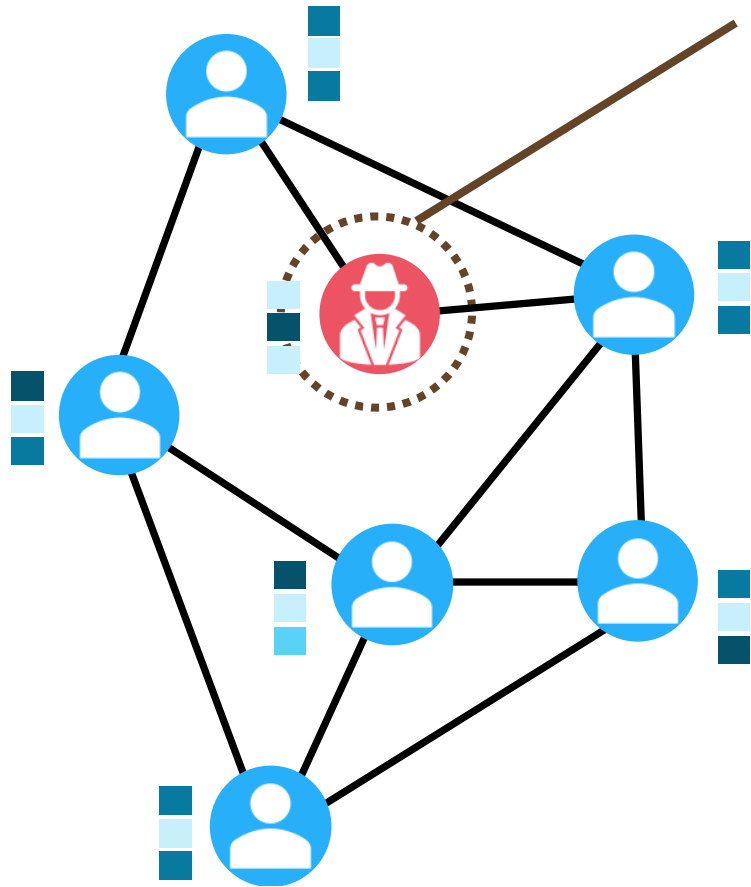# ARC: A Generalist Graph Anomaly Detector with In-Context Learning

**Presenter: Yixin Liu**
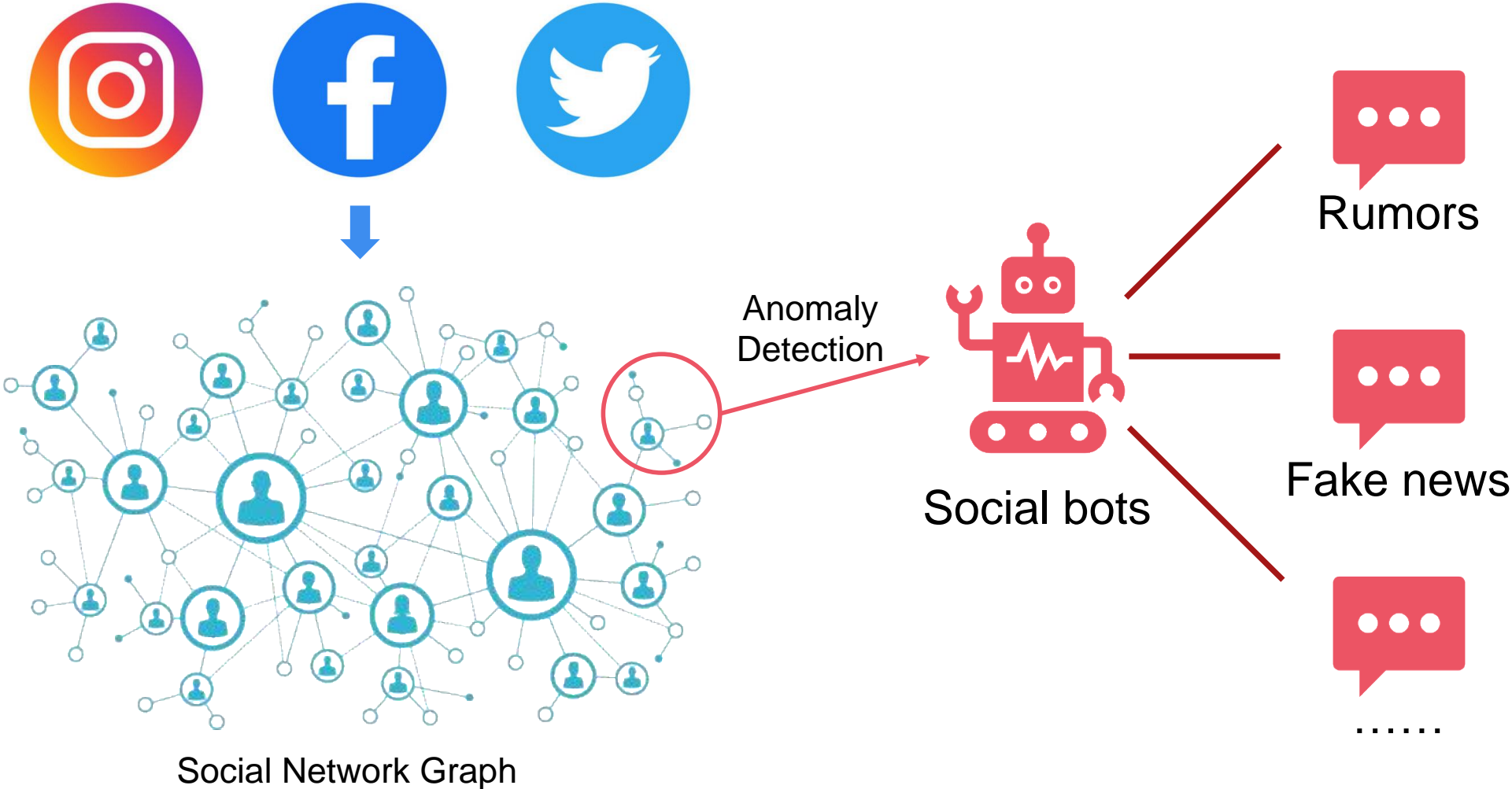
# Graph Anomaly Detection (GAD)
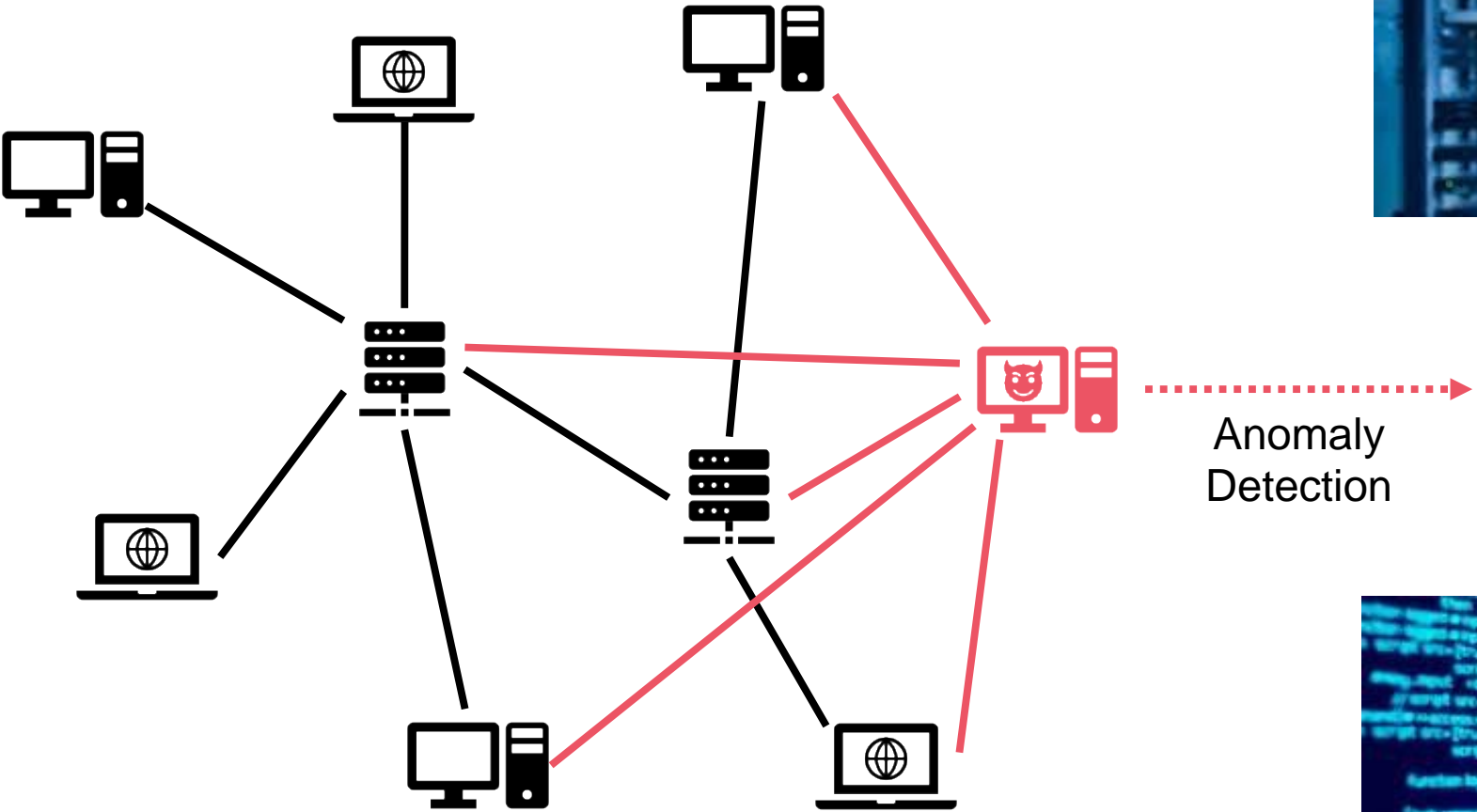


To detect the abnormal nodes that are <u>different</u> from the majority.

# GAD's Application: Social Networks

Anomaly Detection

Social bots

Rumors

Fake news

......

Social Network Graph

# GAD's Application: Cybersecurity



Hackers

Anomaly Detection

Cyber Attacks

# GAD's Application: Traffic Networks



Anomaly Detection

Spatial-temporal Graph

Traffic sensors displayed on GoogleMaps

Accident?

Congestion?

# GAD: Existing Solutions

Mainstream solution: Graph neural networks (GNNs)



Graph data

Predicted
Anomaly Scores

# GAD: Existing Solutions

Graph neural networks (GNNs) based methods



Supervised GAD methods

Unsupervised GAD methods

# GAD: Existing Solutions

Graph neural networks (GNNs) based methods

**Supervised GAD methods:**
training GAD model with labels (normal/anomaly)

Unsupervised GAD methods



CARE-GNN[1]

GHRN[2]

[1] Dou, Yingtong, et al. "Enhancing graph neural network-based fraud detectors against camouflaged fraudsters." *Proceedings of the 29th ACM international conference on information & knowledge management*. 2020.
[2] Gao, Yuan, et al. "Addressing heterophily in graph anomaly detection: A perspective of graph spectrum." *Proceedings of the ACM Web Conference 2023*. 2023.

# GAD: Existing Solutions

Graph neural networks
(GNNs) based methods



**Supervised GAD methods:**
training GAD model with labels (normal/anomaly)
**Unsupervised GAD methods:**
training GAD model without labels



DOMINANT[3]



CoLA[4]

[3] Ding, Kaize, et al. "Deep anomaly detection on attributed networks." Proceedings of the 2019 SIAM international conference on data mining. Society for Industrial and Applied Mathematics, 2019.
[4] Liu, Yixin, et al. "Anomaly detection on attributed networks via contrastive self-supervised learning." IEEE transactions on neural networks and learning systems 33.6 (2021): 2378-2392.

# GAD: Existing Solutions

Graph neural networks (GNNs) based methods

**Supervised GAD methods**

**Unsupervised GAD methods**

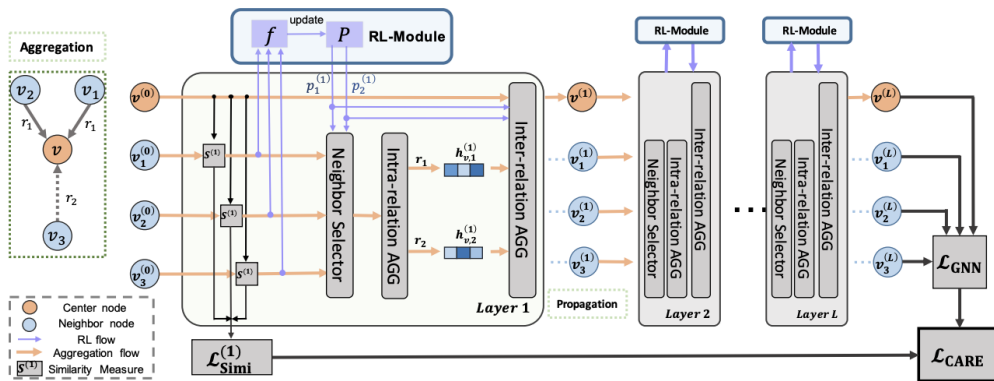Learning paradigm: one model for one dataset



(a) Supervised GAD Paradigm

(b) Unsupervised GAD Paradigm

# GAD: Existing Solutions

Graph neural networks (GNNs) based methods
{ **Supervised GAD methods**
**Unsupervised GAD methods**



(a) Supervised GAD Paradigm

(b) Unsupervised GAD Paradigm

Learning paradigm:
one model for one dataset

✕ Expensive training cost

Devices            Time

# GAD: Existing Solutions

Graph neural networks (GNNs) based methods

**Supervised GAD methods**

**Unsupervised GAD methods**



✓ Normal Node  🚨 Abnormal Node  📄🛍️👤 Unlabeled Node

Training Stage | Inference Stage

Train → Supervised GAD Model → GAD Model →

(a) Supervised GAD Paradigm

Train → Unsupervised GAD Model → GAD Model →

(b) Unsupervised GAD Paradigm

Learning paradigm:
one model for one dataset

✗ Expensive training cost

✗ Data requirements

Labels

Full Data

# GAD: Existing Solutions

Graph neural networks
(GNNs) based methods

{

**Supervised GAD methods**

**Unsupervised GAD methods**



✓ Normal Node    🔔 Abnormal Node    Unlabeled Node

Training Stage          Inference Stage

Train → Supervised GAD Model → GAD Model

(a) Supervised GAD Paradigm

Train → Unsupervised GAD Model → GAD Model

(b) Unsupervised GAD Paradigm

## Learning paradigm:
one model for one dataset

✖ Expensive training cost

✖ Data requirements

✖ Poor generalizability

**Can't be transferred to new datasets!**

# GAD: Existing Solutions

Graph neural networks
(GNNs) based methods {
**Supervised GAD methods**
**Unsupervised GAD methods**



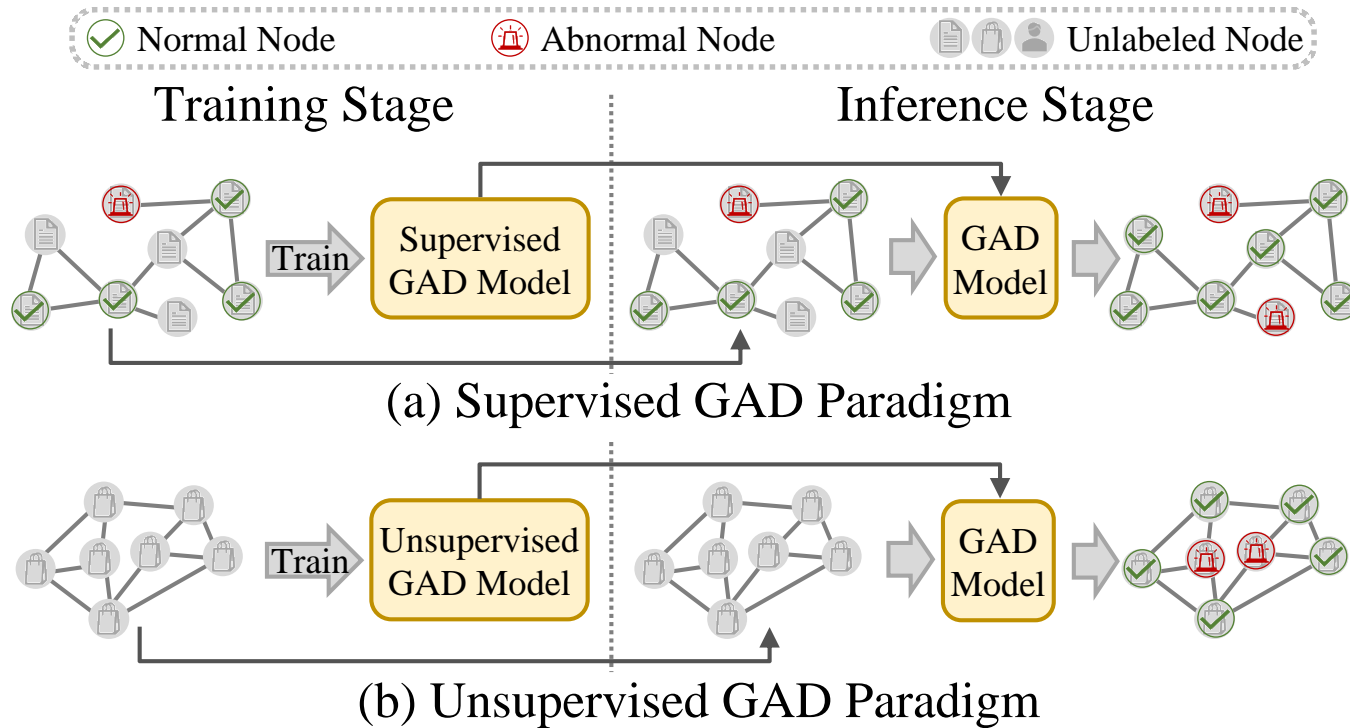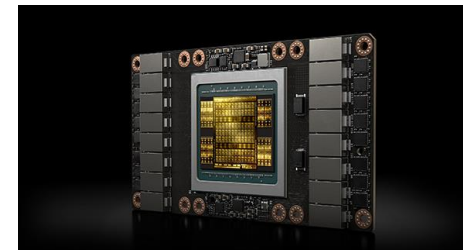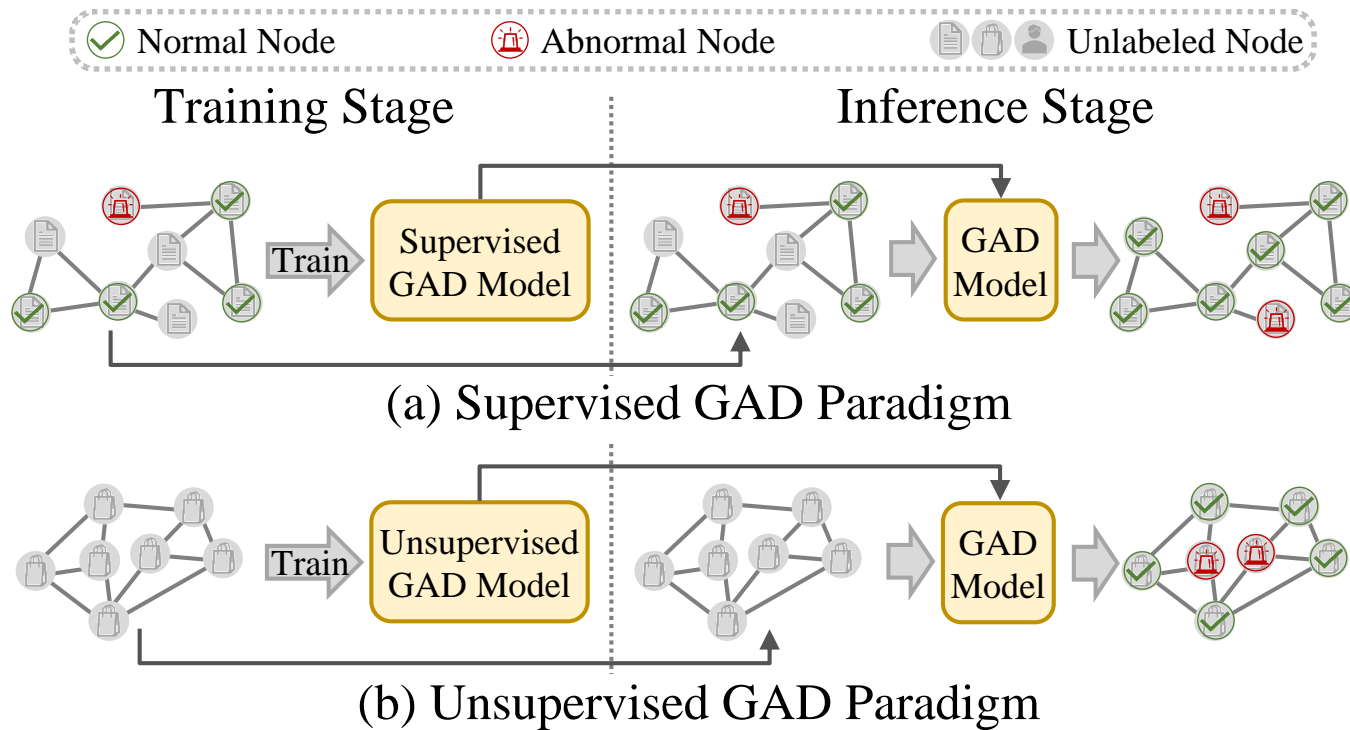○ Normal Node    🔔 Abnormal Node    Unlabeled Node

Training Stage          Inference Stage

Train    Supervised          GAD
         Model               Model

(a) Supervised GAD Paradigm

Train → Unsupervised GAD Model → GAD Model

(b) Unsupervised GAD Paradigm

Learning paradigm:
one model for one dataset

✖ Expensive training cost

✖ Data requirements

✖ Poor generalizability

Can't be transferred to new datasets!

**Can we develop a <u>one-for-all</u> GAD model that can be trained once and effectively applied across various datasets?**

# Generalist GAD: a New Paradigm



Training on multiple datasets ➡️ Directly inference on various datasets without re-training or fine-tuning

The "foundation model" of GAD!

# Generalist GAD: a New Paradigm



Ours "generalist GAD" paradigm

Training on multiple datasets

⬇

Directly inference on various datasets

☑ No fine-tuning

→ low application costs

☑ Only need few-shot normal

→ low data requirement

☑ Great generalizability

→ one-for-all model

# Generalist GAD: a New Paradigm

Training on multiple datasets

Directly inference on various datasets

Training Stage

Inference Stage

Train

Generalist GAD Model

GAD Model

☑ No fine-tuning

Ours "generalist GAD" paradigm

☑ Only need few-shot normal

→ low data requirement

☑ Great generalizability

→ one-for-all model

# How to design a generalist GAD model?

# The proposed generalist GAD method - ARC



Smoothness-Based Feature **A**lignment

Ego-Neighbor **R**esidual Graph Encoder

Cross-Attentive In-**C**ontext Anomaly Scoring

Alignment ➤ Encoding ➤ Scoring

# The proposed generalist GAD method - ARC



**Step 1**: Smoothness-Based Feature **A**lignment

- **Feature projection**

$$\tilde{\mathbf{X}}^{(i)} \in \mathbb{R}^{n^{(i)} \times d_u} = \text{Proj}\left(\mathbf{X}^{(i)}\right) = \mathbf{X}^{(i)}\mathbf{W}^{(i)},$$

Linear projection – PCA

# The proposed generalist GAD method - ARC



**Step 1**: Smoothness-Based Feature **A**lignment

- **Feature projection**

$$\tilde{\mathbf{X}}^{(i)} \in \mathbb{R}^{n^{(i)} \times d_u} = \mathrm{Proj}\left(\mathbf{X}^{(i)}\right) = \mathbf{X}^{(i)}\mathbf{W}^{(i)},$$

Linear projection – PCA

- **Smoothness-based feature sorting**

$$s_k(\mathbf{X}) = -\frac{1}{|\mathcal{E}|} \sum_{(v_i, v_j) \in \mathcal{E}} (\mathbf{X}_{ik} - \mathbf{X}_{jk})^2$$

Reorder the projected features according to s

# The proposed generalist GAD method - ARC



**Motivation**:

The contributions of features with low/high smoothness are similar across datasets!



(a) Cora      (b) Facebook

- **Smoothness-based feature sorting**

$$s_k(\mathbf{X}) = -\frac{1}{|\mathcal{E}|} \sum_{(v_i, v_j) \in \mathcal{E}} (\mathbf{X}_{ik} - \mathbf{X}_{jk})^2$$

Reorder the projected features according to s

# The proposed generalist GAD method - ARC



**Step 2**: Ego-Neighbor **R**esidual Graph Encoder

- **Propagation**

$$\mathbf{X}^{[l]} = \tilde{\mathbf{A}}\mathbf{X}^{[l-1]}$$

- **Transformation**

$$\mathbf{Z}^{[l]} = \mathrm{MLP}\left(\mathbf{X}^{[l]}\right)$$

- **Residual operation**

$$\mathbf{R}^{[l]} = \mathbf{Z}^{[l]} - \mathbf{Z}^{[0]}$$

- **Concatenation**

$$\mathbf{H} = [\mathbf{R}^{[1]}||\cdots||\mathbf{R}^{[L]}]$$

# The proposed generalist GAD method - ARC



**Step 2**: Ego-Neighbor **R**esidual Graph Encoder

- **Propagation**

$$\mathbf{X}^{[l]} = \tilde{\mathbf{A}}\mathbf{X}^{[l-1]}$$

- **Transformation**

$$\mathbf{Z}^{[l]} = \mathrm{MLP}\left(\mathbf{X}^{[l]}\right)$$

- **Residual operation**

$$\mathbf{R}^{[l]} = \mathbf{Z}^{[l]} - \mathbf{Z}^{[0]}$$

- **Concatenation**

$$\mathbf{H} = [\mathbf{R}^{[1]}||\cdots||\mathbf{R}^{[L]}]$$

**Motivation**:
- Residual $\rightarrow$ Local Affinity[5]

$$h(v_i) = \frac{1}{|\mathcal{N}(v_i)|} \sum_{v_j \in \mathcal{N}(v_i)} \mathrm{sim}\left(\mathbf{x}_i, \mathbf{x}_j\right)$$

- Residual $\rightarrow$ Heterophily and High-Frequency Signals

$$\mathbf{R}^{[1]} = \mathbf{Z}^{[1]} - \mathbf{Z}^{[0]} = \tilde{\mathbf{A}}\mathbf{X}\mathbf{W} - \mathbf{X}\mathbf{W} = -\mathbf{L}\mathbf{X}\mathbf{W}$$

[5] Qiao, Hezhe, and Guansong Pang. "Truncated affinity maximization: One-class homophily modeling for graph anomaly detection." Advances in Neural Information Processing Systems 36 (2023).

# The proposed generalist GAD method - ARC



**Step 3**: Cross-Attentive In-**C**ontext Anomaly Scoring

- **Cross-attention**

Key: labelled normal nodes $\mathbf{H}_k$
Query: unlabelled nodes $\mathbf{H}_q$

$$\mathbf{Q} = \mathbf{H}_q \mathbf{W}_q$$
$$\mathbf{K} = \mathbf{H}_k \mathbf{W}_k$$

$$\tilde{\mathbf{H}}_q = \mathrm{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^{\top}}{\sqrt{d_e}}\right)\mathbf{H}_k$$

Training objective: Reconstruct $\mathbf{H}_q$ with $\mathbf{H}_k$

$$\mathcal{L} = \begin{cases} 1 - \cos\left(\mathbf{H}q_i, \tilde{\mathbf{H}}q_i\right), & \text{if } \mathbf{y}_i = 0 \\ \max\left(0, \cos\left(\mathbf{H}q_i, \tilde{\mathbf{H}}q_i\right) - \epsilon\right), & \text{if } \mathbf{y}_i = 1 \end{cases}$$

# The proposed generalist GAD method - ARC



**Step 3**: Cross-Attentive In-**C**ontext Anomaly Scoring

- **Cross-attention**

Key: labelled normal nodes $\mathbf{H}_k$
Query: unlabelled nodes $\mathbf{H}_q$

$$\mathbf{Q} = \mathbf{H}_q \mathbf{W}_q$$
$$\mathbf{K} = \mathbf{H}_k \mathbf{W}_k$$

$$\tilde{\mathbf{H}}_q = \mathrm{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_e}}\right)\mathbf{H}_k$$

- **Anomaly scoring**

Reconstruction errors as anomaly scores

$$f(v_i) = d(\mathbf{H}q_i, \tilde{\mathbf{H}}q_i) = \sqrt{\sum_{j=1}^{d_e}\left(\mathbf{H}q_{ij} - \tilde{\mathbf{H}}q_{ij}\right)^2}$$

# The proposed generalist GAD method - ARC



**Step 3**: Cross-Attentive In-**C**ontext Anomaly Scoring

**Motivation**:
normal query nodes can be easily reconstructed by the key nodes (other normal nodes)



(a) Case I          (b) Case II

- **Anomaly scoring**

Reconstruction errors as anomaly scores

$$f(v_i) = d(\mathbf{H}q_i, \tilde{\mathbf{H}}q_i) = \sqrt{\sum_{j=1}^{d_e} \left( \mathbf{H}q_{ij} - \tilde{\mathbf{H}}q_{ij} \right)^2}$$

# Experiments: Settings

- 4 groups of datasets
- the largest dataset → training datasets; the rest → testing datasets

| Dataset | Train | Test | #Nodes | #Edges | #Features | Avg. Degree | #Anomaly | %Anomaly |
|---|---|---|---|---|---|---|---|---|
| Citation network with injected anomalies | | | | | | | | |
| Cora | - | ✓ | 2,708 | 5,429 | 1,433 | 3.90 | 150 | 5.53 |
| CiteSeer | - | ✓ | 3,327 | 4,732 | 3,703 | 2.77 | 150 | 4.50 |
| ACM | - | ✓ | 16,484 | 71,980 | 8,337 | 8.73 | 597 | 3.62 |
| PubMed | ✓ | - | 19,717 | 44,338 | 500 | 4.50 | 600 | 3.04 |
| Social network with injected anomalies | | | | | | | | |
| BlogCatalog | - | ✓ | 5,196 | 171,743 | 8,189 | 66.11 | 300 | 5.77 |
| Flickr | ✓ | - | 7,575 | 239,738 | 12,047 | 63.30 | 450 | 5.94 |
| Social network with real anomalies | | | | | | | | |
| Facebook | - | ✓ | 1,081 | 55,104 | 576 | 50.97 | 25 | 2.31 |
| Weibo | - | ✓ | 8,405 | 407,963 | 400 | 48.53 | 868 | 10.30 |
| Reddit | - | ✓ | 10,984 | 168,016 | 64 | 15.30 | 366 | 3.33 |
| Questions | ✓ | - | 48,921 | 153,540 | 301 | 3.13 | 1,460 | 2.98 |
| Co-review network with real anomalies | | | | | | | | |
| Amazon | - | ✓ | 10,244 | 175,608 | 25 | 17.18 | 693 | 6.76 |
| YelpChi | ✓ | - | 23,831 | 49,315 | 32 | 2.07 | 1,217 | 5.10 |

# Experiments: Main Results

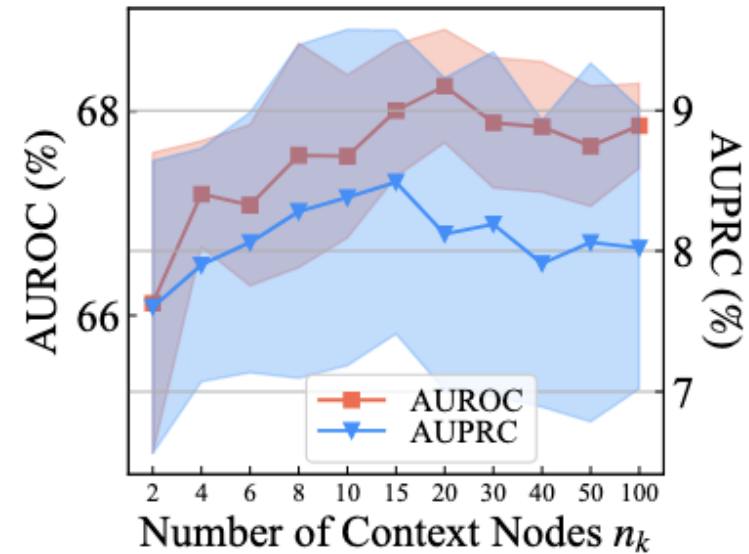| Method | Cora | CiteSeer | ACM | BlogCatalog | Facebook | Weibo | Reddit | Amazon | Rank |
|---|---|---|---|---|---|---|---|---|---|
| Supervised - Pre-Train Only | | | | | | | | | |
| GCN | $59.64_{\pm 8.30}$ | $60.27_{\pm 8.11}$ | $60.49_{\pm 9.65}$ | $56.19_{\pm 6.39}$ | $29.51_{\pm 4.86}$ | $76.64_{\pm 17.69}$ | $50.43_{\pm 4.41}$ | $46.63_{\pm 3.47}$ | 8.9 |
| GAT | $50.06_{\pm 2.65}$ | $51.59_{\pm 3.49}$ | $48.79_{\pm 2.73}$ | $50.40_{\pm 2.80}$ | $51.88_{\pm 2.16}$ | $53.06_{\pm 7.48}$ | $51.78_{\pm 4.04}$ | $50.52_{\pm 17.22}$ | 10.0 |
| BGNN | $42.45_{\pm 11.57}$ | $42.32_{\pm 11.82}$ | $44.00_{\pm 13.69}$ | $47.67_{\pm 8.52}$ | $54.74_{\pm 25.29}$ | $32.75_{\pm 35.35}$ | $50.27_{\pm 3.84}$ | $52.26_{\pm 3.31}$ | 11.1 |
| BWGNN | $54.06_{\pm 3.27}$ | $52.61_{\pm 2.88}$ | $67.59_{\pm 0.70}$ | $56.34_{\pm 1.21}$ | $45.84_{\pm 4.97}$ | $53.38_{\pm 1.61}$ | $48.97_{\pm 5.74}$ | $55.26_{\pm 16.95}$ | 9.0 |
| GHRN | $59.89_{\pm 6.57}$ | $56.04_{\pm 9.19}$ | $55.65_{\pm 6.37}$ | $57.64_{\pm 3.48}$ | $44.81_{\pm 8.06}$ | $51.87_{\pm 14.18}$ | $46.22_{\pm 2.33}$ | $49.48_{\pm 17.13}$ | 9.8 |
| Unsupervised - Pre-Train Only | | | | | | | | | |
| DOMINANT | $66.53_{\pm 1.15}$ | $69.47_{\pm 2.02}$ | $70.08_{\pm 2.34}$ | $74.25_{\pm 0.65}$ | $51.01_{\pm 0.78}$ | $92.88_{\pm 0.32}$ | $50.05_{\pm 4.92}$ | $48.94_{\pm 2.69}$ | 5.8 |
| CoLA | $63.29_{\pm 8.88}$ | $62.84_{\pm 9.52}$ | $66.85_{\pm 4.43}$ | $50.04_{\pm 3.25}$ | $12.99_{\pm 11.68}$ | $16.27_{\pm 5.64}$ | $52.81_{\pm 6.69}$ | $47.40_{\pm 7.97}$ | 9.5 |
| HCM-A | $54.28_{\pm 4.73}$ | $48.12_{\pm 6.80}$ | $53.70_{\pm 4.64}$ | $55.31_{\pm 0.57}$ | $35.44_{\pm 13.97}$ | $65.52_{\pm 12.58}$ | $48.79_{\pm 2.75}$ | $43.99_{\pm 0.72}$ | 11.4 |
| TAM | $62.02_{\pm 2.39}$ | $72.27_{\pm 0.83}$ | $74.43_{\pm 1.59}$ | $49.86_{\pm 0.73}$ | $65.88_{\pm 6.66}$ | $71.54_{\pm 0.18}$ | $55.43_{\pm 0.33}$ | $56.06_{\pm 2.19}$ | 5.6 |
| Unsupervised - Pre-Train & Fine-Tune | | | | | | | | | |
| DOMINANT | $72.23_{\pm 0.34}$ | $74.69_{\pm 0.32}$ | $74.34_{\pm 0.12}$ | $74.61_{\pm 0.04}$ | $49.92_{\pm 0.55}$ | $92.21_{\pm 0.10}$ | $52.14_{\pm 5.06}$ | $59.06_{\pm 2.80}$ | 3.6 |
| CoLA | $67.62_{\pm 4.26}$ | $70.75_{\pm 3.42}$ | $69.11_{\pm 0.67}$ | $62.49_{\pm 3.38}$ | $64.70_{\pm 18.86}$ | $31.55_{\pm 6.02}$ | $58.12_{\pm 0.67}$ | $52.51_{\pm 6.66}$ | 5.4 |
| HCM-A | $56.45_{\pm 4.93}$ | $55.54_{\pm 4.07}$ | $57.69_{\pm 3.59}$ | $55.10_{\pm 0.29}$ | $36.57_{\pm 10.72}$ | $71.89_{\pm 2.79}$ | $49.15_{\pm 2.72}$ | $42.20_{\pm 0.55}$ | 10.1 |
| TAM | $62.56_{\pm 2.10}$ | $76.54_{\pm 1.33}$ | $86.29_{\pm 1.57}$ | $57.69_{\pm 0.88}$ | $76.26_{\pm 3.70}$ | $71.73_{\pm 0.16}$ | $56.62_{\pm 0.49}$ | $57.13_{\pm 1.59}$ | 3.4 |
| Ours | | | | | | | | | |
| ARC | $87.45_{\pm 0.74}$ | $90.95_{\pm 0.59}$ | $79.88_{\pm 0.28}$ | $74.76_{\pm 0.06}$ | $67.56_{\pm 1.60}$ | $88.85_{\pm 0.14}$ | $60.04_{\pm 0.69}$ | $80.67_{\pm 1.81}$ | 1.5 |

🎯 Strong detection capability without fine-tuning

🌍 Generalizability in different datasets/domains

# Experiments: Sensitivity In Terms of #shots



(a) Cora

(b) Facebook

✦ Works well with extremely few shots
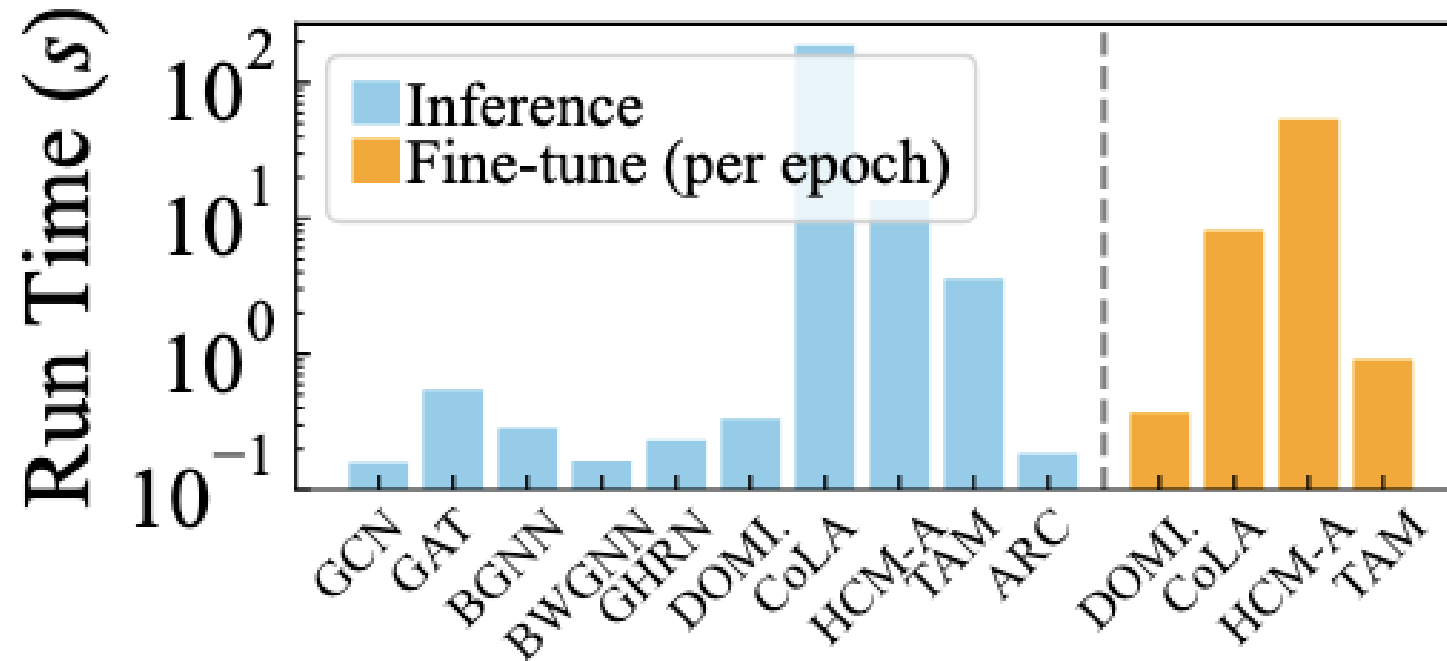
⊞ More labelled normal samples bring better performance

# Experiments: Ablation Study

| Variant | Cora | CiteSeer | ACM | BlogCatalog | Facebook | Weibo | Reddit | Amazon |
|---|---|---|---|---|---|---|---|---|
| ARC w/o A | $80.65_{\pm0.71}$ | $83.35_{\pm0.64}$ | $79.29_{\pm0.16}$ | $73.86_{\pm0.18}$ | $62.80_{\pm2.06}$ | $89.69_{\pm0.17}$ | $54.60_{\pm1.92}$ | $64.76_{\pm2.13}$ |
| ARC w/o R | $37.44_{\pm1.40}$ | $31.52_{\pm0.71}$ | $61.83_{\pm1.16}$ | $49.30_{\pm2.06}$ | $20.38_{\pm9.63}$ | $97.72_{\pm0.59}$ | $52.94_{\pm0.96}$ | $50.15_{\pm0.24}$ |
| ARC w/o C | $47.39_{\pm0.42}$ | $53.98_{\pm0.72}$ | $54.24_{\pm1.32}$ | $60.46_{\pm1.23}$ | $48.86_{\pm0.97}$ | $42.84_{\pm3.01}$ | $51.03_{\pm0.86}$ | $69.02_{\pm0.97}$ |
| ARC | $87.45_{\pm0.74}$ | $90.95_{\pm0.59}$ | $79.88_{\pm0.28}$ | $74.76_{\pm0.06}$ | $67.56_{\pm1.60}$ | $88.85_{\pm0.14}$ | $60.04_{\pm0.69}$ | $80.67_{\pm1.81}$ |

♻ Each component has a significant contribution to the final performance

# Experiments: Efficiency Analysis



🕐 High inference efficiency – comparable to GCN

# Experiments: Visualization



(a) Cora  (b) Facebook

(a) Case I  (b) Case II

🗣 Interpretability – attention score

Case 1: uniform attention weights
→ "Single-class normal ": Reconstructed embeddings that closely to the average embedding of the context nodes

Case 2: two fixed patterns for normal queries
→ "Multi-class normal": Two cluster centers

# Summary

**New paradigm:** generalist GAD: one model for all datasets!

**Effective solution:** ARC – a simple yet effective methods

**Extensive experiments:** ARC enjoys superior performance, great generalizability, high running efficiency, and potential explainability

**Full paper**

**GitHub**

# Thanks for your listening!