

# Inductive Attributed Community Search: to Learn Communities across Graphs

Shuheng Fang<sup>†</sup>, Kangfei Zhao<sup>♠</sup>, Yu Rong<sup>\*</sup>, Zhixun Li<sup>†</sup>, Jeffrey Xu Yu<sup>†</sup>

The Chinese University of Hong Kong<sup>†</sup>

Beijing Institute of Technology<sup>♠</sup>

Alibaba DAMO Academy<sup>\*</sup>

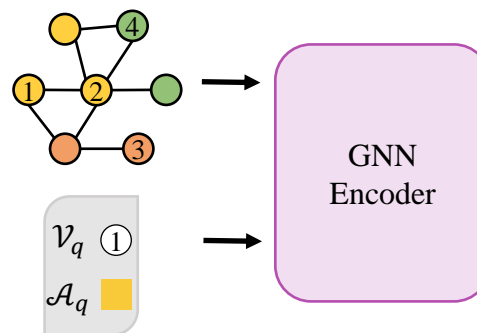
VLDB 2024

# Background: Attributed Community Search

- Community Search (CS)
  - For a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , given a node set  $\mathcal{V}_q \subseteq \mathcal{V}$  as a query  $q$
  - Find the query-dependent community  $C_q \subseteq \mathcal{V}$ , where the nodes in  $C_q$  are intensively intra-connected
- Attributed Community Search (ACS)
  - Satisfy both **structure cohesiveness** and **attributes homogeneity** for a given query that consists of query nodes and query attributes.
- Real applications
  - Social network analysis
  - Recommendation systems
  - Bioinformatics and fraud detection

# Background: Learning-based Community Search Approaches

- GNN-based method: recasting the community membership determination to a classification task
- AQD-GNN/ICS-GNN
  - Their trained models are tailored for specific graph/community
- ICS-GNN/COCLEP/CommunityAF/CGNP
  - Only support single-node query
  - COCLEP & CommunityAF have a limited inductive ability as they rely on the natural generalization of GNN
  - CGNP utilizes meta-learning approach and has inductive ability



# Background: Learning-based Community Search Approaches

**Table 1: Learning-based Community Search Approaches**

Approaches	Single-node Query	Multi-node Query	Attributed Query	Induction
<i>AQD-GNN</i> [29]	✓	✓	✓	✗
<i>ICS-GNN</i> [21]	✓	✗	✗	✗
<i>CommunityAF</i> [11]	✓	✗	✗	✗
<i>COCLEP</i> [35]	✓	✗	✗	✗
<i>CGNP</i> [17]	✓	✗	✗	✓
<i>IACS</i> (Ours)	✓	✓	✓	✓

# Problem Statement

- Attributed Community Search (ACS)
  - For an attributed graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A})$ , given a query  $q = (\mathcal{V}_q, \mathcal{A}_q)$ , where  $\mathcal{V}_q \subseteq \mathcal{V}$  and  $\mathcal{A}_q \subseteq \mathcal{A}$
  - Find the query-dependent community  $\mathcal{C}_q \subseteq \mathcal{V}$ , where the nodes in  $\mathcal{C}_q$  are intensively intra-connected and the attributes are similar

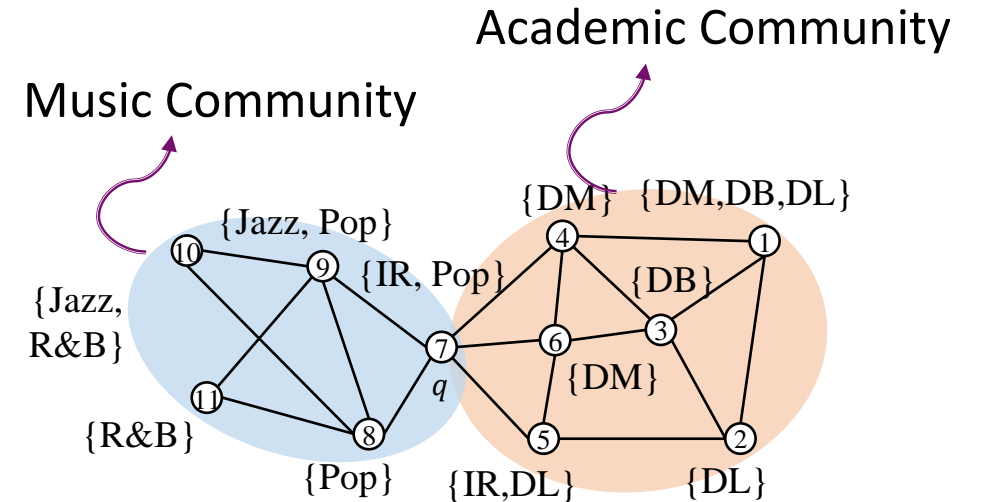
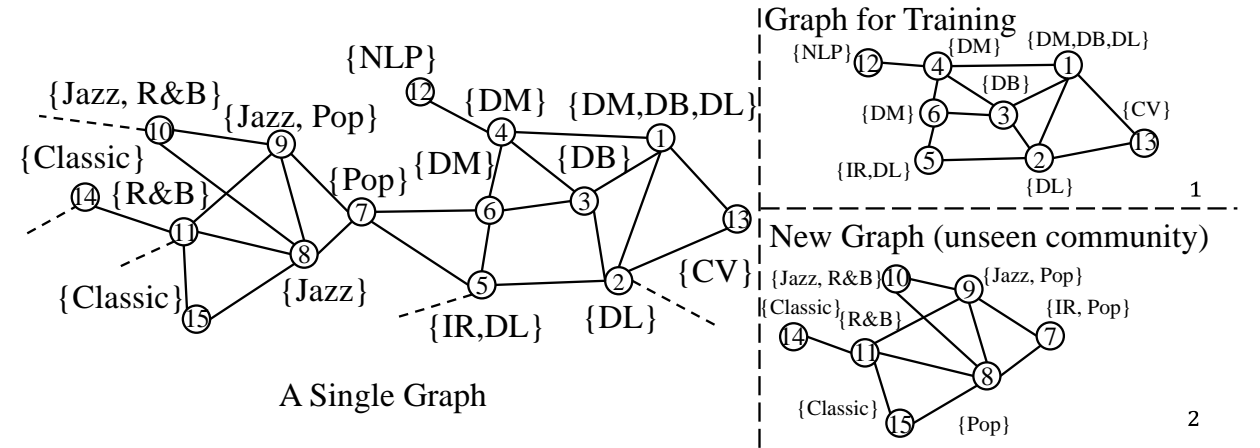


Fig 1. query-dependent community with different attributes.

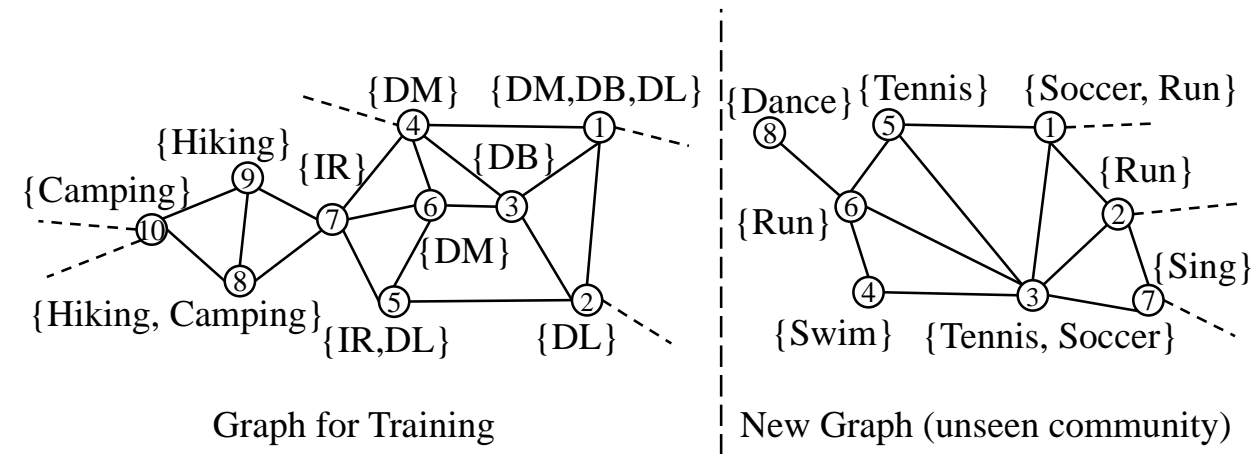
# Problem Statement

Empower the model to generalize and adapt to new **communities** and **graphs** by inductive learning

- For new communities
  - Queries from  $\{C_{q_1}, \dots, C_{q_i}\}$  in graph  $\mathcal{G}$  for training
  - Queries  $q^*$  from a new community  $C_{q^*}$  for test
  - i.e.,  $C_{q_1} \cap C_{q^*} = \emptyset, \dots, C_{q_i} \cap C_{q^*}$
- For new graphs
  - Queries from graph  $\mathcal{G}$  for training
  - Queries from new graph  $\mathcal{G}^*$  for test

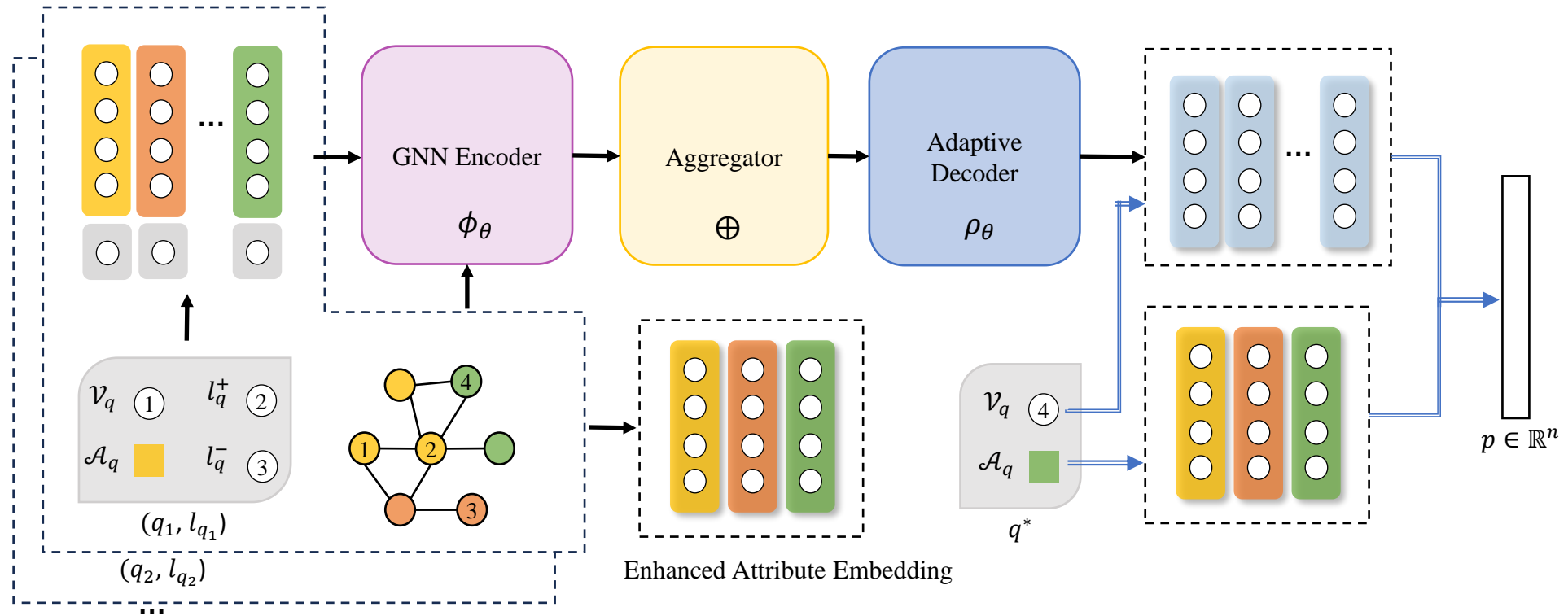


(a) Inductive Setting for Communities



(b) Inductive Setting for Graphs

# IACS Architecture



# IACS

$$p(y_{q^*} | q^*, \mathcal{T}) = \rho_{\theta} \left( q^*, \bigoplus_{(q, l_q) \in (Q, L)} \phi_{\theta}(q, l_q) \right)$$

## GNN Encoder $\phi_{\theta}(q, l_q)$

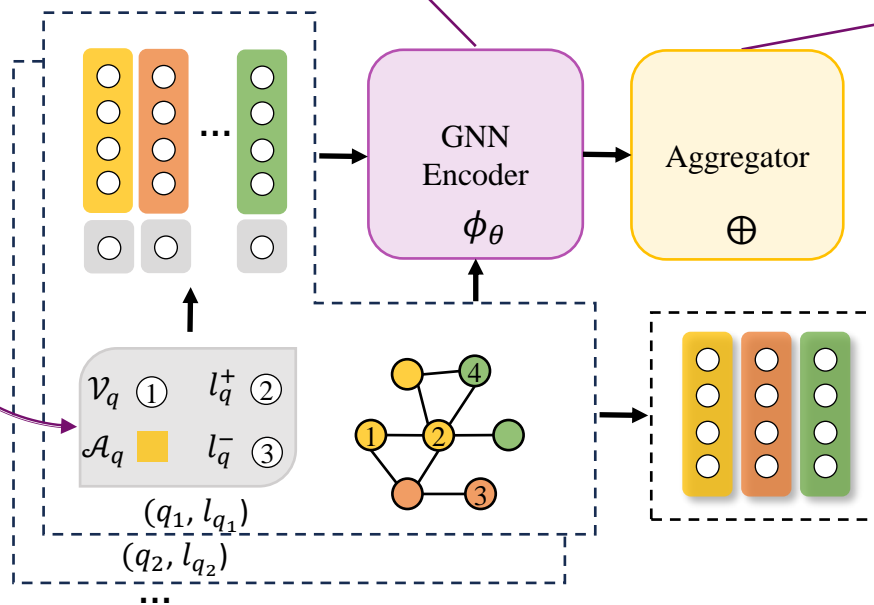
- $k$ -th layer aggregate:  $a^{(k)}(v) \leftarrow f_{\mathcal{A}}^{(k)}(\{h^{(k-1)}(u) \mid u \in N(v)\})$
- $k$ -th layer combine:  $h^{(k)}(v) \leftarrow f_{\mathcal{C}}^{(k)}(h^{(k-1)}(v), a^{(k)}(v))$

Aggregator  $\bigoplus$   
 Permutation invariant operator,  
 average:  $H = \frac{1}{|Q|} \sum_{q \in Q} H_q$

## Inputs

$$h^{(0)}(v) = [I_l(v) \parallel e(v)],$$

$$I_l(v) = \begin{cases} 1 & v \in l_q^+ \cup \{v_q\} \\ 0 & \text{otherwise} \end{cases}$$



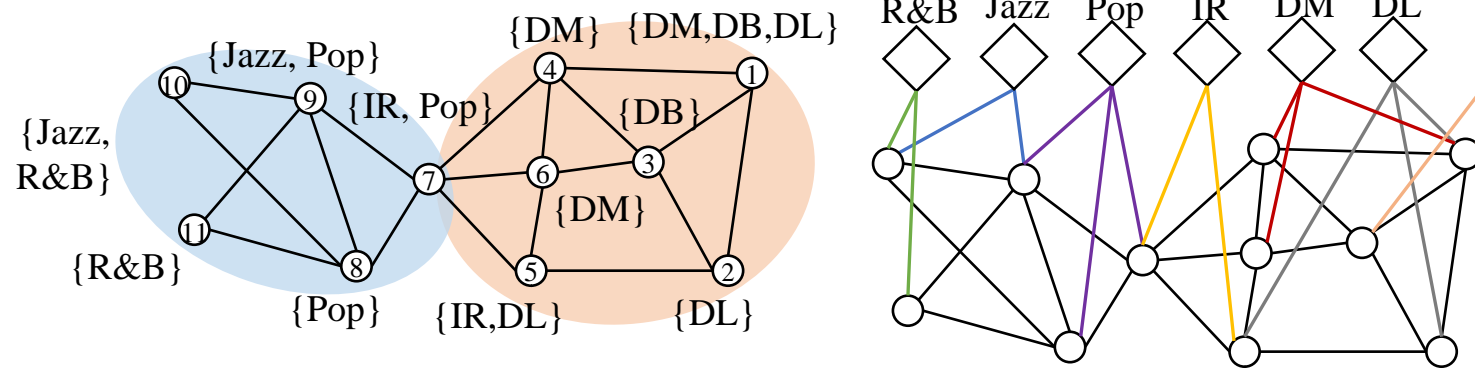


# Enhanced Attribute Encoding

Construct attribute-augmented graph

- $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{A}) \rightarrow \mathcal{G}_{\mathcal{A}} = (\mathcal{V} \cup \mathcal{V}_{\mathcal{A}}, \mathcal{E} \cup \mathcal{E}_{\mathcal{A}})$
- Use a scalable, task-independent graph embedding algorithm, ProNE
- Pretrain a node embedding for the attributed-augmented graph  $\mathcal{G}_{\mathcal{A}}$

$$e(v) = \sum_{a \in \mathcal{A}(v)} e_a$$



# IACS

## Gating mechanism

$$\delta = \text{sigmoid}(W_\delta H), \epsilon = W_\epsilon H,$$

$$\hat{\gamma} = \gamma \odot \delta + \epsilon \odot (1 - \delta), \hat{\beta} = \beta \odot \delta + \epsilon \odot (1 - \delta),$$

$$\hat{H} = \hat{\gamma} \odot H + \hat{\beta}$$

## Adaptive Decoder $\rho_\theta(q^*, H)$

- Feature-wise Linear Modulation (FiLM)

$$\gamma = W_\gamma H, \beta = W_\beta H,$$

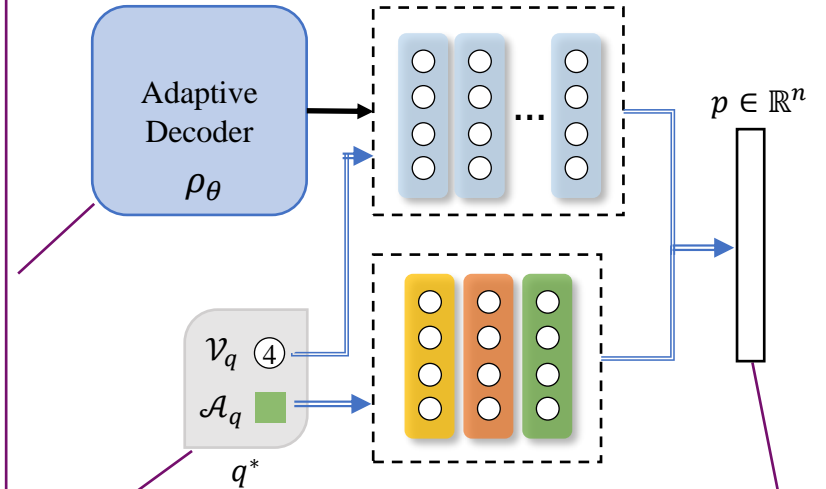
$$\hat{H} = \gamma \odot H + \beta$$

- Concatenate query node embedding and query attribute embedding

$$e_{\mathcal{V}_{q^*}} = \frac{1}{|\mathcal{V}_{q^*}|} \sum_{v \in \mathcal{V}_{q^*}} \hat{H}(v), e_{\mathcal{A}_{q^*}} = \frac{1}{|\mathcal{A}_{q^*}|} \sum_{a \in \mathcal{A}_{q^*}} e_a$$

$$e_{q^*} = \text{MLP}(e_{\mathcal{V}_{q^*}} || e_{\mathcal{A}_{q^*}})$$

- Inner Product Decoder:  $p(\hat{l}_{q^*} | q^*, \mathcal{T}) = \text{sigmoid}(\langle e_{q^*}, \hat{H} \rangle)$



A new ACS query

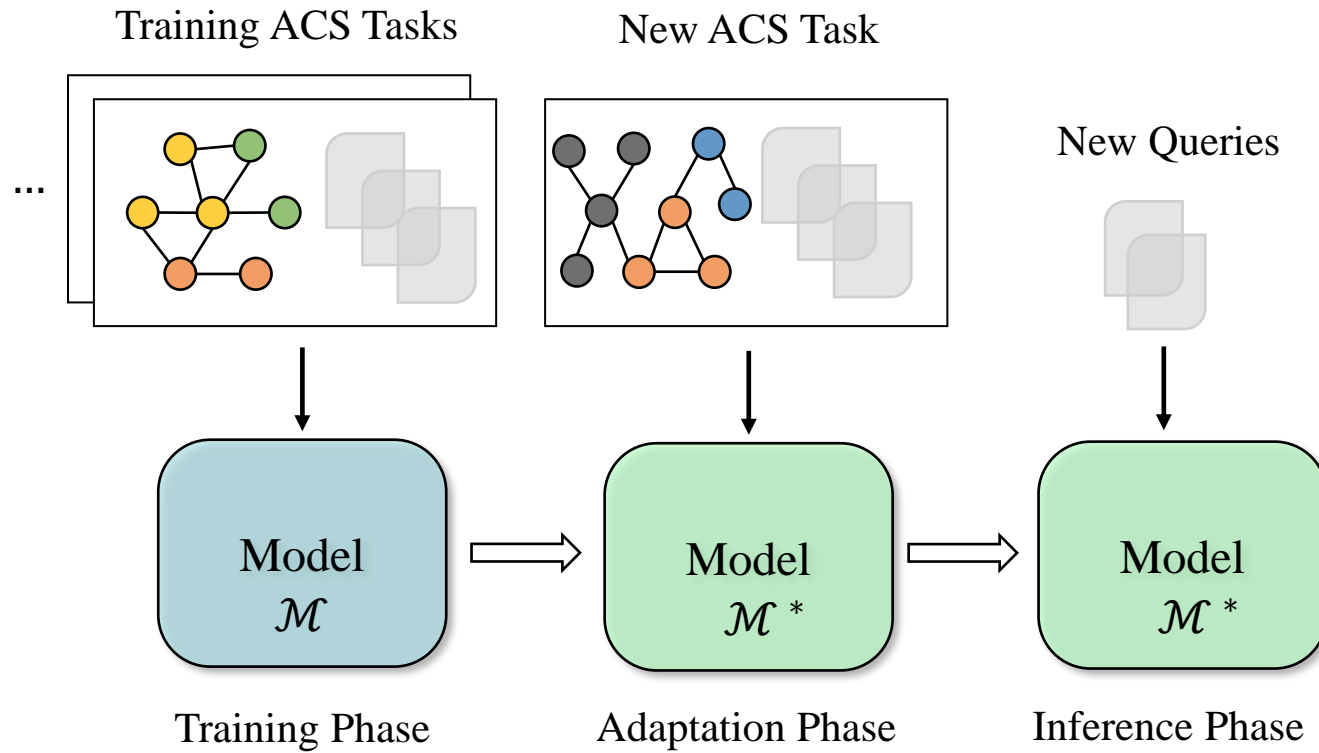
$$q^* = (\mathcal{V}_{q^*}, \mathcal{A}_{q^*})$$

## Binary Cross Entropy loss

$$\mathcal{L} = \sum_{\mathcal{T}_i \in \mathcal{D}} \sum_{(q, l_q) \in (Q_i, L_i)} -\log p(\hat{l}_q | q, \mathcal{T}_i)$$

$$= \sum_{\mathcal{T}_i \in \mathcal{D}} \sum_{(q, l_q) \in (Q_i, L_i)} (-\sum_{v^+ \in l_q^+} \log(\hat{y}(v^+)) - \sum_{v^- \in l_q^-} \log(1 - \hat{y}(v^-)))$$

# IACS Workflow



# Experimental Studies: Setup

- Model
  - encoder: 3-layer GCN, GraphSAGE, GIN, **GAT**.
  - decoder: FiLM + Inner product, FiLM with gating mechanism + Inner product, inner product
- Baselines
  - 3 algorithmic approaches
  - 3 supervised-learning based approaches
  - 2 meta-learning approaches

**Table 3: The Profiles of Dataset**

Dataset	$ \mathcal{G} $	$ \mathcal{V} $	$ \mathcal{E} $	$ \mathcal{A} $	$ \mathcal{C} $	graph des.	attribute des.	community des.	# tasks
Arxiv [54]	1	169,343	1,166,243	N/A	40	paper citation	NA	research topics	1,000
Amazon2M [13]	1	2,449,029	61,859,140	N/A	47	product co-purchasing	NA	product categories	5,000
Cora [57]	1	2,708	5,429	1,433	7	paper citation	paper keywords	research topics	192
Citeseer [57]	1	3,327	4,732	3,703	6	paper citation	paper keywords	research topics	192
Reddit [24]	1	232,965	114,615,892	1,164	50	post co-comment	synthetic	post categories	1,000
Facebook [33]	10	4,039	88,234	2,281	193	social friendship	user profiles	friend circles	10
Twitter [33]	973	81,306	1,768,149	512,985	4,065	social friendship	user profiles	friend circles	973

# Experiential Studies: Effectiveness

**Table 4: Overall Performance on Non-Attributed CS (%)**

Dataset	Approach	4-shot			8-shot		
		Pre	Rec	F1	Pre	Rec	F1
Arxiv	<i>CTC</i>	54.23±0.53	2.16±0.04	4.15±0.09	54.04±0.72	2.16±0.05	4.15±0.09
	<i>ICS-GNN</i>	62.72±0.26	21.09±0.07	31.57±0.10	62.53±0.36	21.12±0.05	31.57±0.08
	<i>QD-GNN</i>	59.97±0.41	83.60±1.18	69.84±0.47	58.91±0.29	89.62±1.14	71.09±0.29
	<i>Supervise</i>	67.99±0.33	69.78±1.49	68.87±0.86	69.09±0.39	74.29±0.98	71.60±0.38
	<i>MAML</i>	63.51±1.07	60.25±2.50	61.81±1.42	62.77±0.72	60.34±3.64	61.48±1.82
	<i>FeatTrans</i>	65.35±0.64	55.18±1.81	59.81±0.88	64.18±0.69	55.42±1.09	59.47±0.74
	<i>IACS</i>	63.65±0.62	89.26±1.05	74.31±0.37	64.14±0.49	90.21±1.31	74.97±0.31
	<i>IACS-G</i>	59.75±0.42	97.99±0.76	74.23±0.13	65.06±0.81	88.12±1.65	74.84±0.36
	<i>IACS-P</i>	61.99±2.56	92.63±5.77	74.12±0.21	65.45±0.42	87.07±0.53	74.72±0.24
Amazon2M	<i>CTC</i>	80.30±0.35	4.06±0.02	7.73±0.04	80.27±0.27	4.06±0.01	7.73±0.02
	<i>ICS-GNN</i>	79.50±0.27	6.55±0.01	12.11±0.02	79.63±0.29	6.55±0.02	12.11±0.03
	<i>QD-GNN</i>	75.46±0.33	95.15±0.53	84.17±0.04	75.33±0.26	96.68±0.13	84.67±0.21
	<i>Supervise</i>	83.86±0.09	77.07±0.44	80.32±0.25	84.46±0.35	80.18±0.52	82.27±0.29
	<i>MAML</i>	78.48±1.62	65.83±8.70	71.38±5.59	79.13±0.88	62.38±4.76	69.66±2.83
	<i>FeatTrans</i>	78.41±0.92	57.89±1.39	66.60±1.14	78.69±0.34	57.18±1.22	66.22±0.72
	<i>IACS</i>	80.52±0.34	93.42±0.83	86.48±0.22	81.44±0.75	93.34±1.07	86.97±0.21
	<i>IACS-G</i>	79.92±0.31	94.25±0.93	86.49±0.21	80.49±0.16	94.60±0.52	86.98±0.24
	<i>IACS-P</i>	79.63±0.88	94.77±1.09	86.54±0.29	80.62±0.81	94.86±0.74	87.16±0.26

Non-attributed CS

- a) IACS models consistently outperform all the baselines.
- b) The superiority of IACS is primarily evident in its significant improvement in recall (+1.28% compared to the best baseline) while maintaining a relatively high precision (59.75% ~ 81.44%).



# Experiential Studies: Effectiveness

**Table 5: Overall Performance on ACS in Single Graph (%)**

Dataset	Approach	4-shot			8-shot		
		Pre	Rec	F1	Pre	Rec	F1
Citeseer	<i>ATC</i>	58.99 $\pm$ 0.87	5.01 $\pm$ 0.17	9.24 $\pm$ 0.29	57.83 $\pm$ 0.36	4.97 $\pm$ 0.25	9.16 $\pm$ 0.42
	<i>ACQ</i>	70.59 $\pm$ 2.14	6.97 $\pm$ 1.32	12.66 $\pm$ 2.19	69.15 $\pm$ 1.43	6.79 $\pm$ 0.99	12.36 $\pm$ 1.64
	<i>AQD-GNN</i>	52.70 $\pm$ 1.04	84.29 $\pm$ 7.59	64.77 $\pm$ 2.78	52.26 $\pm$ 2.06	85.15 $\pm$ 5.57	64.74 $\pm$ 3.01
	<i>Supervise</i>	60.45 $\pm$ 2.00	63.20 $\pm$ 1.53	61.79 $\pm$ 1.67	62.30 $\pm$ 1.88	66.74 $\pm$ 0.84	64.43 $\pm$ 1.16
	<i>MAML</i>	55.53 $\pm$ 1.61	43.18 $\pm$ 4.27	48.46 $\pm$ 2.45	56.15 $\pm$ 0.94	45.02 $\pm$ 3.05	49.93 $\pm$ 1.94
	<i>FeatTrans</i>	58.67 $\pm$ 2.54	37.97 $\pm$ 1.38	46.08 $\pm$ 1.51	58.44 $\pm$ 1.74	39.86 $\pm$ 2.17	47.38 $\pm$ 1.91
	<i>IACS</i>	64.74 $\pm$ 1.55	71.15 $\pm$ 1.45	67.78 $\pm$ 0.94	67.06 $\pm$ 1.45	71.48 $\pm$ 1.84	69.19 $\pm$ 1.41
	<i>IACS-G</i>	65.75 $\pm$ 0.54	70.12 $\pm$ 1.33	<b>67.86<math>\pm</math>0.56</b>	67.48 $\pm$ 1.62	71.90 $\pm$ 1.98	<b>69.59<math>\pm</math>0.87</b>
	<i>IACS-P</i>	65.52 $\pm$ 1.15	70.37 $\pm$ 1.43	<u>67.84<math>\pm</math>0.66</u>	67.25 $\pm$ 1.86	72.08 $\pm$ 0.74	<u>69.57<math>\pm</math>0.98</u>
Reddit	<i>ATC</i>	82.84 $\pm$ 0.77	39.15 $\pm$ 1.68	53.16 $\pm$ 1.63	83.73 $\pm$ 0.87	39.00 $\pm$ 0.91	53.21 $\pm$ 0.79
	<i>ACQ</i>	97.74 $\pm$ 0.27	22.82 $\pm$ 1.02	36.99 $\pm$ 1.35	98.16 $\pm$ 0.03	22.49 $\pm$ 1.20	36.58 $\pm$ 1.59
	<i>AQD-GNN</i>	85.88 $\pm$ 1.63	89.54 $\pm$ 1.89	<b>87.67<math>\pm</math>1.38</b>	85.51 $\pm$ 1.73	92.20 $\pm$ 1.39	<b>88.73<math>\pm</math>1.42</b>
	<i>Supervise</i>	86.16 $\pm$ 1.38	78.14 $\pm$ 1.93	81.95 $\pm$ 1.51	87.14 $\pm$ 0.95	80.15 $\pm$ 0.49	83.50 $\pm$ 0.49
	<i>MAML</i>	88.30 $\pm$ 0.78	64.46 $\pm$ 3.15	74.66 $\pm$ 2.27	OOM	OOM	OOM
	<i>FeatTrans</i>	87.76 $\pm$ 2.41	34.78 $\pm$ 3.22	49.72 $\pm$ 3.34	88.07 $\pm$ 0.59	36.46 $\pm$ 2.43	51.53 $\pm$ 2.40
	<i>IACS</i>	84.03 $\pm$ 1.20	84.11 $\pm$ 3.84	84.01 $\pm$ 1.26	83.50 $\pm$ 1.69	85.07 $\pm$ 4.36	84.20 $\pm$ 1.49
	<i>IACS-G</i>	83.81 $\pm$ 1.73	85.33 $\pm$ 2.23	<u>84.54<math>\pm</math>1.14</u>	86.07 $\pm$ 0.86	84.01 $\pm$ 2.63	<u>85.00<math>\pm</math>1.11</u>
	<i>IACS-P</i>	84.85 $\pm$ 1.28	83.90 $\pm$ 2.53	84.35 $\pm$ 1.31	85.89 $\pm$ 1.84	83.21 $\pm$ 1.72	84.51 $\pm$ 1.33

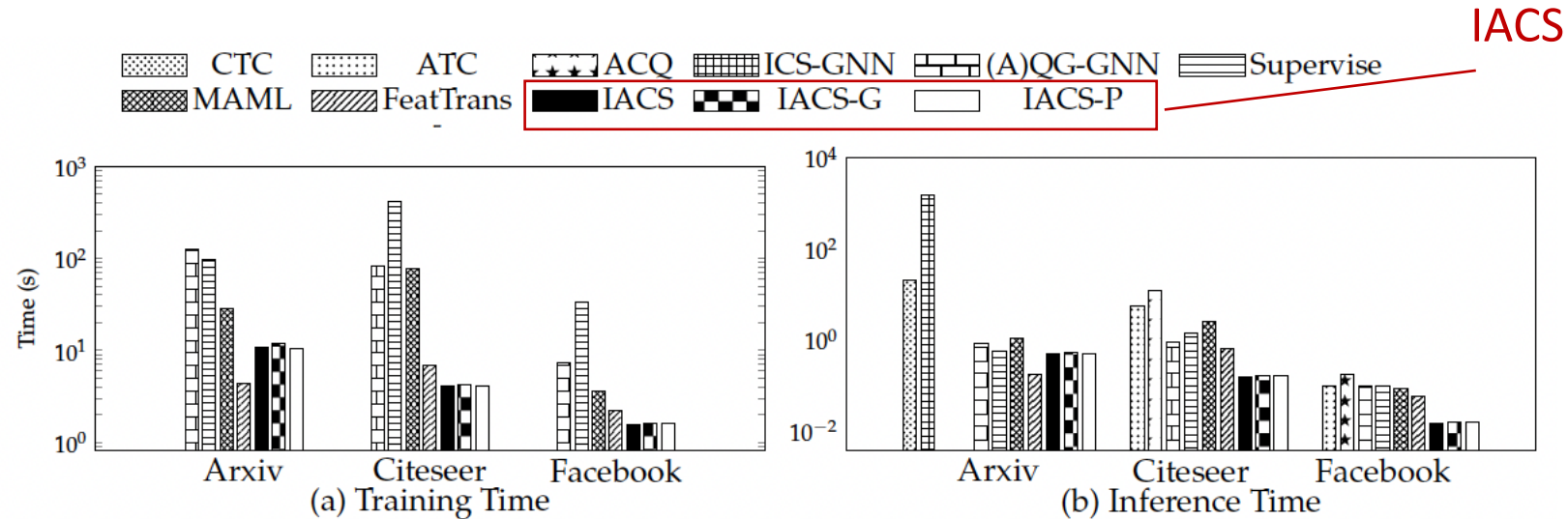
**Table 6: Overall Performance on ACS in Multiple Graphs (%)**

Dataset	Approach	4-shot			8-shot		
		Pre	Rec	F1	Pre	Rec	F1
Facebook	<i>ATC</i>	60.23 $\pm$ 5.10	11.99 $\pm$ 0.88	19.97 $\pm$ 1.26	41.14 $\pm$ 4.18	11.11 $\pm$ 3.27	17.22 $\pm$ 3.71
	<i>ACQ</i>	38.86 $\pm$ 3.52	66.92 $\pm$ 5.79	48.90 $\pm$ 1.92	40.60 $\pm$ 3.06	64.00 $\pm$ 3.33	49.65 $\pm$ 3.07
	<i>AQD-GNN</i>	37.71 $\pm$ 1.65	96.70 $\pm$ 5.93	54.26 $\pm$ 2.59	36.71 $\pm$ 1.03	96.29 $\pm$ 4.68	53.14 $\pm$ 1.75
	<i>Supervise</i>	59.32 $\pm$ 1.22	79.61 $\pm$ 5.37	67.95 $\pm$ 2.73	64.34 $\pm$ 1.63	83.38 $\pm$ 3.20	72.59 $\pm$ 1.08
	<i>MAML</i>	47.06 $\pm$ 6.12	89.04 $\pm$ 3.86	59.85 $\pm$ 8.12	46.12 $\pm$ 3.30	73.30 $\pm$ 5.89	56.44 $\pm$ 2.45
	<i>FeatTrans</i>	50.11 $\pm$ 6.62	68.34 $\pm$ 8.66	56.84 $\pm$ 4.28	50.82 $\pm$ 1.16	59.77 $\pm$ 6.87	72.16 $\pm$ 3.40
	<i>IACS</i>	85.61 $\pm$ 2.16	79.55 $\pm$ 4.13	<b>78.09<math>\pm</math>4.09</b>	66.13 $\pm$ 1.62	84.33 $\pm$ 3.80	74.08 $\pm$ 1.34
	<i>IACS-G</i>	85.75 $\pm$ 2.27	81.31 $\pm$ 4.51	<u>77.42<math>\pm</math>3.82</u>	64.87 $\pm$ 1.67	88.17 $\pm$ 3.04	<u>74.72<math>\pm</math>1.64</u>
	<i>IACS-P</i>	81.82 $\pm$ 4.42	80.79 $\pm$ 2.69	<u>77.37<math>\pm</math>4.62</u>	65.92 $\pm$ 4.20	87.36 $\pm$ 1.66	<b>75.05<math>\pm</math>2.31</b>
Twitter2Facebook	<i>ATC</i>	77.92 $\pm$ 6.75	12.88 $\pm$ 0.74	22.08 $\pm$ 1.04	70.30 $\pm$ 9.71	11.25 $\pm$ 1.91	19.35 $\pm$ 3.05
	<i>ACQ</i>	68.89 $\pm$ 0.81	37.21 $\pm$ 8.31	49.51 $\pm$ 4.76	20.38 $\pm$ 0.58	44.58 $\pm$ 3.51	28.22 $\pm$ 0.47
	<i>AQD-GNN</i>	37.49 $\pm$ 0.49	96.81 $\pm$ 3.49	37.49 $\pm$ 1.03	37.61 $\pm$ 3.45	95.10 $\pm$ 8.97	53.84 $\pm$ 4.54
	<i>Supervise</i>	58.12 $\pm$ 4.20	80.28 $\pm$ 8.12	58.12 $\pm$ 5.27	62.32 $\pm$ 4.41	81.44 $\pm$ 4.99	70.55 $\pm$ 4.15
	<i>MAML</i>	38.12 $\pm$ 0.42	97.57 $\pm$ 1.67	38.12 $\pm$ 0.37	37.97 $\pm$ 1.58	95.01 $\pm$ 4.40	54.20 $\pm$ 1.19
	<i>FeatTrans</i>	38.45 $\pm$ 0.42	98.05 $\pm$ 1.04	38.45 $\pm$ 0.43	38.20 $\pm$ 0.86	97.06 $\pm$ 1.39	54.81 $\pm$ 0.72
	<i>IACS</i>	72.10 $\pm$ 5.43	85.44 $\pm$ 4.30	<u>72.10<math>\pm</math>1.58</u>	66.29 $\pm$ 3.99	84.90 $\pm$ 0.58	73.68 $\pm$ 2.88
	<i>IACS-G</i>	67.76 $\pm$ 3.00	86.17 $\pm$ 1.65	67.76 $\pm$ 1.80	64.86 $\pm$ 1.81	86.36 $\pm$ 3.48	<u>74.05<math>\pm</math>1.94</u>
	<i>IACS-P</i>	72.38 $\pm$ 4.80	83.69 $\pm$ 6.53	72.38 $\pm$ 1.75	65.14 $\pm$ 0.85	86.07 $\pm$ 2.36	<b>74.28<math>\pm</math>1.19</b>

## ACS

- In general, IACS achieves the highest F1 score in most cases (5 out of 6), even when the graphs of training and inference are from different datasets.
- The improvement in the 8-shot setting is relatively lower compared to the 4-shot setting.

# Experiential Studies: Efficiency



**Figure 5: Comparison of Training and Inference Time**

- IACS models exhibit faster training and inference time compared to other baselines in most datasets.

# Experiential Studies: Streaming Model Adaptation

- The streaming adaptation process leads to higher F1 scores for the streaming model than the original model across a wide range of sequential tasks.
- The results indicate that three IACS models exhibit an improvement ratio of 3% in the streaming adaptation model.

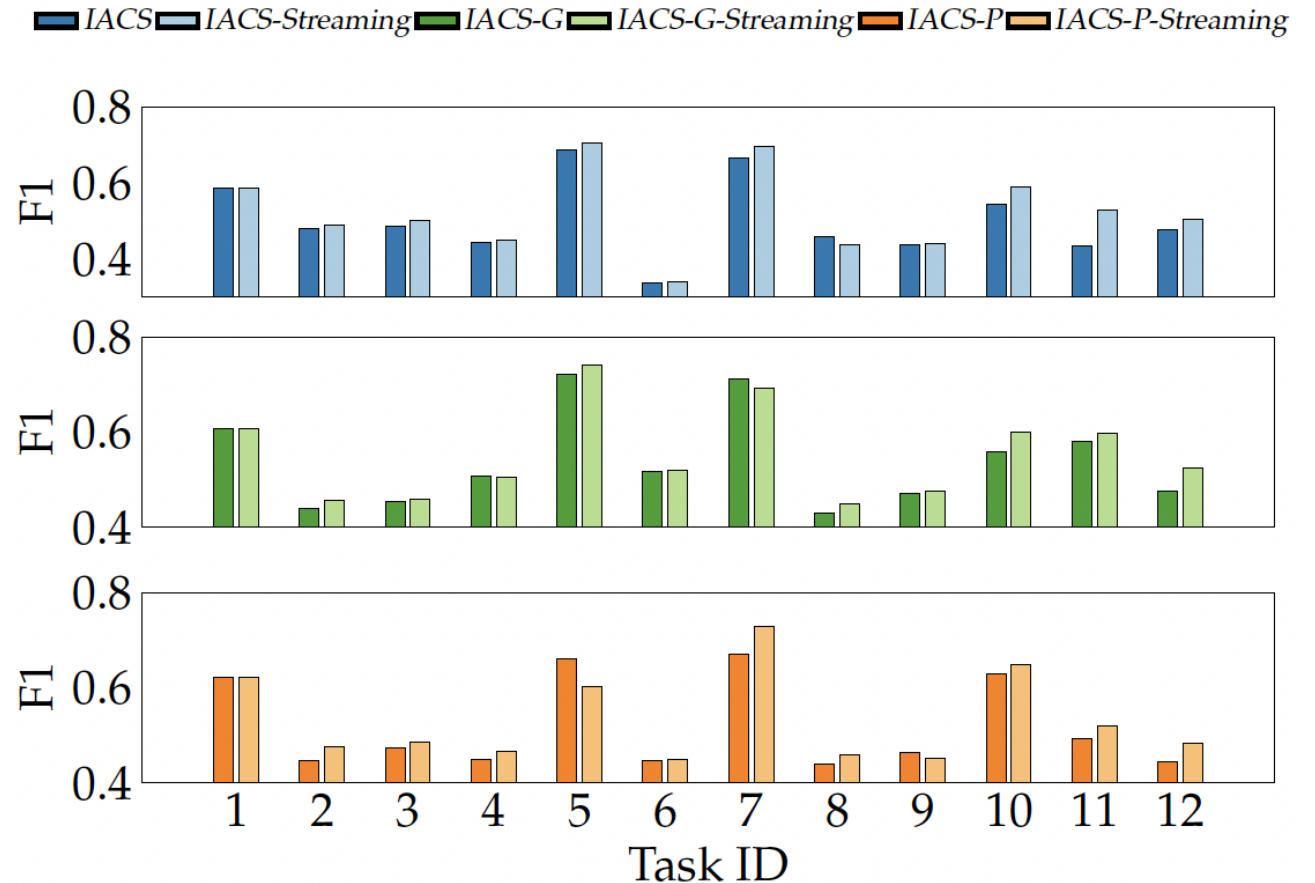
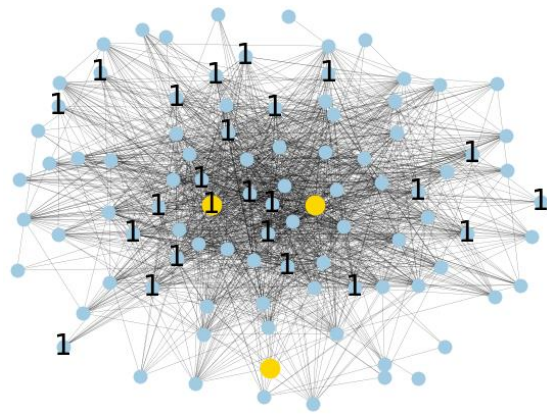


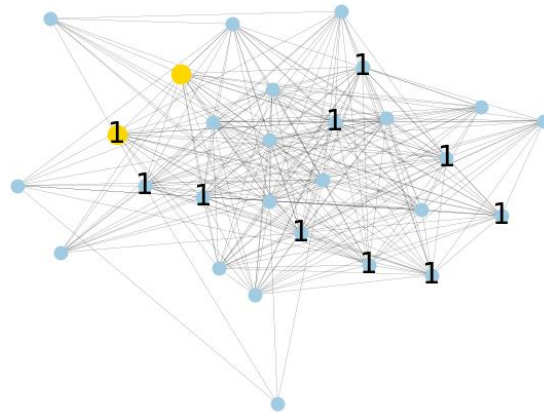
Figure 8: Streaming Model Adaptation on Twitter



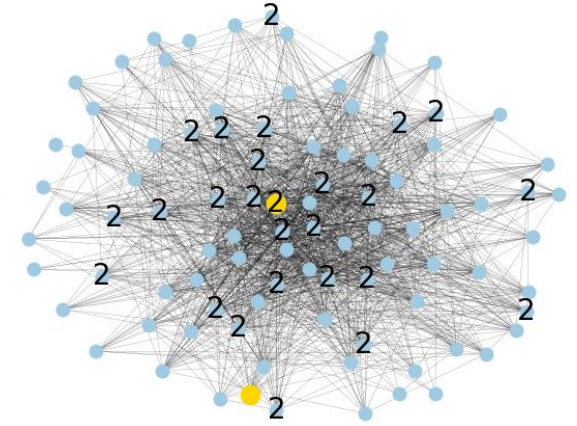
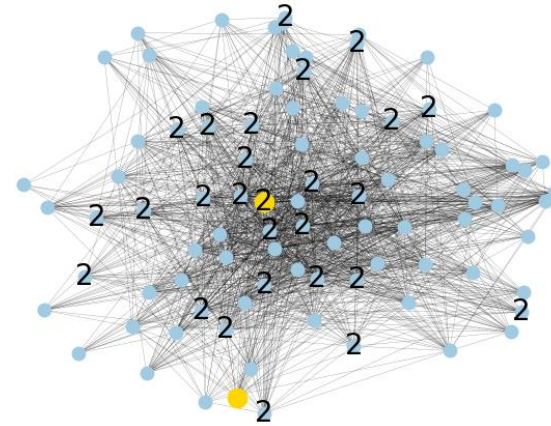
# Experiential Studies: Case Study



(a) Twitter: Training Communities



(b) Ground-truth Community



(c) Predicted Community

- These communities exhibit variations, characterized by heterogeneous topological structures and attribute distributions.
- We observe a notable overlap between the identified communities and the ground-truth, thus confirming the accuracy of our predictions.

# Summary

- Leveraging ML/DL based approaches for attributed community search
- Existing learning approaches have limited inductive ability and cannot deal with complex queries.
- Propose Inductive Attributed Community Search (IACS) to infer new queries for different communities/graphs
- Propose a three-stage workflow to fulfill inductive ACS
- IACS achieves better performance on effectiveness and efficiency



# Thank you!

Shuheng Fang, The Chinese University of Hong Kong.

[shfang@se.cuhk.edu.hk](mailto:shfang@se.cuhk.edu.hk) 😊