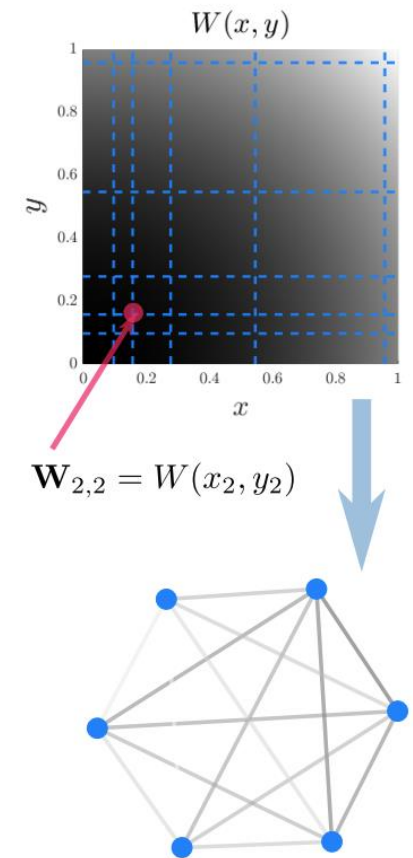
□ **Principle: Root (structures) + Affix (features)**

□ sample and generation (coordinate sampling):

$$\tilde{\mathcal{A}}[i, j] \sim \text{Bern} (W_c^{\mathcal{A}}(\mathbf{u}_i, \mathbf{u}_j))$$

$$\tilde{\mathcal{X}}[i, :] = W_c^{\mathcal{X}}(\mathbf{u}_i)$$

$$\mathbf{u}_1, \dots, \mathbf{u}_{n'} \sim \mathcal{U}[0, 1]$$



GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 2: Vocabulary Modeling with Generative Graphon Experts

□ **Goal:** generate transferable, class-consistent graph vocabularies from structure-feature **graphon** distributions (establish a “**vocabulary bank**”).

□ Why it work?

Proposition 2 (Generation Distributional Convergence). Let $\mathbf{g}_c^{\text{emp}}$ be the empirical distribution over n' -node vocabulary subgraphs of class c , collected from the disentangled vocabulary bank:

$$\mathbf{g}_c^{\text{emp}} := \frac{1}{N_c} \sum_{i=1}^{N_c} \delta_{(\mathcal{A}_i^{(c)}, \mathbf{x}_i^{(c)})}, \quad \text{each } (\mathcal{A}_i^{(c)}, \mathbf{x}_i^{(c)}) \in \{0, 1\}^{n' \times n'} \times \mathbb{R}^{n' \times d}, \quad (12)$$

where $\delta_{(\cdot)}$ denotes Dirac measure, representing a point mass probability distribution. Assume the true underlying structure and feature functions $W_c^{\mathcal{A}^*}, W_c^{\mathbf{x}^*}$ are bounded, Lipschitz, and satisfy permutation equivariance. If the estimators $W_c^{\mathcal{A}}, W_c^{\mathbf{x}}$ converge uniformly in L_∞ norm, then the total variation (TV) distance between $\mathbf{g}_c^{\text{gen}} = (\tilde{\mathcal{A}}_c, \tilde{\mathbf{x}}_c)$ and $\mathbf{g}_c^{\text{emp}}$ vanishes as $N_c \rightarrow \infty$:

$$\|\mathbf{g}_c^{\text{gen}} - \mathbf{g}_c^{\text{emp}}\|_{\text{TV}} \rightarrow 0. \quad (\text{Proof C.4}) \quad (13)$$

TVD: the **strongest** statistical distance measurement (> KL, JS, and Wasserstein)

Key Takeaway: The generated graph vocabularies **match the true** distribution (their gap converges to 0).

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Phase 3: Fine-tuning with Augmented Support Samples

- **Goal:** strengthen support samples with vocabulary-based structural augmentation.
(“supervised augmentation”)

□ MoE-CoE Network for Selective Augmentation

- MoE (Mixture-of-Experts): *where to route*
- CoE (Collaboration-of-Experts): *how to compose*

$$\mathbf{S}_M \in \mathbb{R}^n = \text{Softmax}(\mathbf{W}_M^\top \cdot \phi(\widehat{\mathbf{X}}_i^\top)), \quad \mathbf{S}_C \in \mathbb{R}^c = \text{Softmax}(\mathbf{W}_C^\top \cdot (\phi(\widehat{\mathbf{X}}_i^\top \parallel \tilde{\mathcal{X}})))$$

- vocabulary compose:

$$\tilde{\mathbf{g}}_i^{\text{gen}} \stackrel{\text{def}}{=} (\tilde{\mathcal{A}}_i^{\text{gen}}, \tilde{\mathcal{X}}_i^{\text{gen}}) = \sum_{i=1}^n \mathbf{S}_{M,i}^\top \left(\sum_{c=1}^c \mathbf{S}_{C,c}^\top \cdot \mathbf{g}_c^{\text{gen}} \right), \quad \tilde{G}_i^\top = G_i^\top \oplus \tilde{\mathbf{g}}_i^{\text{gen}}$$

- optimize: $\mathcal{L}_{\text{MoE-CoE}}(\mathbf{S}_M, \mathbf{S}_C) = - \sum_{i=1}^n \mathbf{S}_{M,i}^\top \log(\mathbf{S}_{M,i}) - n \sum_{c=1}^c \mathbf{S}_{C,c}^\top \log(\mathbf{S}_{C,c})$

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Experiment: One-shot Cross-Dataset/-Domain Transfer

background color denotes different dataset domain

Table 1: Accuracy (% \pm std for 20 runs) of **one-shot classification**. Best results are presented **bold** and the runner-ups are underlined. CR = Cora, CS = CiteSeer, PM = PubMed, arXiv = ogbn-arXiv, Tech = ogbn-tech, Home = ogbn-home, Wiki = Wiki-CS. Color denotes domain.

Source	Cross-Dataset								Cross-Domain						
	CS	PM	CR	PM	CR	CS	CR	CS	PM	CR	CS	CR	CS	Home	Tech
Model / Target	Home	Wiki	Home	Wiki	Home	Wiki	Home	Wiki	PM	Wiki	PM	Home	Wiki	arXiv	
Node Classification															
Vanilla GNNs	GCN (bb.) [40]	28.40 \pm 4.62	29.25 \pm 3.39	40.33 \pm 6.90	61.59 \pm 5.13	53.89 \pm 3.35	36.74 \pm 2.53	28.58 \pm 5.39							
	GAT [91]	29.72 \pm 5.17	29.31 \pm 3.47	40.51 \pm 3.95	61.98 \pm 5.23	51.70 \pm 3.96	36.24 \pm 4.19	29.03 \pm 5.75							
Self-supervised Graph Pre-training	GCC [69]	32.47 \pm 4.55	32.78 \pm 3.85	41.66 \pm 3.27	64.74 \pm 3.78	55.36 \pm 5.86	37.66 \pm 3.80	30.42 \pm 5.54							
	DGI [92]	30.77 \pm 3.92	31.41 \pm 4.11	39.97 \pm 5.90	63.76 \pm 3.77	53.31 \pm 2.58	39.20 \pm 5.67	32.15 \pm 4.91							
	GraphCL [116]	33.64 \pm 5.75	28.20 \pm 3.13	39.03 \pm 8.67	62.44 \pm 6.55	51.55 \pm 8.20	38.05 \pm 3.30	31.81 \pm 5.04							
	DSSL [102]	29.76 \pm 4.55	30.84 \pm 7.62	39.99 \pm 6.29	61.27 \pm 5.21	51.90 \pm 6.52	37.58 \pm 6.78	28.14 \pm 5.17							
	GraphACL [103]	35.92 \pm 5.61	33.88 \pm 6.69	42.73 \pm 5.26	66.70 \pm 4.29	58.01 \pm 6.32	40.94 \pm 6.65	35.57 \pm 4.29							
Prompt-based Graph Fine-tuning	GPPT [82]	32.38 \pm 5.74	31.78 \pm 4.61	41.97 \pm 5.48	65.81 \pm 6.56	57.81 \pm 5.36	40.97 \pm 5.41	34.22 \pm 6.51							
	GraphPrompt [58]	38.02 \pm 6.52	34.57 \pm 6.34	46.19 \pm 7.63	70.11 \pm 6.62	60.99 \pm 7.19	44.02 \pm 6.74	40.74 \pm 6.56							
	GPF [14]	40.01 \pm 8.21	40.07 \pm 7.64	47.58 \pm 5.83	70.07 \pm 6.55	59.05 \pm 5.85	44.62 \pm 8.68	40.13 \pm 5.07							
	ProNoG [121]	43.67 \pm 6.33	40.34 \pm 7.11	51.35 \pm 7.16	73.49 \pm 7.34	62.77 \pm 6.10	45.65 \pm 6.44	41.15 \pm 4.21							
Multi-domain GFMs	GCOPE [136]	36.27 \pm 3.93	40.42 \pm 4.64	44.75 \pm 4.67	71.26 \pm 5.95	60.39 \pm 6.08	41.80 \pm 4.96	40.00 \pm 6.53							
	MDGPT [123]	42.55 \pm 6.84	37.92 \pm 7.18	51.03 \pm 8.99	72.10 \pm 7.12	62.69 \pm 7.29	45.93 \pm 6.24	43.33 \pm 5.75							
	SAMGPT [119]	46.79 \pm 6.54	38.65 \pm 6.35	51.92 \pm 9.50	73.60 \pm 7.55	64.32 \pm 7.02	46.03 \pm 6.98	45.27 \pm 5.05							
	GRAVER (ours)	48.00 \pm 2.52	41.13 \pm 3.00	55.30 \pm 3.03	77.62 \pm 2.94	67.12 \pm 3.18	49.33 \pm 2.44	49.25 \pm 3.43							

the **best** performance across all 7 target datasets

the **best average** score, the **lower variation** (more robust / stable)

more **challenging** setting, remains **the best**

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Experiment: Few-shot Node and Graph Classification

consistently **outperforms** all baselines across different m-shot settings.

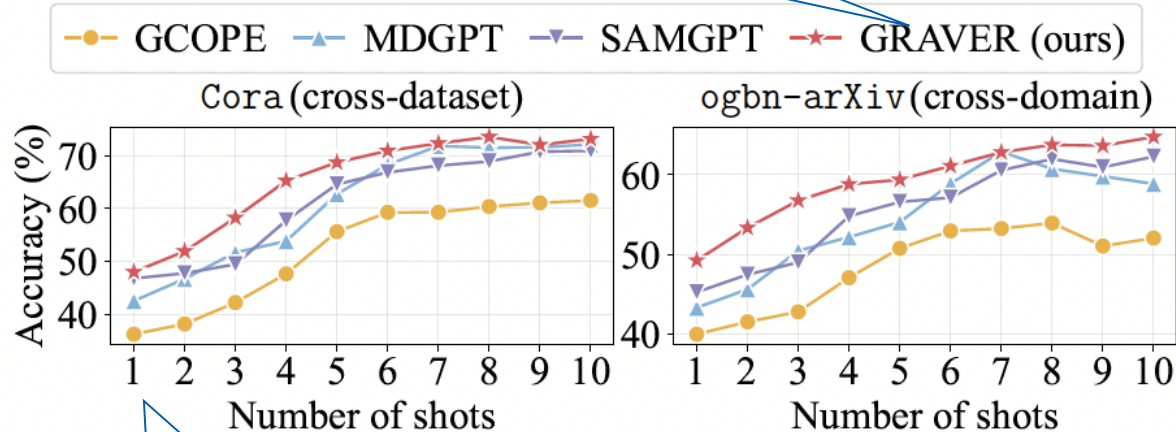


Figure 3: m -shot node classification.

the advantage is **more pronounced** when m is very small.

even performs **better** as the noise becomes **stronger**.

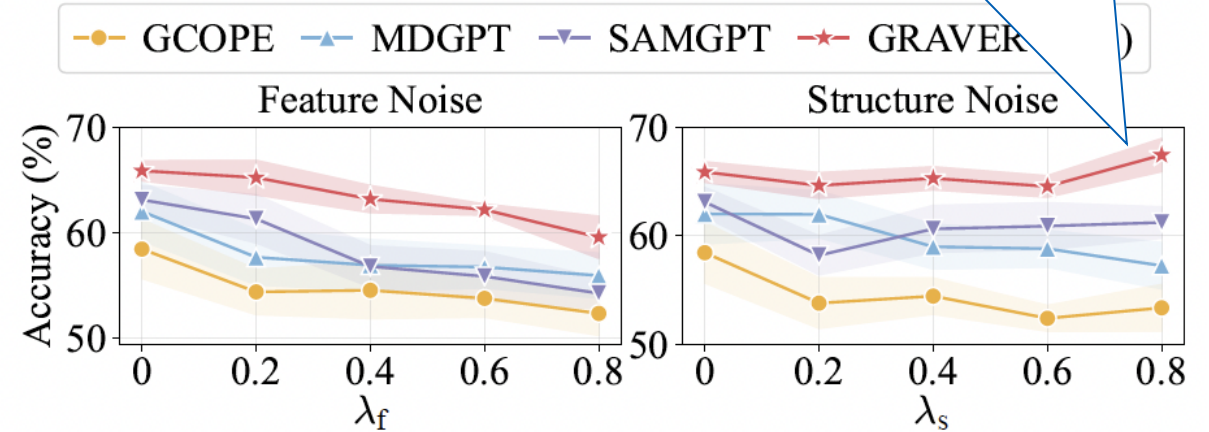


Figure 4: 5-shot graph classification (Wiki-CS).

remains **robust** under both **feature** noise and **structure** noise.

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

Experiment: Ablation Studies

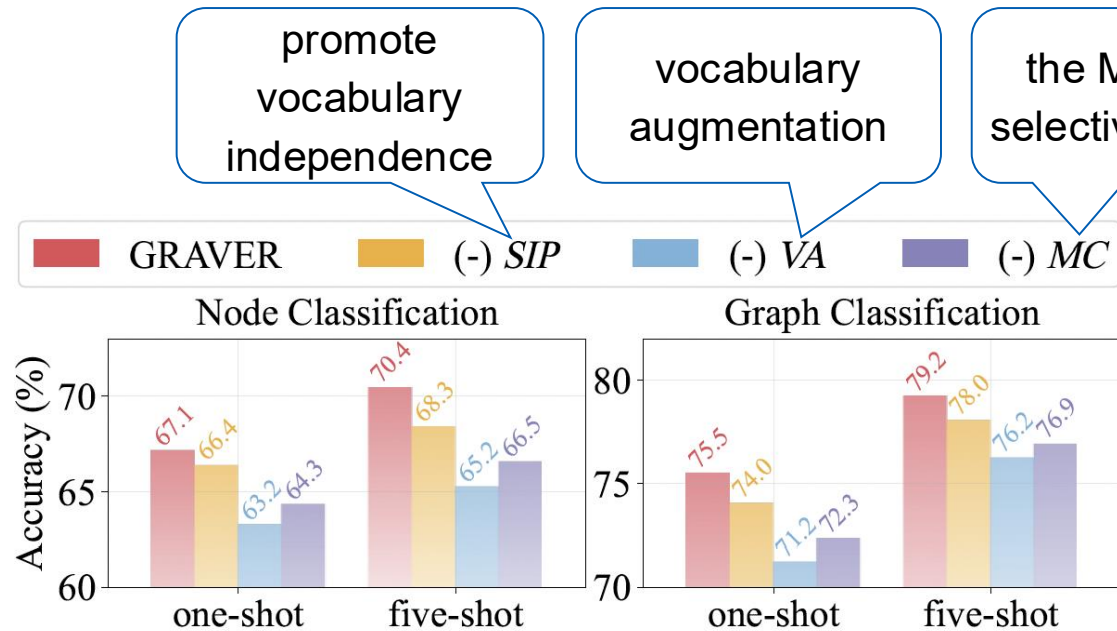


Figure 5: Ablation studies (ogbn-Home).

- All three components contribute to overall performance.
- **Importance:** VA > MC > SIP

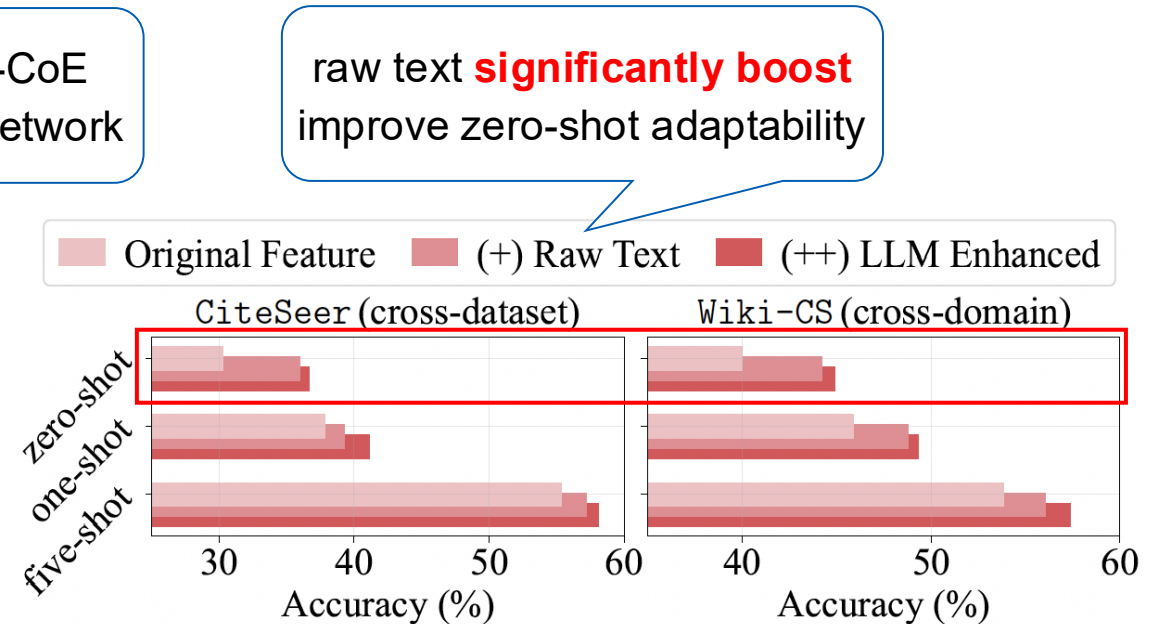


Figure 7: Analysis on LLM (node classification).

more effective under the cross-domain settings (more challenging)

GRAVER: Generative Graph Vocabularies for Robust GFM Fine-tuning

■ Experiment: Efficiency and Sensitivity

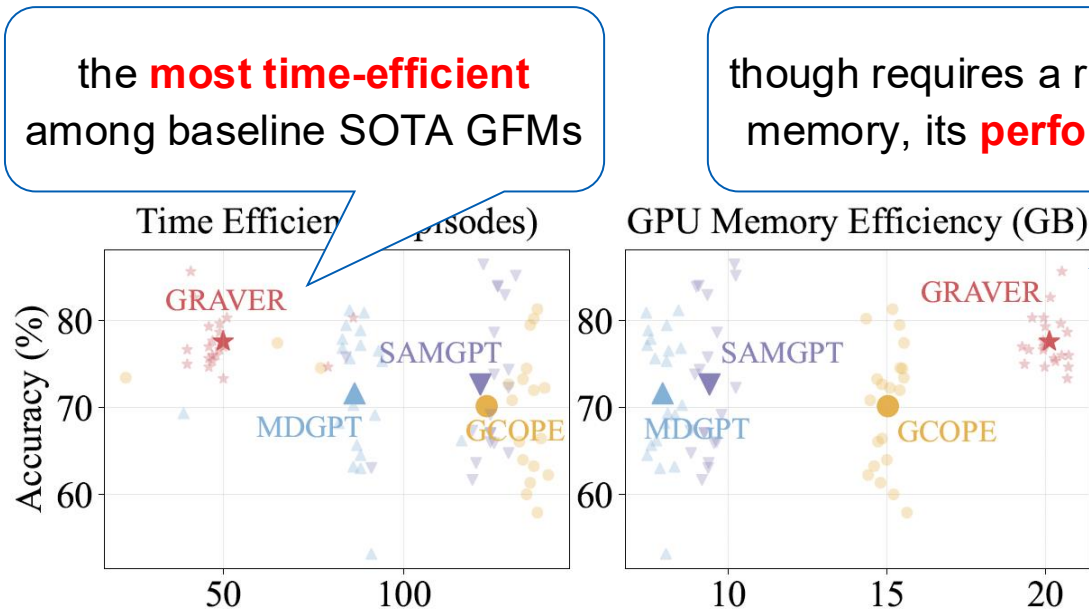


Figure 6: Efficiency analysis (ogbn-Tech).

GRAVER offers the best **accuracy-memory trade-off.**

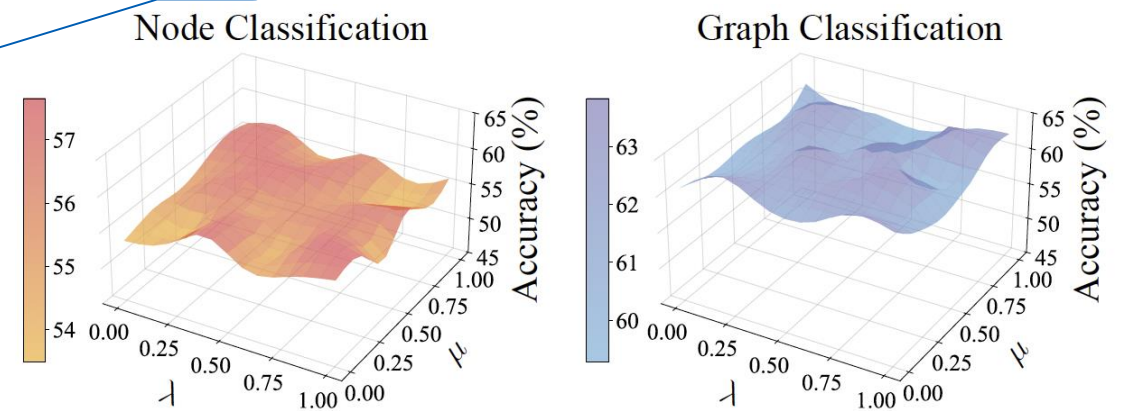


Figure 8: Hyperparameter sensitivity (PubMed).

(To ensure continuity, interpolation smoothing is applied.)

GRAVER exhibits low hyperparameter sensitivity
→ **easy to fine-tune** in practice