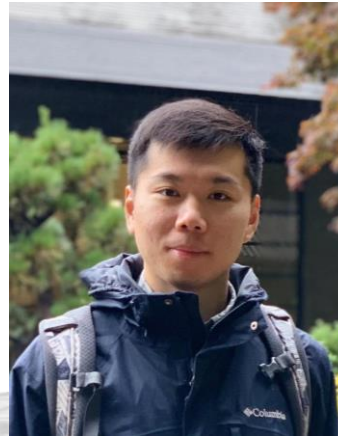




# JuryGCN: Quantifying Jackknife Uncertainty on Graph Convolutional Networks



Jian Kang\*

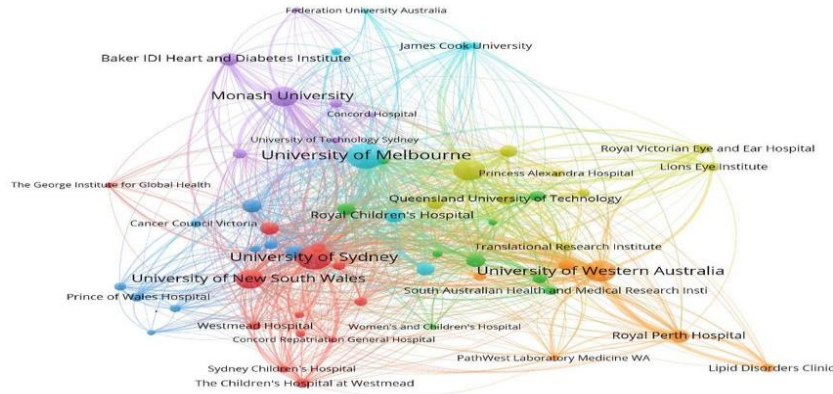


Qinghai Zhou\*

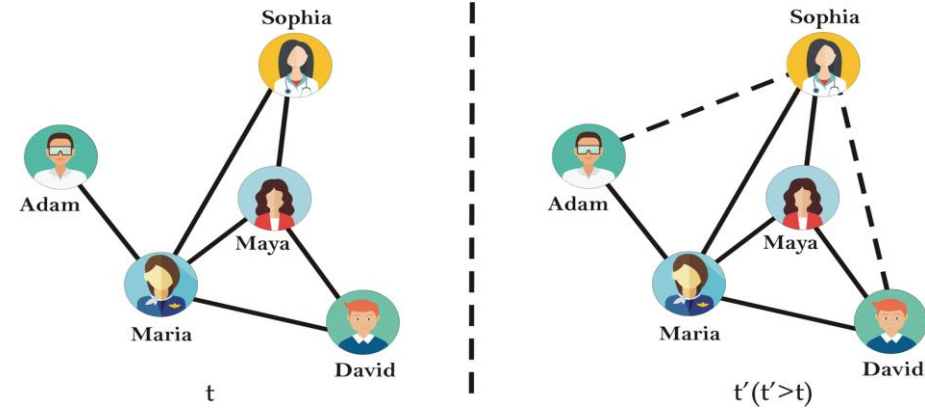


Hanghang Tong

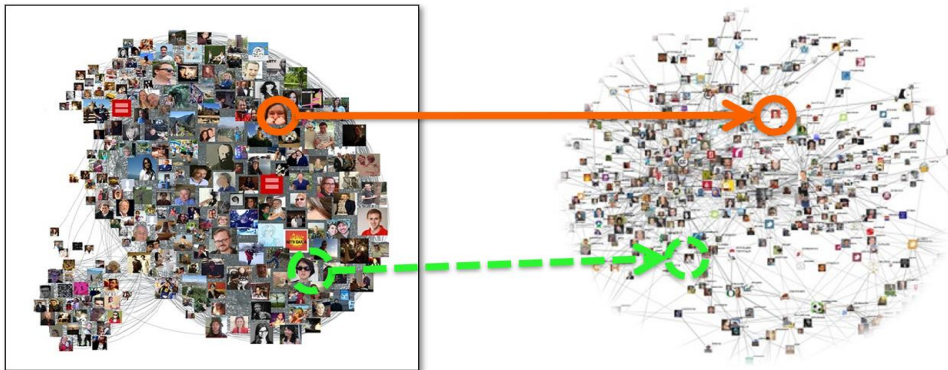
# Applications of Graph Neural Networks



Node classification [1]



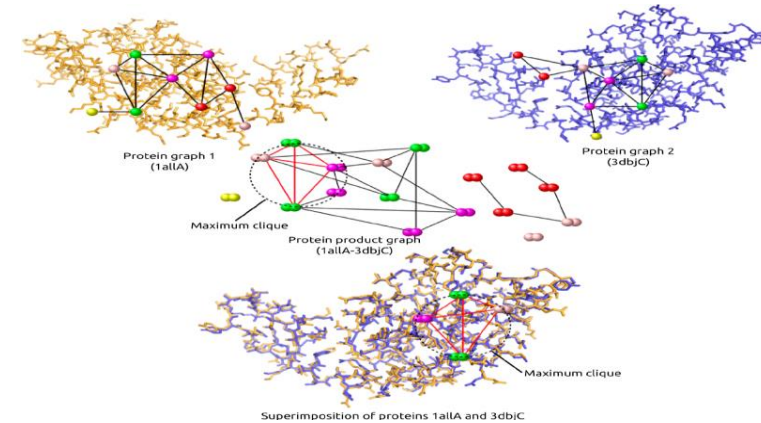
Link prediction [2]



Network A: Facebook

Network B: Twitter

Network alignment [3]



Graph classification [4]

[1] Kipf, T. N., & Welling, M. (2016). Semi-supervised classification with graph convolutional networks. In *arXiv 2016*.

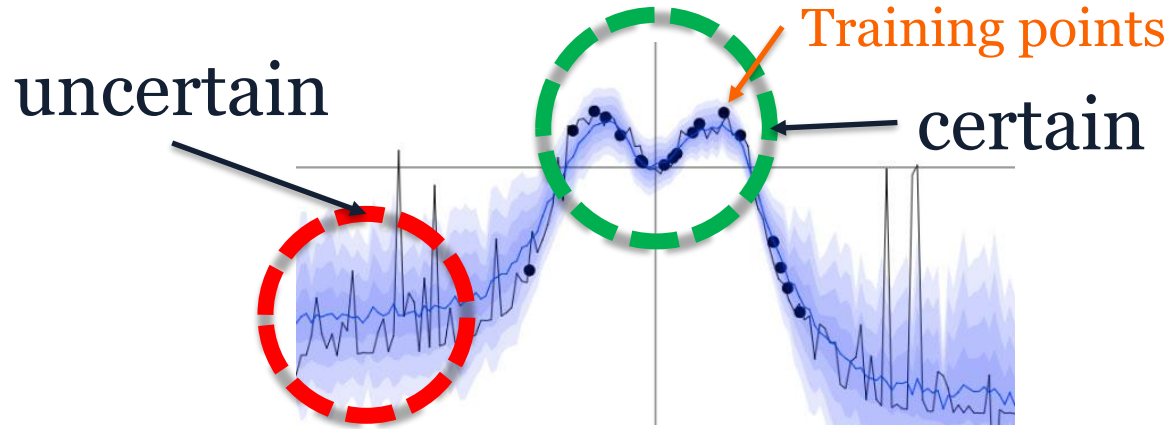
[2] Zhang, M., & Chen, Y. (2018). Link prediction based on graph neural networks. In *NeurIPS 2018*.

[3] Zhang, S., Tong, H., Xia, Y., Xiong, L., & Xu, J. (2020, August). Nettrans: Neural cross-network transformation. In *KDD 2020*.

[4] Errica, F., Podda, M., Bacciu, D., & Micheli, A. (2019). A fair comparison of graph neural networks for graph classification. In *arXiv 2019*.

# Uncertainty in Model Prediction

## Examples



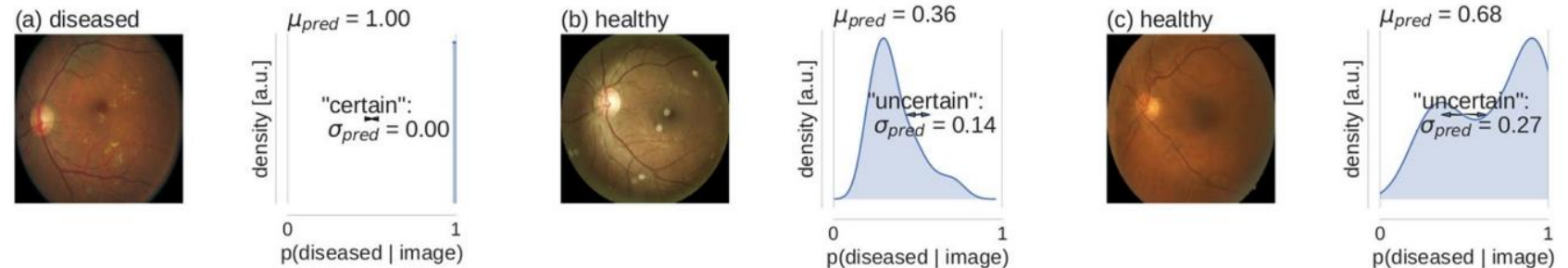
Regression



Classification

## Quantifying the uncertainty is important in high-risk applications

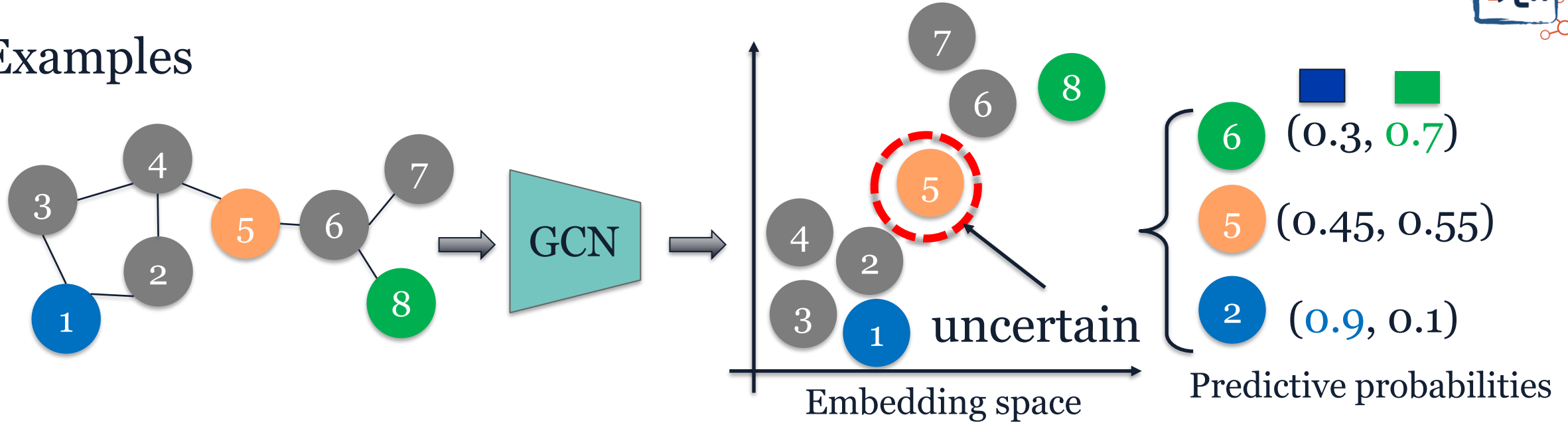
- E.g., medical



# Uncertainty in Graph Learning



## □ Examples



## □ Questions:

Q1: How uncertain is a GCN in its own predictions?

→ Uncertainty quantification (UQ)

Q2: How to improve GCN predictions by leveraging uncertainty?

→ Application of UQ



# Existing Solutions: Bayesian-based

- ❑ Motivation: address over-smoothing/fitting
- ❑ Key idea:
  - ❑ adaptively drop edges
  - ❑ Monte Carlo estimation for posterior uncertainty [1].
- ❑ Limitations: not explicitly quantify the uncertainty on model prediction (ad-hoc)



# Existing Solutions:

## Deterministic Quantification-based

- ❑ Motivation: estimate multi-source uncertainty for GNNs
- ❑ Key idea: a graph-based Dirichlet distribution to reduce errors in quantifying uncertainties [2].
- ❑ Limitations: changing the training procedure, e.g., additional parameters (e.g., Dirichlet distribution) or architectures (e.g., teacher network)



[1] Hasanzadeh, A. et al. Bayesian graph neural networks with adaptive connection sampling. In ICML 2020.

[2] Zhao, X., Chen, F., Hu, S., & Cho, J. H. (2020). Uncertainty aware semi-supervised learning on graph data. In NeurIPS 2020.

# Roadmap

- Background & Motivation ✓
- **JuryGCN Formulation** ←
- JuryGCN Algorithms
- JuryGCN Applications
- Experimental Results
- Conclusion

# Problem Definition

□ Given:

- (1) an undirected graph  $G = \{V, A, X\}$ ;
- (2) an  $L$ -layer GCN with parameter  $\Theta$ ;
- (3) a task-specific objective  $R(G, Y, \Theta)$  ( $Y$ : ground-truth)

□ Find:

An uncertainty score  $U_{\Theta}(u)$  for any node  $u$  in graph  $G$  w.r.t. parameters  $\Theta$  and objective  $R(G, Y, \Theta)$ .







# Preliminaries: Jackknife+ Resampling

□ Key idea: leaving out an observation → evaluating prediction error (LOO)

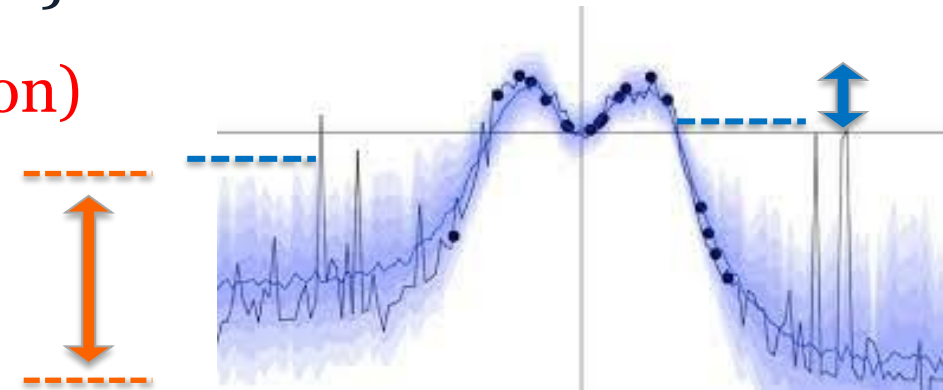
□ Given: training data:  $D = \{(x_i, y_i) | i = 1, \dots, n\}$ ; a test point  $(x^*, y^*)$ ; a trained model  $f_\theta()$ ; target coverage  $\alpha$ ;

□ Confidence interval:  $[C^-(x^*), C^+(x^*)]$

- $C^+(x^*) = Q_{1-\alpha}(P^+)$ ,  $C^-(x^*) = Q_\alpha(P^-)$
- $P^+ = \{ \underbrace{f_{\theta_{-i}}(x^*)}_{\text{LOO prediction}} + \underbrace{|y_i - f_{\theta_{-i}}(x_i)|}_{\text{Error residual (generalization)}} | i = 1, \dots, n \}$
- $P^- = \{ \underbrace{f_{\theta_{-i}}(x^*)}_{\text{LOO prediction}} - \underbrace{|y_i - f_{\theta_{-i}}(x_i)|}_{\text{Error residual (generalization)}} | i = 1, \dots, n \}$

LOO prediction    Error residual (generalization)

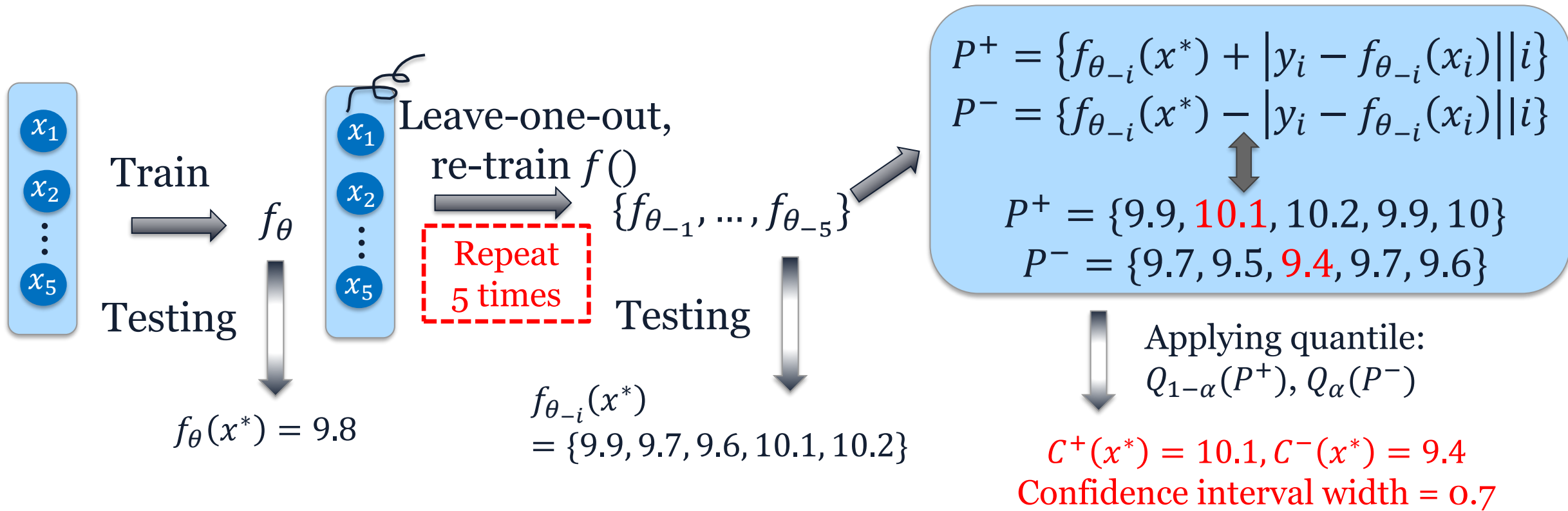
Larger interval → less confident





# Jackknife+ Resampling: A Numerical Example

□ Regression task: training set,  $\{(x_1, y_1), \dots, (x_5, y_5)\}$ , a test point,  $(x^*, y^*)$  where  $y^* = 10$ , coverage,  $\alpha = 0.2$



[1] Barber, Rina Foygel, et al. "Predictive inference with the jackknife+." The Annals of Statistics 49.1 (2021): 486-507.



# Challenges

□ C1: How to formally define the Jackknife uncertainty for GNNs?

- Non-IID graph data

□ C2: How to efficiently compute the node uncertainty?

- Avoid re-training

w.l.o.g, considering a node-level tasks (e.g., node classification)

$$\Theta^* = \operatorname{argmin}_{\Theta} R(G, Y_{\text{train}}, \Theta) = \operatorname{argmin}_{\Theta} \frac{1}{|V_{\text{train}}|} \sum_v r(v, y_v, \Theta)$$

Training labels ←  $Y_{\text{train}}$ 
Training set ←  $|V_{\text{train}}|$ 
Node-specific loss (cross-entropy) ←  $r(v, y_v, \Theta)$

$$r(v, y_v, \Theta) = - \sum_{i=1}^c y_v[i] \log(\text{GCN}(v, \Theta)[i])$$



# Jackknife Uncertainty: Definition

Confidence interval:  $U_{\Theta}(u) = C_{\Theta}^{+}(u) - C_{\Theta}^{-}(u)$

Error residual:  
 $err_i = \left\| y_i - GCN(i, \Theta_{\epsilon,i}^{*}) \right\|_2$

Compute  $C^{+}, C^{-}$   $C_{\Theta}^{-}(u) = Q_{\alpha}(\left\{ \left\| GCN(u, \Theta_{\epsilon,i}^{*}) \right\|_2 - err_i \mid \forall i \in V_{train} \right\})$

$C_{\Theta}^{+}(u) = Q_{1-\alpha}(\left\{ \left\| GCN(u, \Theta_{\epsilon,i}^{*}) \right\|_2 + err_i \mid \forall i \in V_{train} \right\})$

Why Jackknife+: stable coverage

Upweighting the loss of node  $i$ :  
 $\Theta_{\epsilon,i}^{*} = \operatorname{argmin}_{\Theta} \epsilon r(i, y_i, \Theta) \frac{1}{|V_{train}|} \sum_v r(v, y_v, \Theta)$

Question: How to obtain  $\Theta_{\epsilon,i}^{*}$  without re-training?



# Roadmap

- Background & Motivation ✓
- JuryGCN Formulation ✓
- **JuryGCN Algorithms** ←
- JuryGCN Applications
- Experimental Results
- Conclusion



# Jackknife Uncertainty: Efficient Computation

□ Key idea: efficiently estimate  $\Theta_{\epsilon,i}^*$  with influence function [1]

□ Taylor expansion over parameters

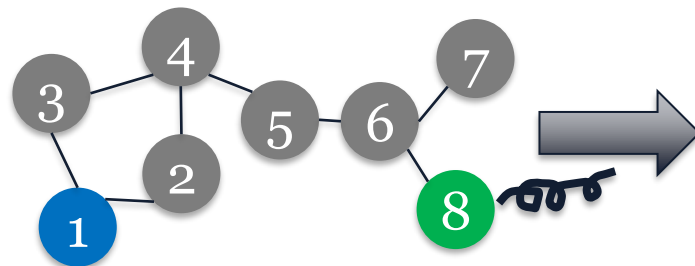
$$\Theta_{\epsilon,i}^* \approx \Theta^* + \epsilon \mathbf{I}_{\Theta^*}(i) \quad (1) \text{ where } \mathbf{I}_{\Theta^*}(i) = \left. \frac{d\Theta_{\epsilon,i}^*}{d\epsilon} \right|_{\epsilon \rightarrow 0}$$

□ The influence function can be further computed as [2],

$$\mathbf{I}_{\Theta^*}(i) = \mathbf{H}_{\Theta^*}^{-1} \nabla_{\Theta} r(i, y_i, \Theta^*) \quad (2) \text{ Hessian matrix w.r.t. model parameters}$$

$$\mathbf{H}_{\Theta^*} = \frac{1}{|V_{\text{train}}|} \nabla_{\Theta}^2 R(G, Y_{\text{train}}, \Theta^*)$$

By setting  $\epsilon = -\frac{1}{|V_{\text{train}}|}$ , the leave-one-out parameters,  $\Theta_{\epsilon,i}^*$  (Eq. (1)) can be computed efficiently.



Adjusting the weights by  $\epsilon$



# Jackknife Uncertainty: Efficient Computation (Cont.)

□ Proposition: First-order derivative of GCN [1] w.r.t. the parameters in the  $l$ -th layer, i.e.,  $\mathbf{W}^{(l)} \leftarrow \nabla_{\mathbf{W}^{(l)}} r(i, y_i, \Theta)$

$$\mathbf{I}_{\Theta^*}(i) = \mathbf{H}_{\Theta^*}^{-1} \nabla_{\Theta} r(i, y_i, \Theta^*)$$

- Key idea: apply chain rule on layer parameters.

$$\nabla_{\mathbf{W}^{(l)}} r(i, y_i, \Theta) = (\widehat{\mathbf{A}} \mathbf{E}^{l-1})^T \left( \frac{\partial r(i, y_i, \Theta)}{\partial \mathbf{E}^{(l)}} \odot \sigma'(\widehat{\mathbf{A}} \mathbf{E}^{(l-1)} \mathbf{W}^{(l)}) \right)$$

Hidden representations

$$\mathbf{E}^{(l)} = \sigma(\widehat{\mathbf{A}} \mathbf{E}^{(l-1)} \mathbf{W}^{(l)})$$

Normalized graph Laplacian



# Jackknife Uncertainty: Efficient Computation (Cont.)

$$\mathbf{I}_{\Theta^*}(i) = \mathbf{H}_{\Theta^*}^{-1} \nabla_{\Theta} r(i, y_i, \Theta^*)$$

□ Theorem: Computing the Hessian tensor of GCN

(the  $i$ -th and  $l$ -th layer)  $\rightarrow S_{l,i} = \frac{\partial^2 R}{\partial \mathbf{W}^{(l)} \partial \mathbf{W}^{(i)}}$



- vectorize the first-order
- compute the element-wise second-order

4-D tensor

- Flattened Hessian matrix
- Applying Hessian-vector product [2] using power iteration

THEOREM 1. (The Hessian tensor of GCN) Following the settings of Proposition 2, denoting the overall loss  $R(\mathcal{G}, \mathcal{Y}_{train}, \Theta)$  as  $R$  and  $\sigma'_l$  as  $\sigma'(\hat{\mathbf{A}}\mathbf{E}^{(l-1)}\mathbf{W}^{(l)})$ , the Hessian tensor  $\mathfrak{H}_{l,i} = \frac{\partial^2 R}{\partial \mathbf{W}^{(l)} \partial \mathbf{W}^{(i)}}$  of  $R$  with respect to  $\mathbf{W}^{(l)}$  and  $\mathbf{W}^{(i)}$  has the following forms.

- Case 1.  $i = l, \mathfrak{H}_{l,i} = 0$
- Case 2.  $i = l - 1$

$$\mathfrak{H}_{l,i}[:, :, c, d] = \left( \hat{\mathbf{A}} \frac{\partial \mathbf{E}^{(l-1)}}{\partial \mathbf{W}^{(i)}[c, d]} \right)^T \left( \frac{\partial R}{\partial \mathbf{E}^{(l)}} \circ \sigma'_l \right) \quad (11)$$

where  $\frac{\partial \mathbf{E}^{(l-1)}}{\partial \mathbf{W}^{(i)}[c, d]}$  is the matrix whose entry at the  $a$ -th row and the  $b$ -th column is

$$\frac{\partial \mathbf{E}^{(l-1)}[a, b]}{\partial \mathbf{W}^{(l-1)}[c, d]} = \sigma'_{l-1}[a, b] (\hat{\mathbf{A}} \mathbf{E}^{(l-2)})[a, c] \mathbf{I}[b, d] \quad (12)$$

- Case 3.  $i < l - 1$
- Apply Eq. (12) for the  $i$ -th hidden layer.
- Forward to the  $(l - 1)$ -th layer iteratively with

$$\frac{\partial \mathbf{E}^{(l-1)}}{\partial \mathbf{W}^{(i)}[c, d]} = \sigma'_{l-1} \circ \left( \hat{\mathbf{A}} \frac{\partial \mathbf{E}^{(l-2)}}{\partial \mathbf{W}^{(i)}[c, d]} \mathbf{W}^{(l-1)} \right) \quad (13)$$

- Apply Eq. (11).
- Case 4.  $i = l + 1$

$$\mathfrak{H}_{l,i}[:, :, c, d] = (\hat{\mathbf{A}} \mathbf{E}^{(l-1)})^T \left( \frac{\partial^2 R}{\partial \mathbf{E}^{(l)} \partial \mathbf{W}^{(i)}[c, d]} \circ \sigma'_l \right) \quad (14)$$

where  $\frac{\partial^2 R}{\partial \mathbf{E}^{(l)}[a, b] \partial \mathbf{W}^{(l+1)}[c, d]} = \mathbf{I}[b, c] \left[ \hat{\mathbf{A}}^T \left( -\frac{\partial R}{\partial \mathbf{E}^{(l+1)}} \circ \sigma'_{l+1} \right) \right][a, d]$ .

- Case 5.  $i > l + 1$
- Compute  $\frac{\partial^2 R}{\partial \mathbf{E}^{(l-1)} \partial \mathbf{W}^{(i)}[c, d]}$  whose  $(a, b)$ -th entry has the form  $\frac{\partial^2 R}{\partial \mathbf{E}^{(l-1)}[a, b] \partial \mathbf{W}^{(i)}[c, d]} = \mathbf{I}[b, c] \left( \hat{\mathbf{A}}^T \left( \frac{\partial R}{\partial \mathbf{E}^{(i)}} \circ \sigma'_i \right) \right)[a, d]$
- Backward to  $(l + 1)$ -th layer iteratively with

$$\frac{\partial^2 R}{\partial \mathbf{E}^{(l)} \partial \mathbf{W}^{(i)}[c, d]} = \hat{\mathbf{A}}^T \left( \frac{\partial^2 R}{\partial \mathbf{E}^{(l+1)} \partial \mathbf{W}^{(i)}[c, d]} \circ \sigma'_{l+1} \right) (\mathbf{W}^{(l+1)})^T \quad (15)$$

- Apply Eq. (14).

[1] Kang, J. et al. RawlsGCN: Towards Rawlsian Difference Principle on Graph Convolutional Network. In WWW 2022.  
 [2] Alaa, A. et al. Discriminative Jackknife: Quantifying Uncertainty in Deep Learning via Higher-Order Influence Functions. In ICML 2020.







# Algorithm: JuryGCN

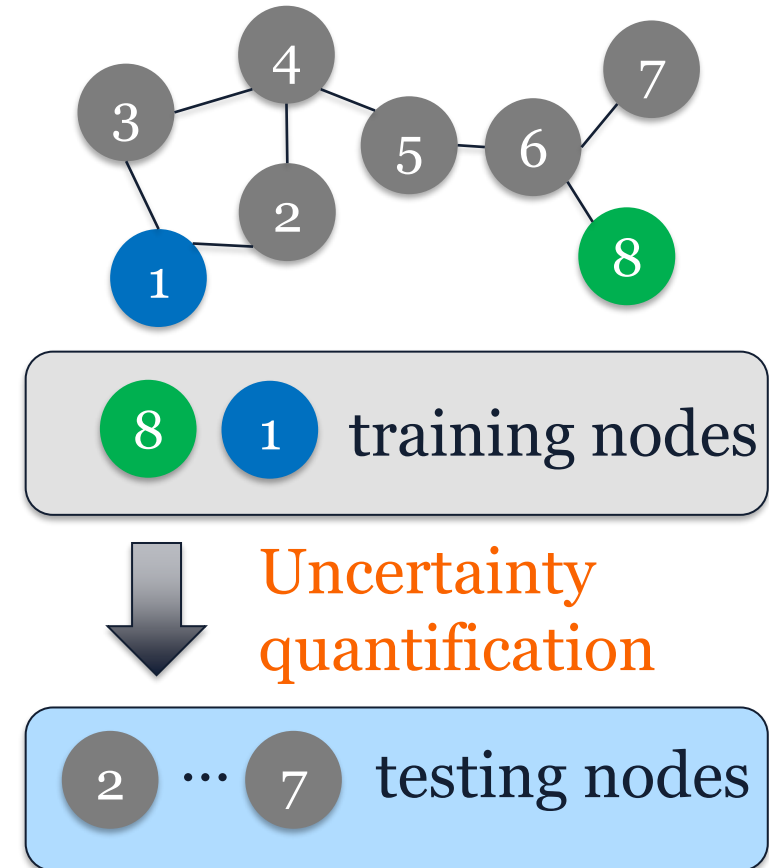
□ Goal: to estimate uncertainty  $U_{\Theta}(u)$  of node  $u$ .

□ Initialize:  $\epsilon = -\frac{1}{|V_{\text{train}}|}$ , a GCN with parameter  $\Theta$

□ Key steps (for each training node):

- Compute node-wise loss  $r_{i,\Theta}$  and derivative  $\nabla_{\Theta} r_{i,\Theta}$
- Evaluate the influence w.r.t. training node
- Compute LOO parameters/predictions/errors
- Compute lower and upper bound

□ Return: confidence interval of node  $u$ .



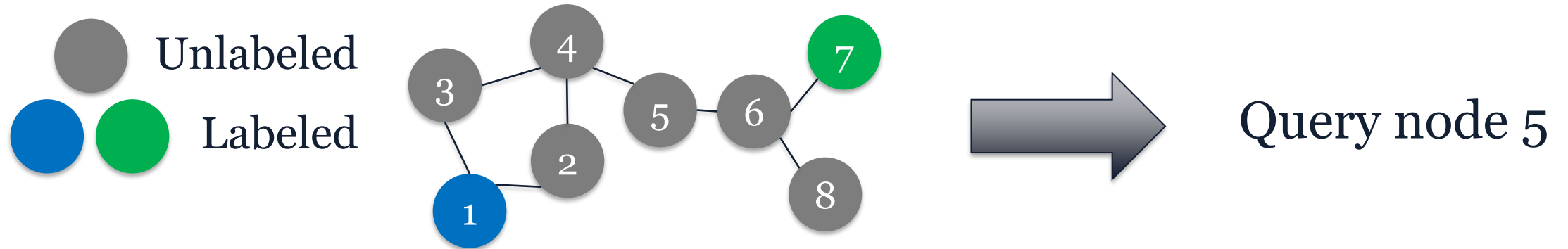
# Roadmap

- Background & Motivation ✓
- JuryGCN Formulation ✓
- JuryGCN Algorithms ✓
- **JuryGCN Applications** ←
- Experimental Results
- Conclusion

# Applications: Active Learning on Node Classification



- ❑ Task: query the nodes for true labels → node classifier
- ❑ General idea: select the most informative nodes



- ❑ Our idea: iteratively query the nodes with the largest uncertainty

$$Acq(V_{\text{train}}) = \operatorname{argmax}_{u \in V_{\text{train}}} U_{\Theta}(u)$$

# Applications: Semi-supervised Node Classification



- Existing objective: mean of loss from all training nodes

$$R = \frac{1}{|V_{\text{train}}|} \sum_{i \in V_{\text{train}}} r(i, y_i, \Theta)$$

- Uncertainty-aware node-specific objective

$$r_u = -\beta_u^\tau \log(p_u^{(i)})$$

*i*-th class predictive probability

$$\beta_u = \frac{|U_\Theta(u)|}{\sqrt{\sum_{i \in V_{\text{train}}} |U_\Theta(u)|^2}}$$

normalizing over all training nodes

- 
- (1) *u* is misclassified &  $p_u^{(i)}$  is small.
  - (2) *u* is well classified &  $p_u^{(i)}$  is large.

# Roadmap

- Background & Motivation ✓
- JuryGCN Formulation ✓
- JuryGCN Algorithms ✓
- JuryGCN Applications ✓
- **Experimental Results** ←
- Conclusion

# Experiment Settings



❑ Datasets: 4 widely-adopted datasets

❑ Evaluation metric: micro-F1

❑ Comparison methods

- Active learning-based: AGE[1], ANRMAB[2], Coreset[3], SOPT-GCN[4], Centrality, Random
- Semi-supervised: S-GNN[5], GPN[6], GCN[7], GAT[8]

❑ Parameters

- Active node classification (Cora, Citeseer, Pubmed and Reddit)
  - Query budget: 100, 100, 50, 250, step size: 20, 20, 10, 50
- Semi-supervised node classification
  - hyperparameter:  $\tau = 2$ , coverage:  $\alpha = 0.025$

Datasets	Cora	Citeseer	PubMed	Reddit
# nodes	2,708	3,327	19,717	232,965
# edges	5,429	4,732	44,338	114,615,892
# features	1,433	3,703	500	602
# classes	7	6	3	41

# Experimental Results:

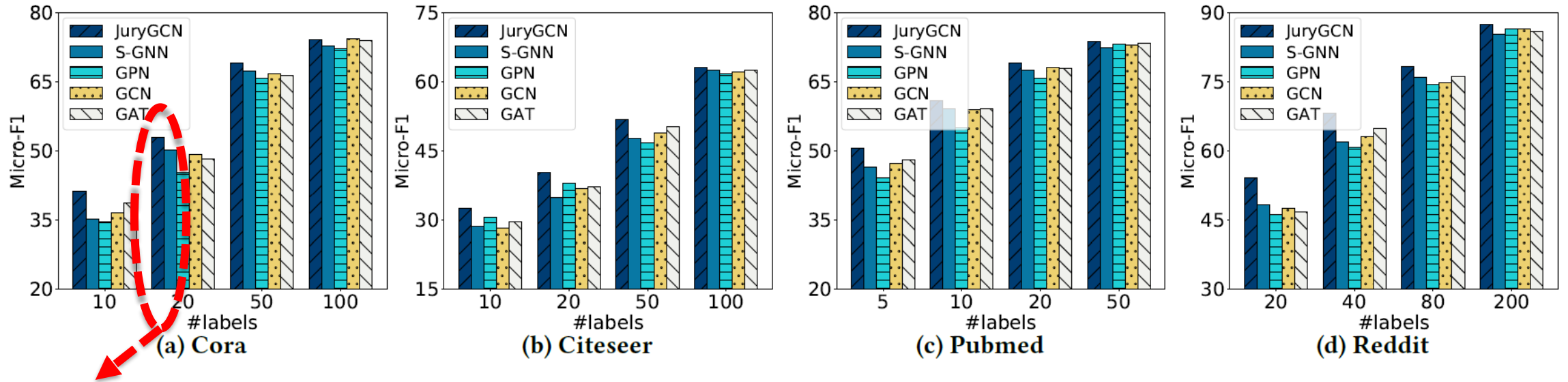
## Active Learning on Node Classification

proposed method

Data	Query size	JURYGCN (Ours)	ANRMAB	AGE	Coreset	Centrality	Degree	Random	SOPT-GCN
Cora	20	51.1 ± 1.2	46.8 ± 0.5	<u>49.4 ± 1.0</u>	43.8 ± 0.8	41.9 ± 0.6	38.5 ± 0.7	40.5 ± 1.6	48.8 ± 0.7
	40	64.7 ± 0.8	61.2 ± 0.8	58.2 ± 0.7	55.4 ± 0.5	57.3 ± 0.7	48.4 ± 0.3	56.8 ± 1.3	<u>62.6 ± 0.8</u>
	60	69.9 ± 0.9	67.8 ± 0.7	65.7 ± 0.8	62.2 ± 0.6	63.1 ± 0.5	58.8 ± 0.6	64.5 ± 1.5	<u>67.9 ± 0.6</u>
	80	74.2 ± 0.7	73.3 ± 0.6	72.5 ± 0.4	70.2 ± 0.5	69.1 ± 0.4	67.6 ± 0.4	69.7 ± 1.6	<u>73.6 ± 0.5</u>
	100	75.5 ± 0.6	74.9 ± 0.4	74.2 ± 0.3	73.8 ± 0.4	74.1 ± 0.3	73.0 ± 0.2	74.2 ± 1.2	<u>75.5 ± 0.7</u>
Citeseer	20	38.4 ± 1.5	35.9 ± 1.0	33.1 ± 0.9	30.2 ± 1.2	35.6 ± 1.1	31.5 ± 0.9	30.3 ± 2.3	<u>36.1 ± 0.7</u>
	40	51.1 ± 0.9	46.7 ± 1.3	49.5 ± 0.6	42.1 ± 0.8	<u>49.8 ± 1.3</u>	39.8 ± 0.7	41.1 ± 1.8	49.2 ± 0.5
	60	58.2 ± 0.8	55.2 ± 0.9	56.1 ± 0.5	52.1 ± 0.9	<u>57.1 ± 0.7</u>	50.1 ± 1.1	49.8 ± 1.3	56.4 ± 0.5
	80	63.8 ± 1.1	<u>63.2 ± 0.7</u>	61.5 ± 0.8	59.9 ± 0.6	63.3 ± 1.0	58.8 ± 0.6	58.1 ± 1.1	63.2 ± 0.8
	100	64.3 ± 1.2	<u>64.1 ± 0.5</u>	63.2 ± 0.7	62.8 ± 0.4	63.9 ± 0.6	61.8 ± 0.5	62.9 ± 0.8	63.8 ± 0.6
Pubmed	10	61.8 ± 0.9	60.5 ± 1.3	58.9 ± 1.1	53.1 ± 0.7	55.8 ± 1.2	56.4 ± 1.5	52.4 ± 1.7	59.5 ± 0.6
	20	70.2 ± 0.6	66.8 ± 1.1	<u>68.7 ± 0.7</u>	62.8 ± 0.5	67.2 ± 1.4	64.3 ± 1.0	60.5 ± 1.4	67.9 ± 0.9
	30	73.9 ± 0.3	71.6 ± 0.8	<u>72.8 ± 1.0</u>	68.9 ± 0.3	73.5 ± 0.9	70.1 ± 0.7	68.9 ± 1.1	72.3 ± 0.8
	40	<u>74.6 ± 0.4</u>	73.2 ± 0.6	<u>74.7 ± 0.8</u>	72.8 ± 0.8	74.1 ± 0.7	72.0 ± 0.8	71.8 ± 1.2	73.8 ± 0.7
	50	75.4 ± 0.5	74.7 ± 0.4	75.1 ± 0.5	73.5 ± 0.6	74.2 ± 0.6	72.9 ± 0.5	73.1 ± 1.0	<u>75.2 ± 0.5</u>
Reddit	50	69.7 ± 1.7	67.8 ± 0.9	64.2 ± 1.1	62.1 ± 0.6	65.5 ± 1.2	62.5 ± 1.4	63.7 ± 2.4	<u>68.1 ± 1.2</u>
	100	82.9 ± 1.5	<u>81.3 ± 1.0</u>	79.5 ± 0.8	81.2 ± 1.0	78.2 ± 0.9	81.1 ± 1.2	80.5 ± 1.6	80.4 ± 1.3
	150	86.0 ± 1.4	84.3 ± 0.7	83.2 ± 0.4	84.8 ± 0.9	84.1 ± 1.1	82.5 ± 1.2	81.5 ± 1.4	<u>85.0 ± 1.5</u>
	200	88.1 ± 0.9	86.1 ± 0.8	85.8 ± 0.5	85.5 ± 0.8	87.5 ± 0.8	85.4 ± 0.7	83.1 ± 1.8	<u>87.2 ± 0.9</u>
	250	89.2 ± 0.8	87.6 ± 0.7	87.1 ± 0.4	86.6 ± 1.1	<u>88.7 ± 0.6</u>	86.1 ± 1.0	87.3 ± 1.5	87.8 ± 1.1

Observation: JuryGCN achieves the best query performance

# Experimental Results: Semi-supervised Node Classification



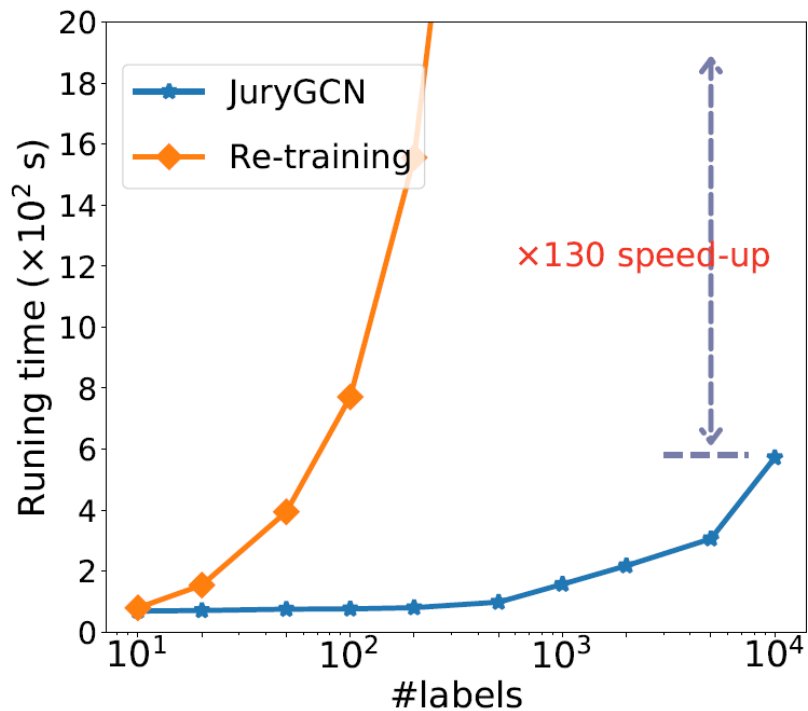
proposed method

Observation: achieving better performance when #labels is smaller

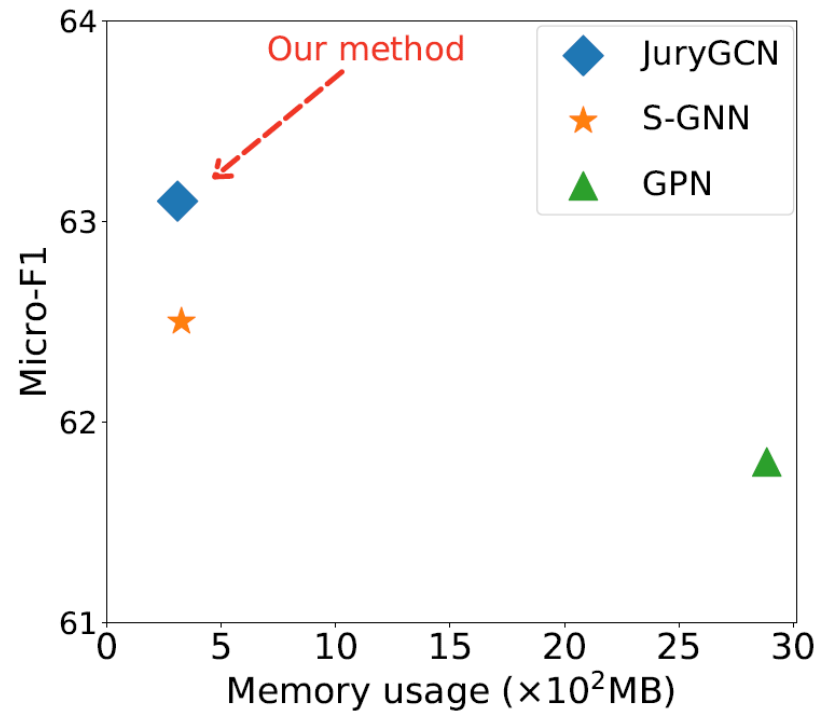


# Experimental Results: Efficiency

□ Metrics: running time, memory usage



Running time vs. #training labels

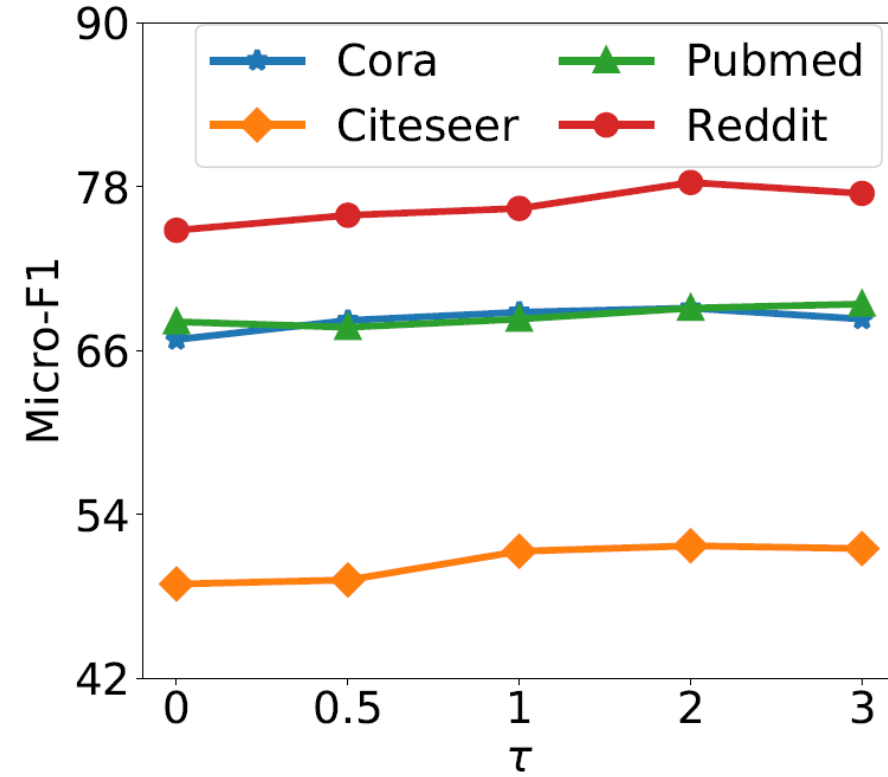
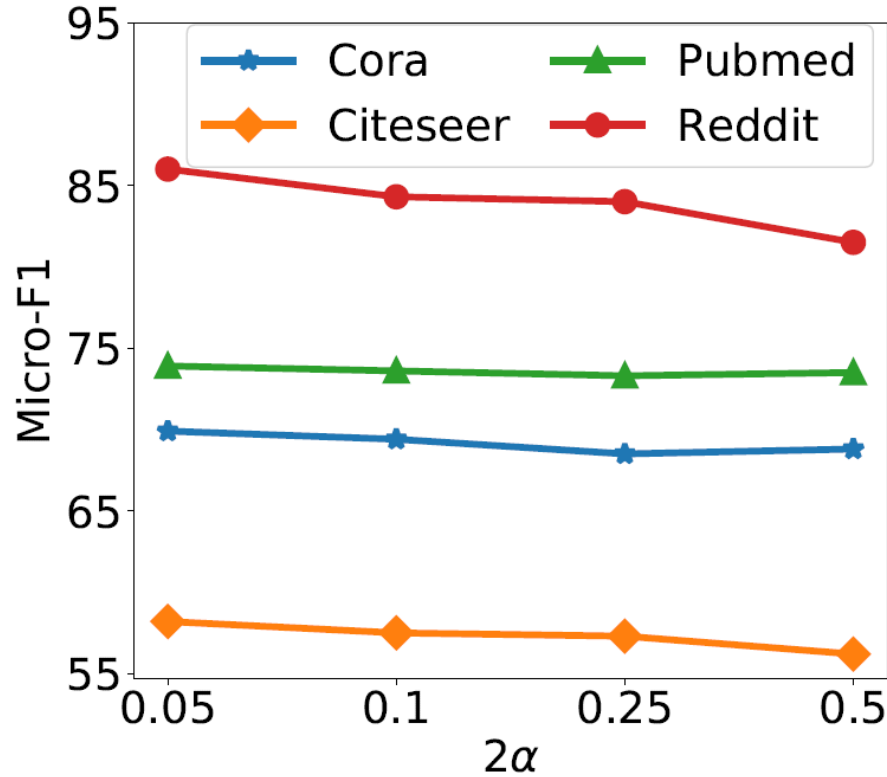


Micro-F1 vs. memory usage

Observation: JuryGCN can achieve the best efficiency performance.

# Experimental Results: Parameter Study

□ coverage,  $\alpha$ ; hyperparameter,  $\tau$



Observation: constantly achieving good performance.

# Roadmap

- Background & Motivation ✓
- JuryGCN Formulation ✓
- JuryGCN Algorithms ✓
- JuryGCN Applications ✓
- Experimental Results ✓
- **Conclusion** ←

# Conclusion



❑ Problem: Jackknife Uncertainty Quantification on GCN

❑ Solution:

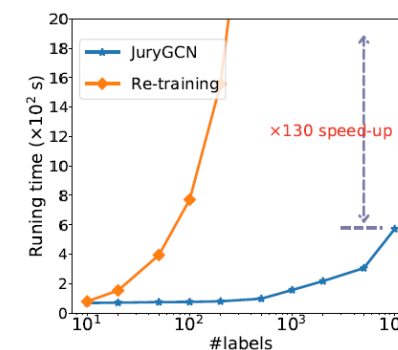
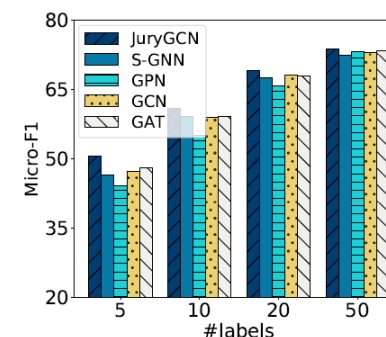
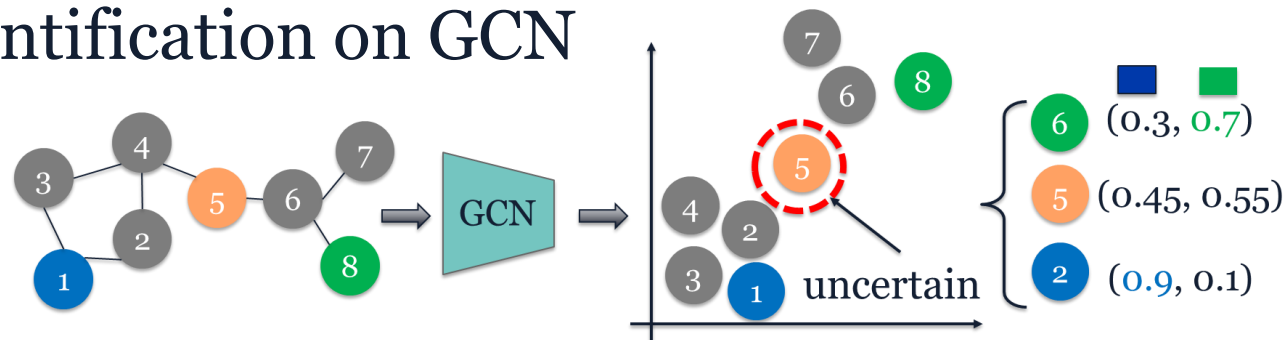
- Jackknife+ estimation
- Influence-based approach

❑ Applications:

- Active learning on node classification
- Semi-supervised node classification

❑ Results: outperforming other comparison method

- Improve node classification accuracy
- Select the most informative nodes
- Efficient computation compared to re-training



Title: JuryGCN: Quantifying Jackknife Uncertainty on Graph Convolutional Networks

Authors: Jian Kang, Qinghai Zhou, Hanghang Tong

Email: jiank2@illinois.edu, qinghai2@illinois.edu