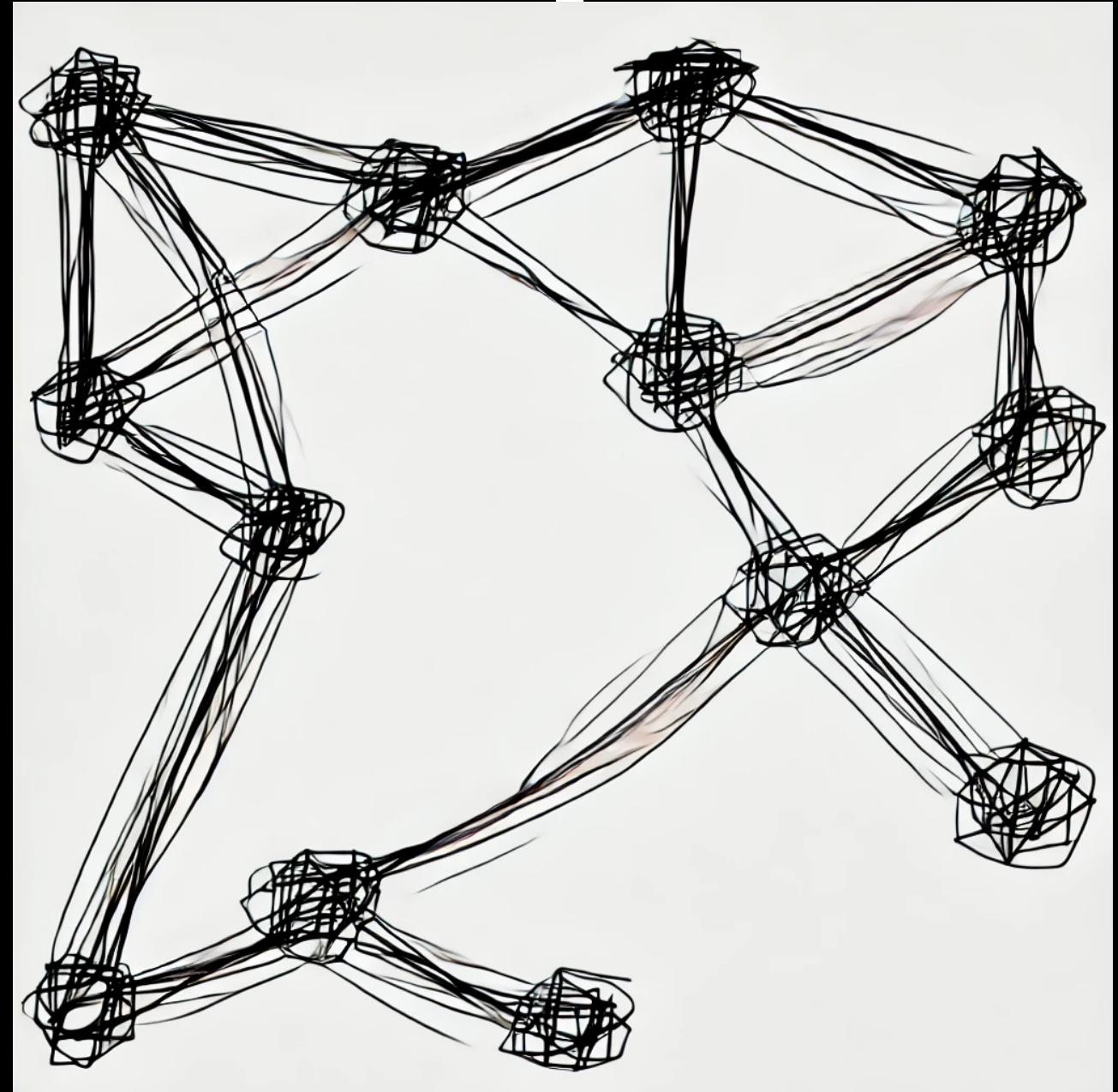


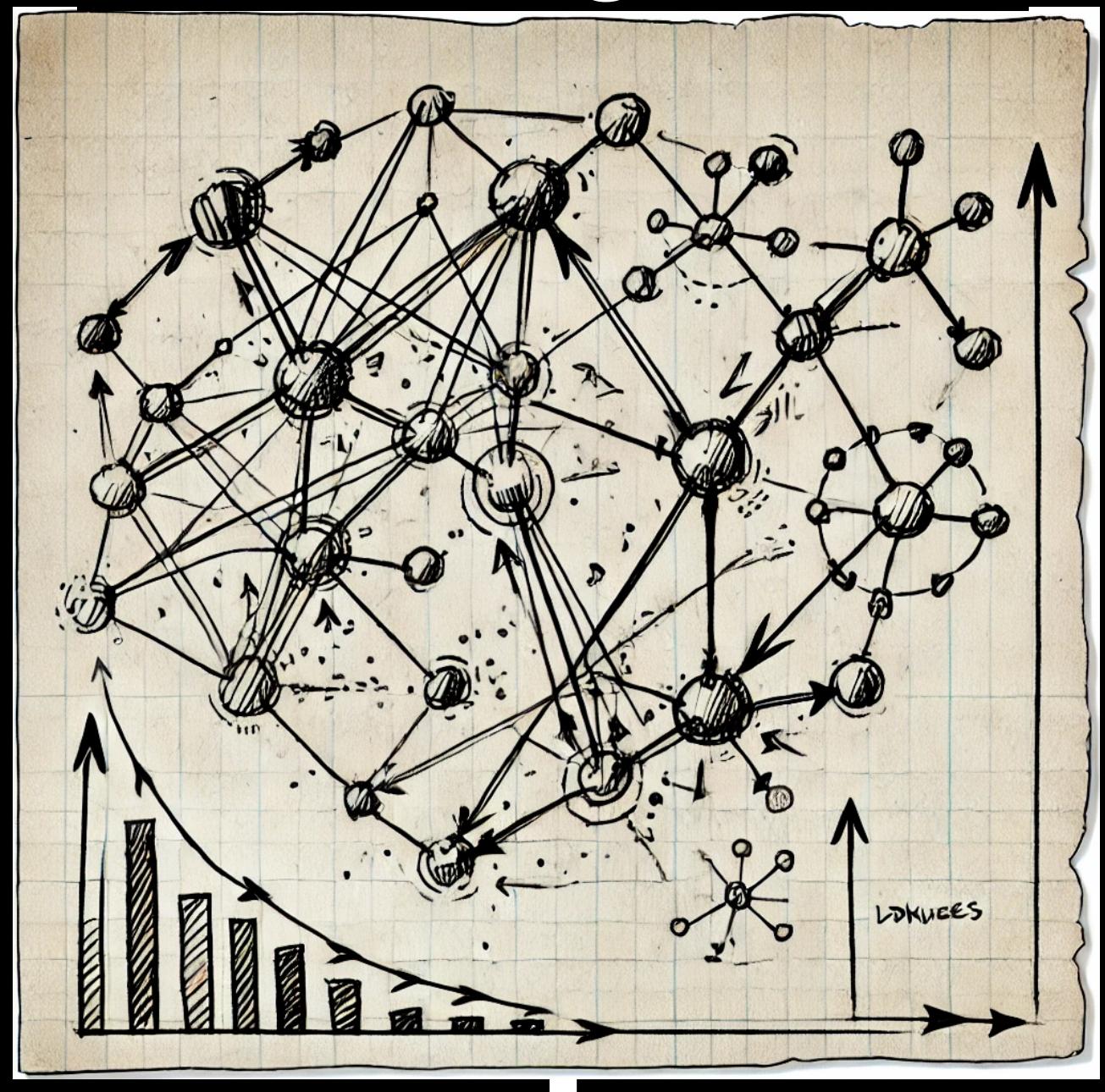
# Have We Gone Too Far?

(Pics generated by ChatGPT)

## Notion of Graph

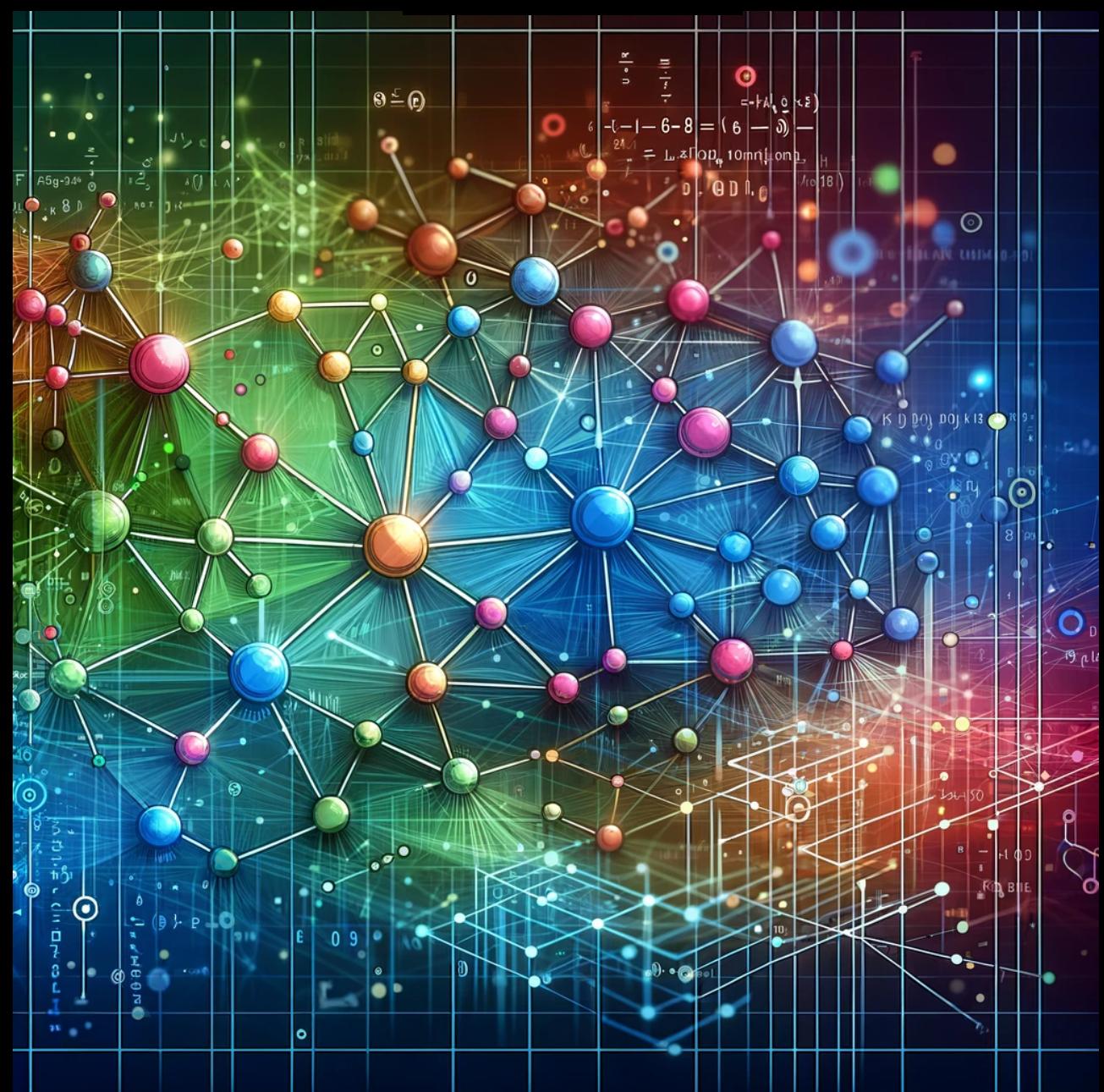


## Classic Algorithms



## (What Happens?)

### GNNs



## Foundation Models



# How Graph Neural Networks Learn: Lessons from Training Dynamics

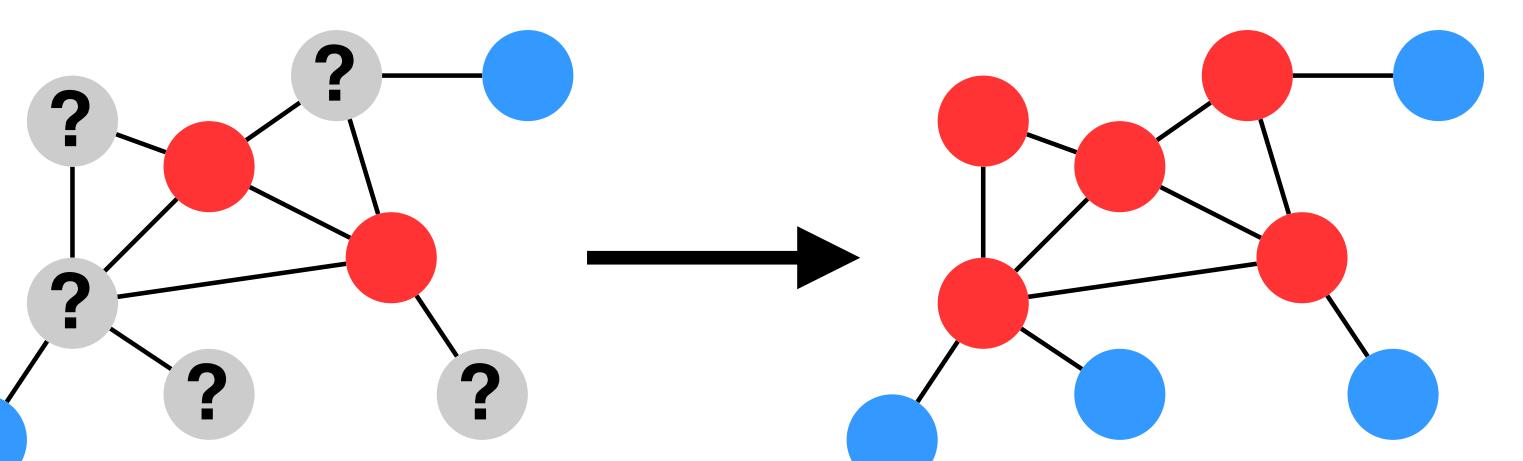
Chenxiao Yang<sup>1</sup>, Qitian Wu<sup>1</sup>, David Wipf<sup>2</sup>, Ruoyu Sun<sup>3</sup>, Junchi Yan<sup>1</sup>

<sup>1</sup> Shanghai Jiao Tong University <sup>2</sup> Amazon Web Services <sup>3</sup> The Chinese University of Hong Kong



## Embarrassingly Simple Algorithm > GNNs

### • Setup



Adjacency  $\mathbf{A} \in \mathbb{R}^{n \times n}$  ( $n = \# \text{ nodes}$ ), labels  $\mathbf{Y} \in \mathbb{R}^{n_l}$  ( $n_l = \# \text{ labeled nodes}$ ), predictions for labeled/unlabeled nodes  $[\mathbf{F}, \mathbf{F}'] \in \mathbb{R}^n$ .

### • Residual Propagation Algorithm

Forward message passing is all this algorithm does:

$$[\mathbf{R}_{t+1}, \mathbf{R}'_{t+1}] = -\eta \mathbf{A}^K [\mathbf{R}_t, \mathbf{0}] + [\mathbf{R}_t, \mathbf{R}'_t]$$

Residuals, a.k.a. errors, defined as  $\mathbf{R} = \mathbf{F} - \mathbf{Y} \in \mathbb{R}^{n_l}$  (for training set),  $\mathbf{R}' = \mathbf{F}' - \mathbf{0} \in \mathbb{R}^{n-n_l}$  (for testing set). At initialization,  $\mathbf{R}_0 = -\mathbf{Y}$ ,  $\mathbf{R}'_0 = \mathbf{0}$ .

#### 1. Interpretation: the prediction for an unlabeled node $x'$

$$(Ground-truth) label propagation$$

$$f_{t+1}(x') = f_t(x') + \eta \sum_{x \in X} \mathbf{A}^k(x, x') (y(x) - f_t(x))$$

#### 2. Extensions: $\mathbf{A}^k$ could be replaced by other matrices (not necessarily PSD or symmetric) and incorporate node feature, e.g. $\mathbf{A}^2 \mathbf{K}(\mathbf{X}, \mathbf{X}) \mathbf{A}^5$ .

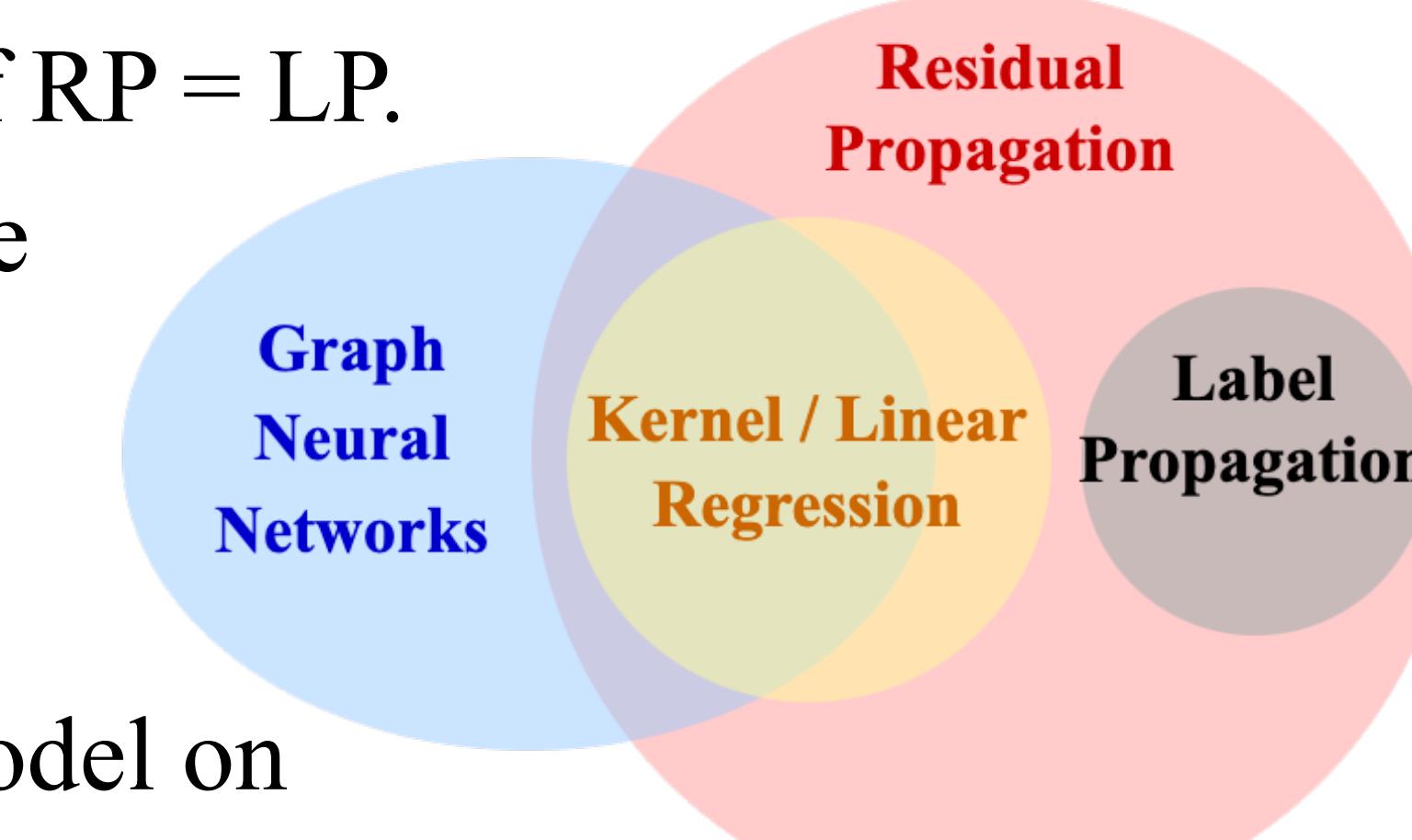
### • Connection with Classic Algorithms & GNNs

#### 1. Label Propagation: first step of RP = LP.

#### 2. Kernel Regression: RP relaxes the constraints on kernel matrices; RP does not need to converge; complexity $\mathcal{O}(n^3) \rightarrow \mathcal{O}(e)$ .

#### 3. Network Embedding: linear model on spectral decomposition of $\mathbf{A}$ (which is computationally very expensive) induces RP.

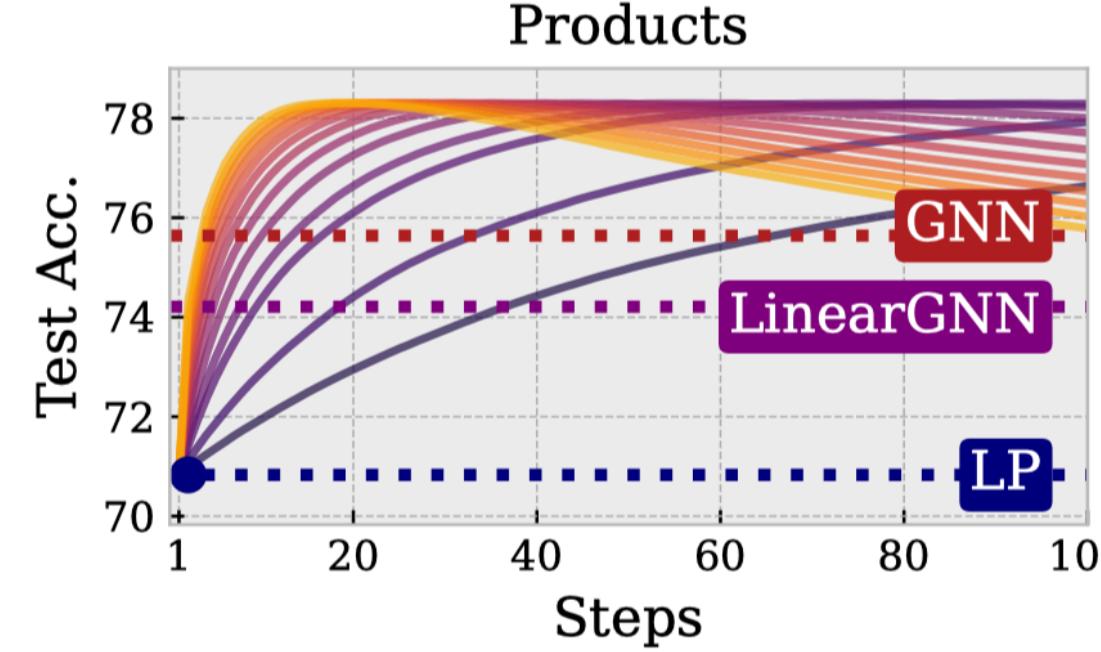
#### 4. GNN: discussed later ....



### • Partial Result

### • Implications

- 1. Performance of classic algorithms can be boosted to the same level as GNNs.
- 2. Explanation of generalization can be potentially decoupled with the model.



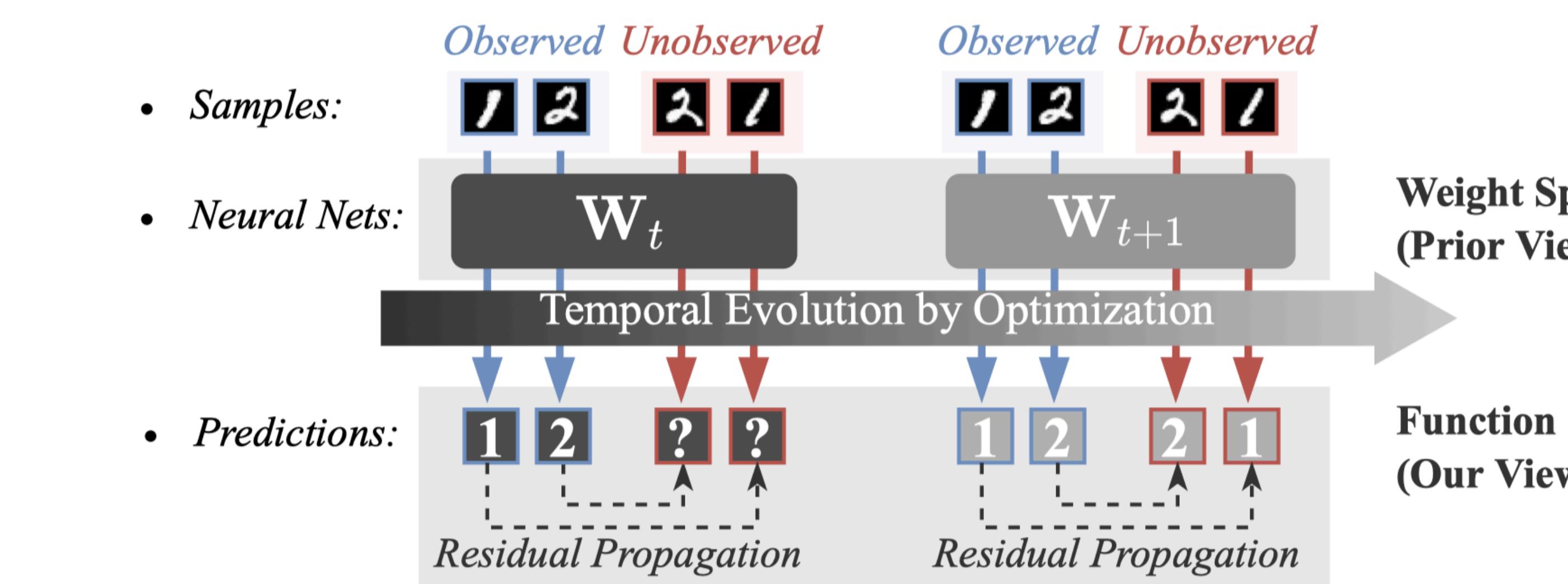
## Optimization and Generalization of GNNs

Q: What is the central commonality of these effective methods?  
→ The particular way the graph structure is leveraged in training matters!

### • Graph Implicit Bias in GNNs' Optimization

$$\text{GD Dynamics: } \frac{\partial [\mathbf{R}_t, \mathbf{R}'_t]}{\partial t} = -\eta \Theta_t(\bar{\mathbf{X}}, \bar{\mathbf{X}}; \mathbf{A}) [\mathbf{R}_t, \mathbf{0}]$$

$$\text{Def (Node-Level GNTK)} \quad \Theta_t^{(\ell)}(\mathbf{x}, \mathbf{x}'; \mathbf{A}) = \nabla_{\mathbf{W}} f(\mathbf{x}; \mathbf{A})^\top \nabla_{\mathbf{W}} f(\mathbf{x}'; \mathbf{A})$$



### Role of Graph Structure in Optimization

#### (Transformation)

$$\bar{\Theta}^{(\ell)} = \Theta^{(\ell-1)} \odot \dot{\Sigma}^{(\ell)} + \bar{\Sigma}^{(\ell)}$$

$$\begin{cases} \Sigma^{(\ell)} = \mathbf{A} \bar{\Sigma}^{(\ell)} \mathbf{A} \\ \Theta^{(\ell)} = \mathbf{A} \bar{\Theta}^{(\ell)} \mathbf{A}. \end{cases}$$

and illustrative examples showing cases where the GNTK is equivalent to some special forms of the adjacency.

### • Generalization Analysis

#### Def (Alignment & Optimal Kernel)

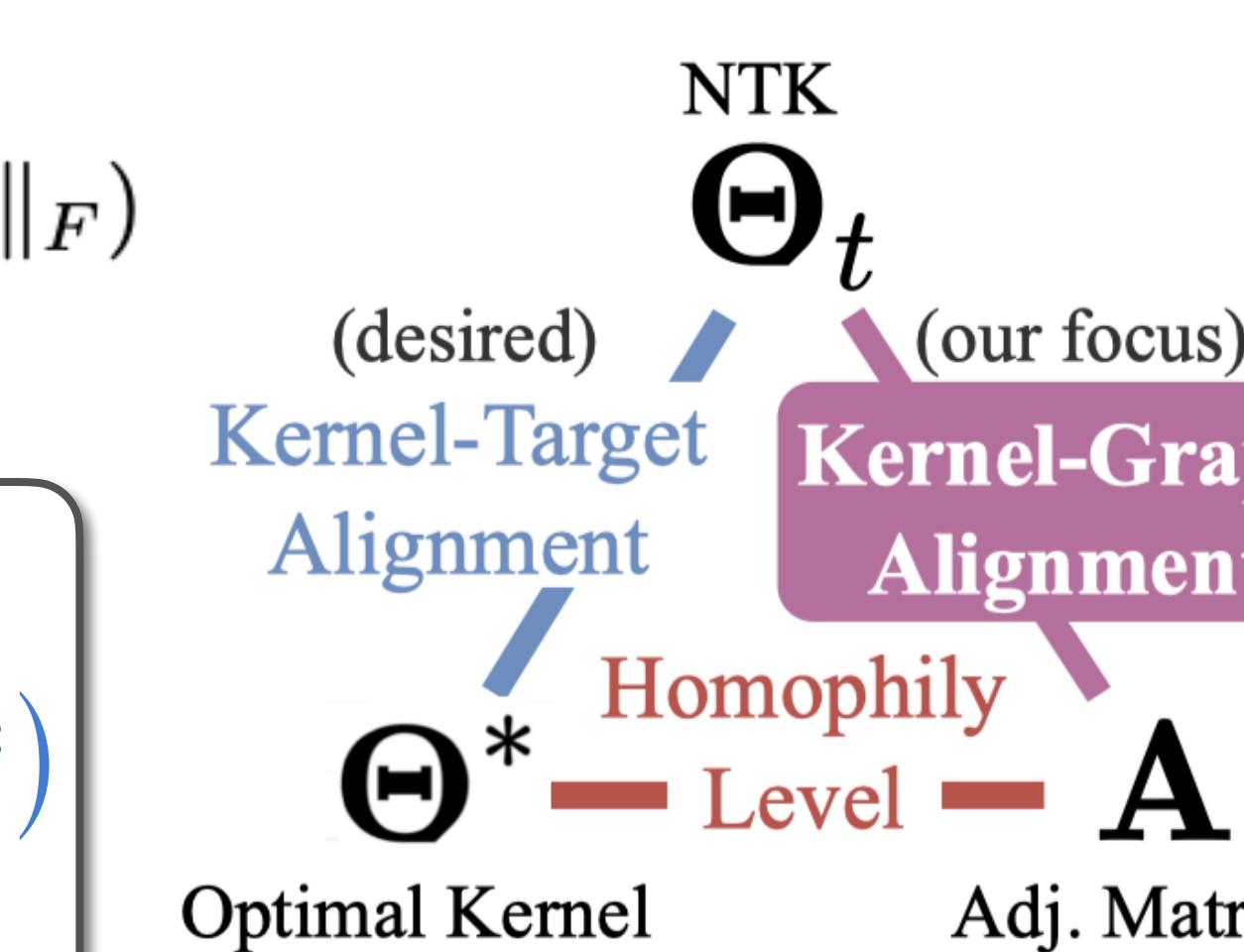
$$A(\mathbf{K}_1, \mathbf{K}_2) \triangleq \langle \mathbf{K}_1, \mathbf{K}_2 \rangle_F / (\|\mathbf{K}_1\|_F \|\mathbf{K}_2\|_F)$$

$$\Theta^* \triangleq \bar{\mathbf{Y}} \bar{\mathbf{Y}}^\top$$

#### Homophily Level: $A(\mathbf{A}, \Theta^*)$

#### Kernel-Target Alignment: $A(\Theta_t, \Theta^*)$

#### Kernel-Graph Alignment: $A(\Theta_t, \mathbf{A})$



Result 1 (following standard analysis)

$$\mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [l(f(\mathbf{x}), y)] = \mathcal{O} \left( \sqrt{1 - cn_l^{-1} A(\mathbf{A}, \Theta^*)} + \sqrt{n_l^{-1} \log(\delta^{-1})} \right)$$

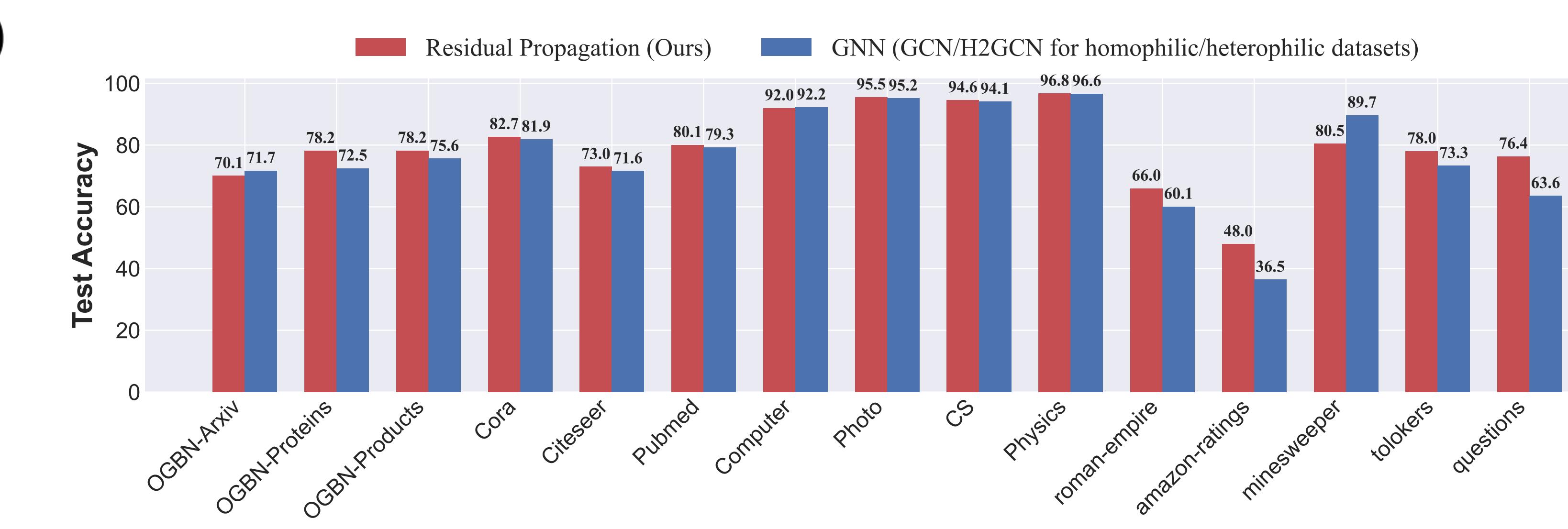
Result 2 (removing the iid assumption)

When the data generating distribution satisfies certain homophily properties, GNN is the Bayesian optimal model that minimizes the population risk.

## Experimental Verification

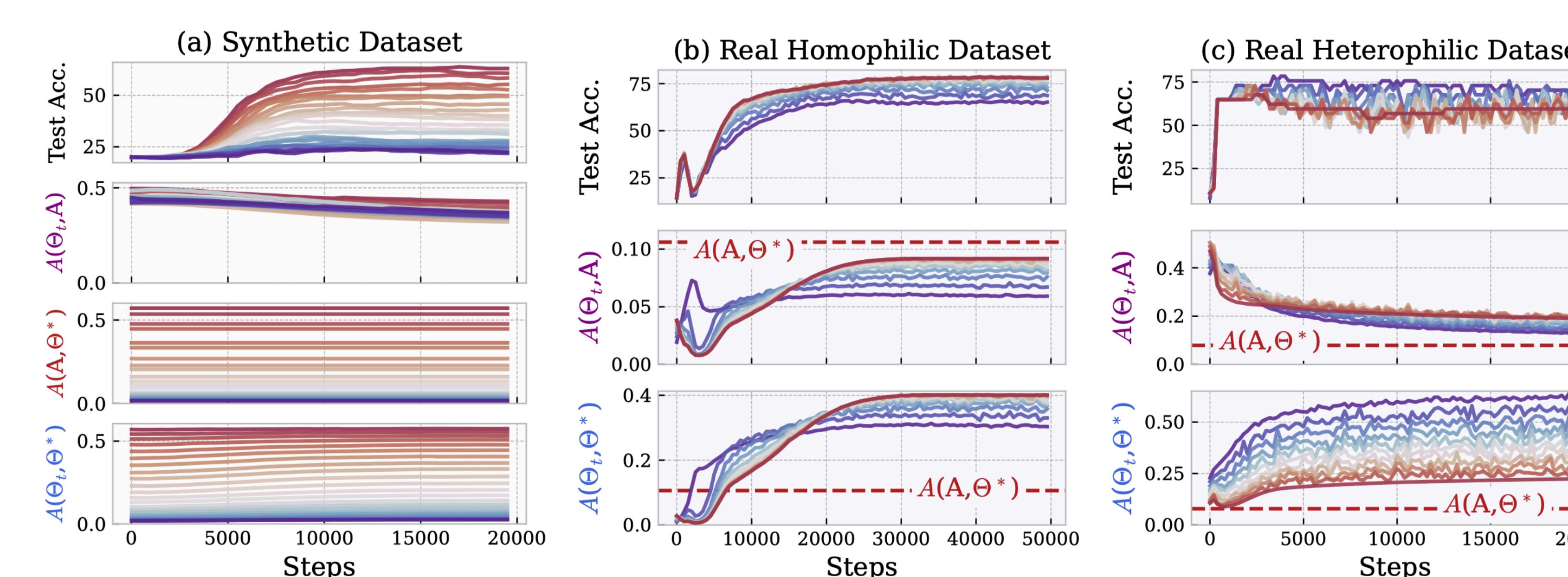
### • Evaluation of Residual Propagation

Model	Feat.	Arxiv		Proteins		Products		# Param.
		Validation	Test	Validation	Test	Validation	Test	
MLP	X	57.65 ± 0.12	55.50 ± 0.23	77.06 ± 0.14	72.04 ± 0.48	75.54 ± 0.14	71.06 ± 0.08	$\mathcal{O}(dn^2)$
LinearGNN	X, A	70.67 ± 0.02	69.39 ± 0.11	66.11 ± 0.87	62.89 ± 0.11	88.97 ± 0.01	74.21 ± 0.04	$\mathcal{O}(dc)$
GNN	X, A	73.00 ± 0.17	71.74 ± 0.29	79.21 ± 0.18	72.51 ± 0.35	92.00 ± 0.03	75.64 ± 0.21	$\mathcal{O}(dn^2)$
LP	A	70.14 ± 0.00	68.32 ± 0.00	83.02 ± 0.00	74.73 ± 0.00	90.91 ± 0.00	74.34 ± 0.00	0
RP (ours)	A	71.37 ± 0.00	70.06 ± 0.00	85.19 ± 0.00	78.17 ± 0.00	91.31 ± 0.00	78.25 ± 0.00	0
Speedup / step		$\times 14.48$		$\times 14.00$		$\times 12.46$		
Time to Acc.		$\times 0.01461$		$\times 0.00008$		$\times 0.00427$		
Memory		$\times 0.094$		$\times 0.363$		$\times 0.151$		



1. Simple non-parametric algorithm rivals with GNNs on 15 benchmarks (including homophilic and heterophilic ones).
2. RP is significantly more efficient and scalable (up to 10,000× speedup, 10× less memory, easily run on million-scale dataset).
3. RP is especially effective on large datasets (where it secretly has more parameters than most GNNs).

### • Training Dynamics: How GNNs Handle Heterophily?



1. GNN's NTK inherently aligns with the graph structure (i.e. kernel-graph alignment or KGA).
2. On homophilic graphs, better KGA implies better generalization.
3. GNN's NTK evolves during training and also gradually aligns with the optimal kernel (i.e. the ideal graph), indicating how GNN adaptively handles heterophily.