

## Motivation

- Realistic Graphs often display non-uniform patterns such as local homophily or heterophily.
- Most GNNs overlook these variations since they focus on *global* properties of the graph.
- Node-specific adaptations could boost performance.

## Mowst

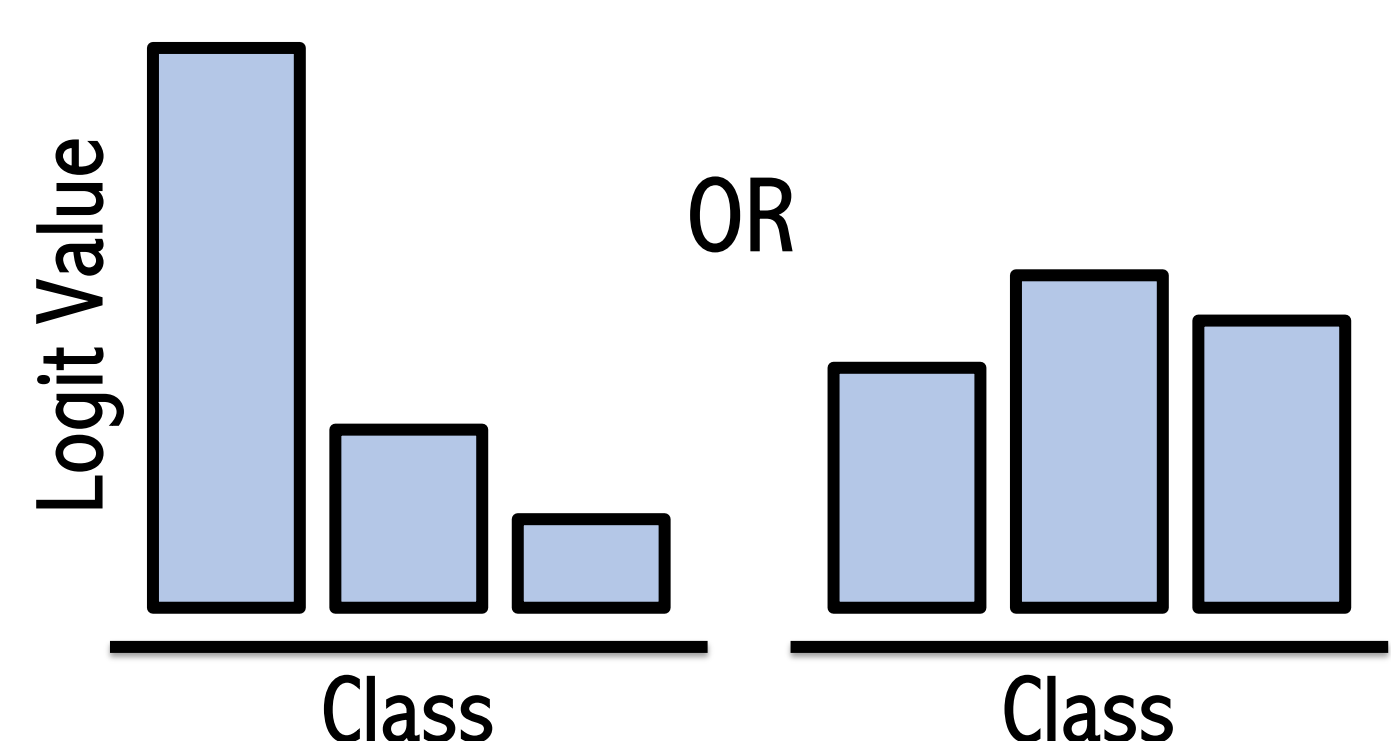
$$L_{\text{Mowst}} = \frac{1}{|V|} \sum_{v \in V} (C(\mathbf{p}_v) \cdot L(\mathbf{p}_v, \mathbf{y}_v) + (1 - C(\mathbf{p}_v)) \cdot L(\mathbf{p}'_v, \mathbf{y}_v))$$

Target node:  $v$

MLP's prediction:  $\mathbf{p}_v$  GNN's prediction:  $\mathbf{p}'_v$

How confident is MLP:  $C(\mathbf{p}_v)$

High Dispersion OR Low Dispersion



Confidence

Nodes are split based on the confidence of the weak expert

Weak Expert (MLP)

Strong Expert (GNN)

### Algorithm 1 Mowst inference

**Input:**  $\mathcal{G}(\mathcal{V}, \mathcal{E}, \mathbf{X})$ ; target node  $v$   
**Output:** prediction of  $v$   
 Run the trained MLP expert on  $v$   
 Get prediction  $\mathbf{p}_v$  & confidence  $C(\mathbf{p}_v) \in [0, 1]$   
**if** random number  $q \in [0, 1]$  has  $q < C(\mathbf{p}_v)$  **then**  
   Predict  $v$  by MLP's prediction  $\mathbf{p}_v$   
**else**  
   Run the trained GNN expert on  $v$   
   Predict  $v$  by GNN's prediction  $\mathbf{p}'_v$   
**end if**

### Algorithm 2 Mowst training

**Input:**  $\mathcal{G}(\mathcal{V}, \mathcal{E}, \mathbf{X})$ ; training labels  $\{\mathbf{y}_v\}$   
 Initialize MLP & GNN weights as  $\theta_0$  &  $\theta'_0$   
**for** round  $r = 1$  until convergence **do**  
   Fix GNN weights  $\theta'_{r-1}$   
   Update MLP weights to  $\theta_r$  by gradient descent on  $L_{\text{Mowst}}$   
 until convergence  
   Fix MLP weights  $\theta_r$   
   Update GNN weights to  $\theta'_r$  by gradient descent on  $L_{\text{Mowst}}$   
 until convergence  
**end for**

## Mowst\*

$$L_{\text{Mowst}^*} = \frac{1}{|V|} \sum_{v \in V} L(C(\mathbf{p}_v) \cdot \mathbf{p}_v + (1 - C(\mathbf{p}_v)) \cdot \mathbf{p}'_v, \mathbf{y}_v)$$

- Mowst may be easier to optimize, while Mowst\* has a theoretically lower loss.

	Flickr	pokec	twitch-gamer
Mowst-GCN (joint)	53.47±0.36	76.62±0.11	63.44±0.22
Mowst-GCN	<b>54.62±0.23</b>	77.12±0.09	<b>63.74±0.23</b>
Mowst*-GCN	53.94±0.37	<b>77.28±0.08</b>	63.59±0.11

## Expressive Power & Computation Complexity

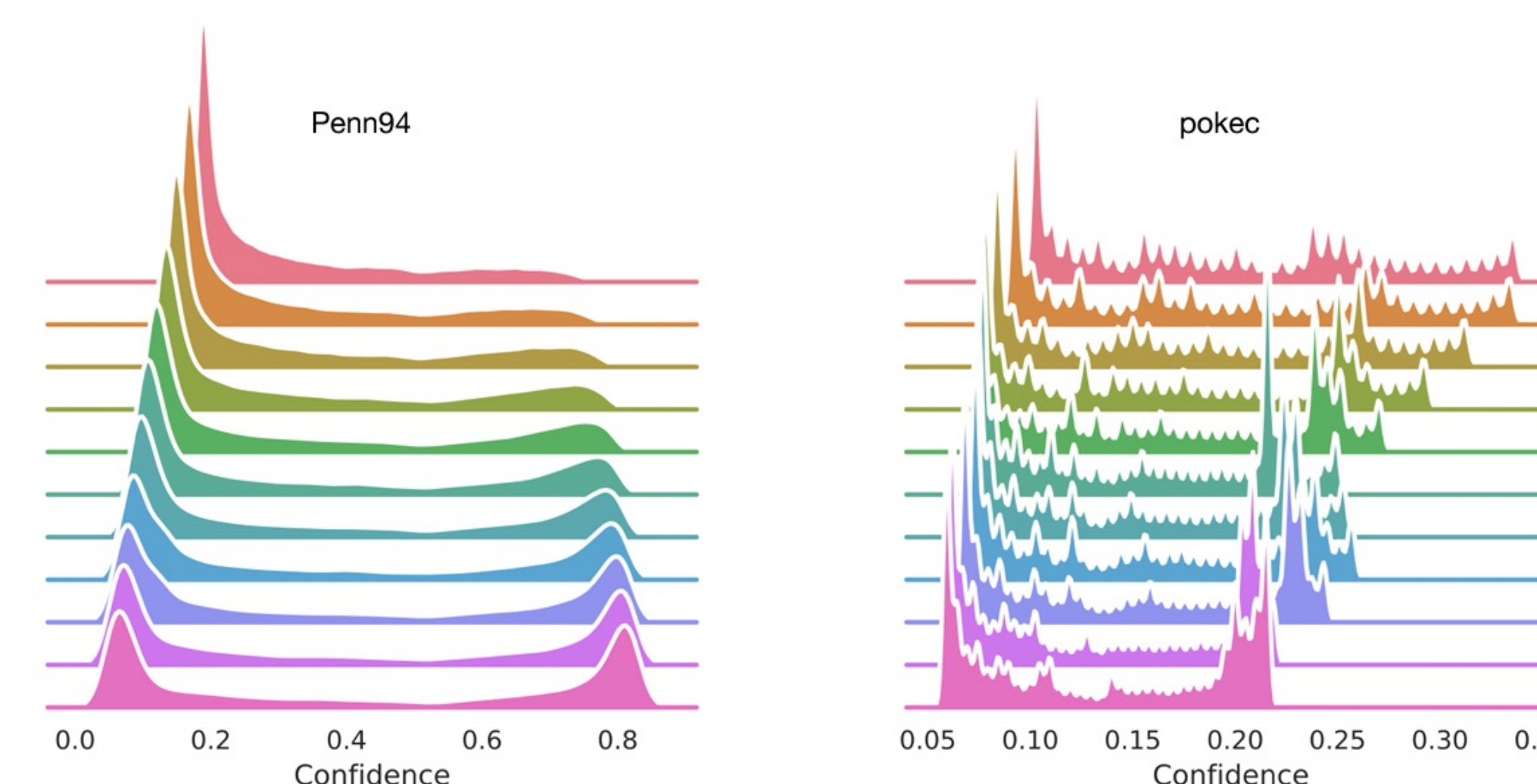
- Mowst and Mowst\* are at least as expressive as the MLP or GNN expert alone.
  - Mowst-GCN and Mowst\*-GCN are more expressive than the GCN expert alone.
- The worst-case cost of Mowst-GNN or Mowst\*-GNN is similar to that of a vanilla GNN.

## Main Results

- Mowst(\*) outperforms all other baselines under the same number of layers and hidden dimensions.
- The decoupling of the self-features and neighbor structures, along with the denoising effect of the weak expert are generally beneficial.

	Flickr	ogbn-products	ogbn-arxiv	Penn94	pokec	twitch-gamer
MLP	46.93 ±0.00	61.06 <sup>†</sup> ±0.08	55.50 <sup>†</sup> ±0.23	73.61 <sup>‡</sup> ±0.40	62.37 <sup>‡</sup> ±0.02	60.92 <sup>‡</sup> ±0.07
GAT	52.47 ±0.14	OOM	71.58 ±0.17	81.53 <sup>‡</sup> ±0.55	71.77 <sup>‡</sup> ±6.18	59.89 <sup>‡</sup> ±4.12
GPR-GNN	53.23 ±0.14	72.41 ±0.04	71.10 ±0.22	81.38 <sup>‡</sup> ±0.16	<b>78.83<sup>‡</sup></b> ±0.05	61.89 <sup>‡</sup> ±0.29
AdaGCN	48.96 ±0.06	69.06 ±0.04	58.45 ±0.50	74.42 ±0.58	55.92 ±0.35	61.02 ±0.14
GCN	53.86 ±0.37	75.64 <sup>†</sup> ±0.21	71.74 <sup>†</sup> ±0.29	82.17 ±0.04	76.01 ±0.49	62.42 ±0.53
GCN-skip	52.98 ±0.00	-	69.56 ±0.00	76.58 ±0.53	73.46 ±0.04	61.05 ±0.23
GraphMoE-GCN	53.03 ±0.14	73.90 ±0.00	71.88 <sup>††</sup> ±0.32	81.61 ±0.27	76.99 ±0.10	62.76 ±0.22
Mowst(*)-GCN	<u>54.62 ±0.23</u> (+0.76)	76.49 ±0.22 (+0.85)	<b>72.52 ±0.07</b> (+0.64)	<u>83.19 ±0.43</u> (+1.02)	77.28 ±0.08 (+0.29)	63.74 ±0.23 (+0.83)
GIN	53.71 ±0.35	-	69.39 ±0.56	82.68 ±0.32	53.37 ±2.15	61.76 ±0.60
Mowst(*)-GIN	<b>55.48 ±0.32</b> (+1.77)	-	71.43 ±0.26 (+2.04)	<b>84.56 ±0.31</b> (+1.88)	76.11 ±0.39 (+22.74)	64.32 ±0.34 (+2.56)
GIN-skip	52.70 ±0.00	-	71.28 ±0.00	80.32 ±0.43	76.29 ±0.51	64.27 ±0.25
Mowst(*)-GIN-skip	53.19 ±0.31 (+0.49)	-	71.79 ±0.23 (+0.51)	81.20 ±0.55 (+0.88)	<b>79.70 ±0.23</b> (+3.41)	<b>64.91 ±0.22</b> (+0.64)
GraphSAGE	53.51 ±0.05	78.50 <sup>†</sup> ±0.14	71.49 <sup>†</sup> ±0.27	76.75 ±0.52	75.76 ±0.04	61.99 ±0.30
GraphMoE-SAGE	52.16 ±0.13	77.79 ±0.00	71.19 ±0.15	77.04 ±0.55	76.67 ±0.08	63.42 ±0.23
Mowst(*)-SAGE	53.90 ±0.18 (+0.39)	<b>79.38 ±0.44</b> (+0.88)	<u>72.04 ±0.24</u> (+0.55)	79.07 ±0.43 (+2.03)	77.84 ±0.04 (+1.33)	<u>64.38 ±0.14</u> (+1.05)

- Mowst can substantially enhance the performance of state-of-the-art heterophilous GNNs like H2GCN, with the help of a relatively simple expert such as a standard MLP.



	Penn94	pokec	twitch-gamer
GCN	82.17 ±0.04	76.01 ±0.49	62.42 ±0.53
Mowst(*)-GCN	83.19 ±0.43 (+1.02)	77.28 ±0.08 (+0.29)	63.74 ±0.23 (+0.83)
H2GCN	82.71 ±0.67	80.89 ±0.16	65.70 ±0.20
Mowst(*)-H2GCN	<b>83.39 ±0.43</b> (+0.68)	<b>83.02 ±0.30</b> (+2.13)	<b>66.03 ±0.16</b> (+0.33)

- Specialization via Data Splitting.** Both Mowst and Mowst\* adapt their expert collaboration based on the confidence-weighted loss across various graphs.

- See Appendix for the empirical findings on denoised fine-tuning.

## Future Work

- Multi-expert (e.g., Mixture of progressively stronger experts, hierarchical mixture)
- Weak and strong experts in non-graph domains (e.g., NLP, computer vision)

## Acknowledgments

We are grateful to Jingyang Lin, Dr. Wei Zhu, and Dr. Wei Xiong for their constructive suggestions. Luo was supported in part by NSF Award #2238208.