

# FIXED AGGREGATION FEATURES CAN RIVAL GNNs

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Graph neural networks (GNNs) are widely believed to excel at node representation learning through trainable neighborhood aggregations. We challenge this view by introducing Fixed Aggregation Features (FAFs), a training-free approach that transforms graph learning tasks into tabular problems. This simple shift enables the use of well-established tabular methods, offering strong interpretability and the flexibility to deploy diverse classifiers. Across 14 benchmarks, well-tuned multilayer perceptrons trained on FAFs rival or outperform state-of-the-art GNNs and graph transformers on 12 tasks – often using only mean aggregation. The only exceptions are the Roman Empire and Minesweeper datasets, which typically require unusually deep GNNs. To explain the theoretical possibility of non-trainable aggregations, we connect our findings to Kolmogorov–Arnold representations and discuss when mean aggregation can be sufficient. In conclusion, our results call for (i) richer benchmarks benefiting from learning diverse neighborhood aggregations, (ii) strong tabular baselines as standard, and (iii) employing and advancing tabular models for graph data to gain new insights into related tasks.

## 1 INTRODUCTION

Graph neural networks (GNNs) have become the standard approach for learning from graph based data, and in particular, for solving node classification. Most models follow the message-passing paradigm (Gilmer et al., 2017), where each node updates its representation by alternating neighborhood aggregation with learned linear combinations across multiple hops. This framework has been remarkably successful at combining node features with graph structure, driving applications in domains ranging from social networks to biology (Bongini et al., 2023; Sharma et al., 2024). Yet, it comes at the cost of high model complexity that poses challenges for interpretation. We ask the question whether this high complexity is really necessary.

Recent evidence (Luo et al., 2024; 2025a) shows that classic models, such as GCN (Kipf & Welling, 2017), GATv2 (Brody et al., 2022), and GraphSAGE (Hamilton et al., 2017), remain surprisingly competitive when equipped with proper hyperparameter tuning and standard optimization techniques. When carefully tuned, they can rival more sophisticated approaches, including state-of-the-art Graph Transformers (Wu et al., 2023; Deng et al., 2024; Kong et al., 2023; Wu et al., 2022; Chen et al., 2023; Rampásek et al., 2022; Shirzad et al., 2023) and models designed for heterophily (Zhu et al., 2020; 2021; Chien et al., 2021; Maurya et al., 2022; Li et al., 2022).

These results invite a closer look at which components of graph learning architectures are essential for strong performance, and thus raise a natural next question: *How relevant is learning the aggregation?* In fact, the field has invested heavily in learning increasingly complex convolution layers and attention mechanisms. In this paper we challenge that premise from first principles. Leveraging the Kolmogorov–Arnold representation theorem (Kolmogorov, 1957; Schmidt-Hieber, 2021), we give an explicit, lossless construction of neighborhood aggregations. Consequently, one can in theory encode neighbor features without discarding information. However, the same construction exposes a crucial gap between expressiveness and learnability: these lossless encoders are numerically brittle (e.g., sensitive to floating-point noise) and tend to produce “rough” embeddings that are ill-suited for standard classifiers on Euclidean space such as MLPs.

Surprisingly, we find that, standard, untrained aggregation operators—sum, mean, max, min—though not information-preserving, yield useful features without any learnable parameters. Building on this observation, we propose Fixed Aggregation Features (FAF) (§ 3): a training-free

aggregation pipeline that applies fixed **aggregation functions**—also referred to as “**reducers**”—over neighborhoods at multiple hops, concatenates the results into a tabular feature matrix, and then trains only a downstream classifier (e.g., an MLP). This data transformation brings several advantages: high interpretability (feature importance and ablations over hops/**reducers**), compatibility with the rich toolbox of tabular learning (designed to handle noise, class imbalance, feature selection, etc.), architectural flexibility, and reduced training compute.

Empirically, FAFs combined with well-tuned MLP classifiers are competitive on 12/14 common node-classification benchmarks, including citation networks (McCallum et al., 2000; Sen et al., 2008; Namata et al., 2012), coauthor and Amazon co-purchase graphs (Shchur et al., 2018), Wikipedia (Mernyei & Cangea, 2020; Rozemberczki et al., 2021), and other heterophilous datasets (Platonov et al., 2023b). Performance truly trails only on Minesweeper and Roman-Empire, where the best GNNs rely on linear residual connections; in fact, the remaining gap aligns with the gains from residuals reported by Luo et al. (2024). This pattern suggests that these datasets (Minesweeper and Roman-Empire) benefit from hop-specific aggregations. While these GNNs profit from many layers (10–15), the best-performing FAFs use only 2–4 hops. But why do FAFs work so well in the other cases? Our theoretical analysis of the employed aggregation functions (§ 4) and our empirical findings (§ 5) suggest that, for most benchmarks, the relevant signal is concentrated within hops 0–2; on hops 0–1, sum and mean preserve information. At higher hops, different **reducers** are complementary, but the information gain from min/max diminishes.

FAFs also let us examine datasets from an optimization-first viewpoint without hard-to-interpret architectural factors (§ 3.1). By converting neighborhoods into tabular features, we decouple representation from optimization and enable standard interpretability tools (e.g., feature importance (Lundberg & Lee, 2017)) to identify which hop distances and **reducers** carry signal. Beyond revisiting the homophily–heterophily dichotomy or one-hop informativeness (Platonov et al., 2023a; Zheng et al., 2024; 2025), our method supports a richer characterization of interaction patterns—how signal varies across scales, which effects are additive vs. redundant, and where long-range dependencies matter. The tabular view also makes it natural to augment features with network-science descriptors (Blöcker et al., 2025) and neighborhood-masking features inspired by graph rewiring (Rubio-Madrigal et al., 2025) or computational-graph splitting (Roth et al., 2025).

Together, our results suggest that many benchmarks do not require sophisticated learned aggregations, and that a large portion of GNN performance can be matched by powerful tabular baselines built from fixed, transparent summaries. FAFs can serve both as a strong baseline and as a diagnostic tool for graph benchmarks in this setting. In summary, here are our main contributions:

1. **Theory:** We construct lossless neighborhood aggregations via Kolmogorov–Arnold representations, clarifying that learnability and numeric stability—not just expressiveness—govern practical success. Moreover, we analyze what information common **reducers** extract from neighborhoods, revealing information preservation at 0 and 1 hops and diminishing information with higher depth for min/max.
2. **Method:** We introduce FAFs, which convert graph data into a tabular task by stacking fixed multi-hop aggregations, offering an interpretable framework to study the interplay between graph structure, features, and the task.
3. **Empirical evidence:** FAFs match or exceed classic GNNs on 12/14 standard benchmarks. Our experiments further corroborate our analysis, finding low hop features to be more important and diminishing information at higher depth.
4. **Implications:** Our findings question the necessity of learned neighborhood aggregation on current standard benchmarks, motivate strong tabular baselines for graph data, and open a path to more interpretable, efficient graph learning—and to designing harder benchmarks that genuinely benefit from aggregation learning.

## 2 RELATED WORK

**Simplifying GNNs.** A growing body of work shows that much of a GNN’s power can be retained—even improved for some tasks—when message passing is simplified or fixed. Early evidence comes from Kipf & Welling (2017), inspiring lines of work where aggregation layers are

frozen or randomized: Ramachandra et al. (2025) obtain competitive node classification on relational graphs by aggregating randomly transformed random features; Kelesis et al. (2025) analyze partially frozen GCNs, showing that fixing aggregation can mitigate over-smoothing and ease optimization; and GESN (Gallicchio & Micheli, 2010) compute node embeddings via a dynamical system with randomly initialized reservoir weights, after which only a linear readout is trained for node classification (Micheli & Tortorella, 2024). Another simplification comes from SGC (Wu et al., 2019), which remove nonlinearities from a GCN, yielding a model equivalent to applying a low-pass graph filter followed by a linear readout. For link prediction (Qarkaxhija et al., 2024), SGCNs are found to be better than GCNs, but even removing the classifier linear layer—and thus all trainable parameters—provides a good baseline. A fixed aggregation scheme is also used in APPNP (Gasteiger et al., 2019): an MLP is first trained to produce node embeddings, and a subsequent Personalized PageRank-based propagation is applied; although the propagation itself is fixed, it remains in the computation graph, so gradients flow during backpropagation. And for graph classification on non-attributed graphs, Cai & Wang (2019) show that first-neighborhood statistics with an SVM form a surprisingly strong baseline, but performance lags on attributed graphs. In contrast, we focus on node classification with rich node features. We aggregate across all hops and concatenate these as inputs, and we place a more powerful, and more carefully tuned classifier on top (an MLP), all of which are necessary for our results—as shown in our empirical results, such as in Tables 7 and 8. With it we highlight the value of concatenating dependent but informative hop-wise features (Reddy et al., 2025), and the benefits of overparameterization on graphs—echoing evidence from random-feature models (Donghi et al., 2024).

**Benchmarking GNNs.** Our work also connects to the growing literature on properly benchmarking GNNs and what constitutes a meaningful graph-learning dataset. For graph classification, Errica et al. (2020) show that, under controlled protocols, simple and even structure-agnostic baselines can rival complex GNNs, suggesting that common benchmarks often fail to exploit graph structure. Recent analyses likewise warn that graph learning risks losing relevance without application-grounded benchmarks (Bechler-Speicher et al., 2025) and principled criteria for dataset quality beyond accuracy (Coupette et al., 2025). On the dataset side, Bazhenov et al. (2025) recently introduced industrial node property prediction benchmarks with graph-agnostic baselines. Their neighborhood feature aggregation (NFA), which augments tabular models with one-hop aggregated neighbor statistics, can be seen as a one-hop instance of our FAF construction. Concurrently, Anonymous (2025) propose a suite of benchmarks to standardize graph-learning evaluation. In this context, our FAF approach serves as a simple stress test of whether proposed graph benchmarks genuinely benefit from learned message passing, and we argue that such well-tuned, fixed, multi-hop baselines should be routinely included when assessing new graph models and datasets.

**GNN aggregation functions.** Classical message-passing GNNs differ primarily in how they aggregate neighbor features under permutation invariance. Sum, mean, and max are the canonical choices, with injectivity and stability trade-offs of each of them tied to their multiset representations (Xu et al., 2019). Beyond single operators, principal neighborhood aggregation (PNA) mixes several base **reducers** with degree-aware scalars to boost expressiveness and well-conditioning for continuous **reducers** (Corso et al., 2020). Attention mechanisms instantiate learned weighted sums (Veličković et al., 2018; Brody et al., 2022), although it has been shown that they suffer from trainability problems, including small relative gradients on the attention parameters, slowed-down layer-wise training speed, and the inability to mute neighbors (Mustafa et al., 2023; Mustafa & Burkholz, 2024a;b). Our perspective is complementary: We study fixed **reducers** whose strength comes from (i) their information preservation and (ii) their separability by a powerful downstream classifier. This decoupling clarifies what must be learned (the readout) versus what can be fixed (the propagation), and it aligns with our empirical finding that stronger, well-tuned classifiers capitalize on rich, concatenated neighborhood views. Because of this, we argue that good optimization and learnability is as important as expressivity results. In line with this argument, Gorishniy et al. (2022) argue for tabular data that having the right embeddings for continuous features is key to closing the gap between transformer-like architectures and feed forward networks, proposing a lossless piecewise linear embedding to improve the trainability of the latter.

**Kolmogorov-Arnold theorem.** The Kolmogorov–Arnold representation (Kolmogorov, 1957) admits several equivalent formulations that reduce multivariate functions to compositions of univariate

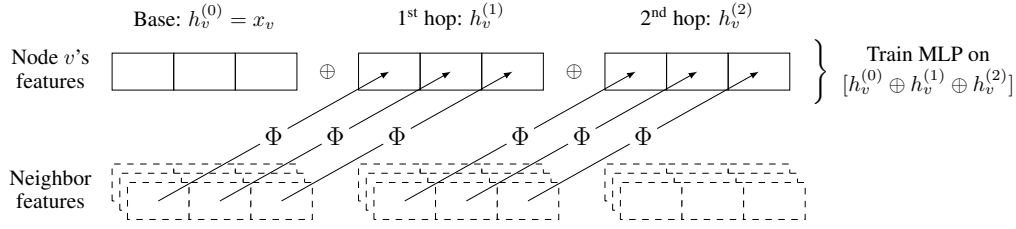


Figure 1: Fixed Aggregation Features (FAFs) are calculated as a pre-processing step, concatenated to the input ( $\oplus$ ), and fed to an MLP. If the aggregation function  $\Phi$  is injective, the neighborhood information is preserved. The Kolmogorov-Arnold representation theorem ensures the existence of such a function, although simple reducers are empirically more amenable for optimization.

functions. Recent architectures explicitly instantiate such decompositions with learnable spline-based univariate components and linear mixing (Liu et al., 2025; Carlo et al., 2025). In contrast, we use a specific fixed-aggregation formulation with predetermined aggregation weights and a fixed univariate encoding (Schmidt-Hieber, 2021), so that any multivariate function  $f$  can be learned exclusively from a single univariate readout  $g$  applied to a fixed weighted sum of univariate transforms.

### 3 FIXED AGGREGATION FEATURES

Node neighborhoods can be compressed into single-node features, eliminating the need to learn feature embeddings before every layer of message passing. Our approach, Fixed Aggregation Features (FAFs), recursively constructs and concatenate features via a set of reducers  $\mathcal{R} \subset \{\text{mean, sum, max, min, std, } \dots\}$  in the following way:

$$h_v^{(0,r)} = x_v, \quad h_v^{(k,r)} = r\left(\{h_u^{(k-1,r)} : u \in N(v)\}\right), \quad (1)$$

where  $k \in \{1, \dots, K\}$  and  $r \in \mathcal{R}$ . We then train an MLP on the concatenated representation

$$z_v = x_v \oplus \left( \bigoplus_{r \in \mathcal{R}} \bigoplus_{k \in \{1, \dots, K\}} \left( h_v^{(k,r)} \right) \right) \quad (2)$$

with input dimensionality  $|x_v| \cdot (1 + |\mathcal{R}| \cdot K)$  per node  $v$ . Figure 1 illustrates the case  $\mathcal{R} = \{\Phi\}$  with  $K = 2$ . If the reducers are injective, then the neighborhood information at each depth is preserved in  $z_v$ . This waives the need for aggregating learned embeddings in GNNs, thus transforming graph data into high-dimensional tabular data. Our analysis explains why this is theoretically possible (§ 4). Additionally, in our experiments (§ 5), we show that MLPs trained on FAFs can match the performance of classic GNNs on most standard node-classification benchmarks and, by comparing with Luo et al. (2024), of Graph Transformers and heterophily-aware models.

#### 3.1 ADVANTAGES OF TABULAR OVER GRAPH DATA REPRESENTATION

We now turn to further benefits of the tabular view: interpretability, optimization, efficiency, and augmentations.

**Interpretability.** Our construction concatenates each node’s original features with  $K$ -hop neighborhood statistics and feeds this expanded representation to a tabular classifier. This setup surfaces feature and hop aggregation factors explicitly, enabling us to assess their contributions using the widely-used toolbox for tabular interpretability. For instance, we can analyze effects across hops by examining feature importance of the MLP. As an illustration, we compute Shapley

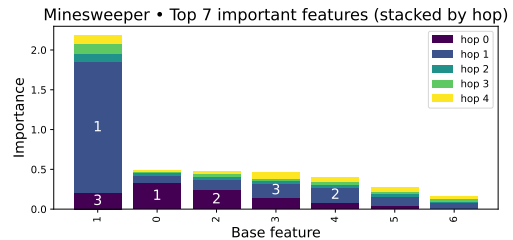


Figure 2: SHAP feature importance for Minesweeper, stacked by hop. Numbers on the stacked bars indicate the ranking of that particular feature on that particular hop.

Additive Explanations (SHAP) (Lundberg & Lee, 2017) on Minesweeper with mean aggregation (Fig. 2), one of the two datasets where FAFs lag. In this dataset, labels are bombs, feature 0 masks other features, and features 1–6 **one-hot** encode the number of neighboring bombs. The top signal is the hop-1 mean of feature 1, i.e., the fraction of neighbors whose local bomb count is **null**. **When this proportion is greater than zero, the model knows that the node cannot have a bomb; when it is zero, all its neighbors observe bombs so there is a possibility of having a bomb.** This heuristic does not completely solve the problem, as neighbors **can be merely observing second-hop bombs**, creating ambiguity that likely underlies residual errors. The model also correctly gives importance to the number of masked neighbors (hop-0 feature 0). These attributions clarify where the model succeeds and where it fails. For comparison, we report SHAP importances for two other datasets in Fig. 5 in the appendix: Pubmed (homophilic) and Amazon-Ratings (heterophilic).

Rather than explaining a particular classifier, one can aim to explain the dataset—localizing which hops and features carry signal independent of model choice. Following Donnelly et al. (2023), noisy tabular datasets often admit a “Rashomon set” of comparably well-performing models. Accordingly, feature importance is better assessed over this set—preferably constrained to simpler or sparser models—than from a single fit. This lens may offer a principled way to interrogate graph data beyond feature homophily–heterophily (Zheng et al., 2024) or strictly one-hop neighborhoods (Zheng et al., 2025), though its current implementation requires binarized data on regression tasks.

**Optimization.** GNNs usually exhibit early overfitting, where training accuracy converges almost immediately while validation and test accuracy plateau or even decay. Thus, the best validation accuracy is often achieved before relevant aggregations are learned. This might partially explain why FAFs can often compete with trained GNNs: They avoid overfitting aggregations. GNNs can also suffer from ineffective aggregation learning (Mustafa et al., 2023; Mustafa & Burkholz, 2024b), so their potential to outcompete FAFs is likely underexplored due to trainability issues. By contrast, optimization on tabular data (like FAFs) is better tractable and understood by standard toolkits.

**Efficiency.** Precomputing aggregation once and then training an MLP on top is far more scalable than repeatedly running message-passing layers and backpropagating through them, as required in GNNs. However, as the number of reducers, original features, and hops in FAF increases, so does the input dimensionality, which in turn enlarges the parameter count of the MLP’s first layer. This issue could be mitigated through common feature reduction techniques. For the original features, we report the average training runtimes of our FAF and GNN models in Table 5. FAFs are generally more efficient, particularly when using a single reducer.

**Augmentations.** Adding more features does not improve accuracy monotonically. Beyond a point, some feature selection is needed. Still, the tabular view lets us concatenate diverse features atop the aggregations. In addition to concatenating multiple **reducers** and hops, one can append structural statistics such as degree, centrality, and other network-science metrics (Blöcker et al., 2025). Our framework is also compatible with pre-processing graph rewiring, where aggregation is computed on a modified adjacency matrix e.g. to fight over-squashing (Topping et al., 2022; Jamadandi et al., 2024). But unlike standard rewiring, we can concatenate the rewired features instead of replacing the originals. As fixed aggregations can suffer from similar issues as trainable GNNs, FAFs can also benefit from proposed remedies. To extract more precise information from complex environments, we examine a feature similarity-based rewiring loosely based on Rubio-Madriral et al. (2025), where edges of negative feature cosine similarity between nodes are dropped. We then append features aggregated on the rewired graph, or split edges into positive/negative sets and aggregate separately, inspired by computational-graph splitting that helps fight over-smoothing (Roth et al., 2025). Results (Table 10) show that on datasets already helped by mean aggregation, adding these features yields larger gains than substituting them; e.g., WikiCS surpasses classic GNNs where other FAF variants do not. These augmentations not only improve performance but also help disentangle where the gains come from: additional extracted signal versus changes to the optimization of graph models. We therefore advocate FAFs as baselines for methods that modify the aggregation component of GNNs—akin to analyses of SGC and GESN (Micheli & Tortorella, 2025), though in our case we obtain benefits from these operations.



## 4 THEORETICAL FOUNDATIONS: DOES AGGREGATION NEED LEARNING?

Let  $G = (V, E)$  be a graph with node features  $\mathcal{X} \in \mathbb{R}^F$ . A neighborhood function is a map  $f : \mathcal{M}(\mathbb{R}^F) \rightarrow \mathbb{R}$  acting on the multiset  $X_v := \{\mathcal{X}_u : u \in N(v)\}$  for  $v \in V$ . We seek a fixed encoder  $\Phi$  such that we can learn any neighborhood function via a univariate map  $g$  with  $f = g \circ \Phi^{-1}$ . This enables tabular classifiers to learn graph data losslessly.

### 4.1 WHAT INFORMATION IS PRESERVED BY STANDARD AGGREGATIONS?

FAFs apply a transformation of neighborhoods that is not learnable, which raises the question of what information gets lost by the aggregation. Permutation invariant aggregations treat graph neighborhoods as multisets consisting of feature vectors of neighbors. Accordingly, they extract distribution information and forget about the identity of specific neighbors. This property is usually regarded as helpful inductive bias and therefore of no concern. Our next theorems analyze which information is preserved by sum and mean aggregations from these multisets. To do so, we first generalize Lemma 5 by Xu et al. (2019) for one-hot encoded discrete features to orthogonal features. Combined with the fact that hop features are concatenated, this insights establishes that information from the 1-hop neighborhood can be preserved.

**Theorem 1** (1-hop aggregation). *Assume the features  $\mathcal{X}$  are orthogonal. Then, the function  $h(X) = \sum_{x \in X} x$  defined on multisets  $X \subseteq \mathcal{X}$  of bounded size is injective. Moreover, any multiset function  $f$  can be decomposed as  $f(X) = g(\sum_{x \in X} x)$  for some function  $g$ .*

The proof is given in the appendix A.2. Note that a multiset  $X \subseteq \mathcal{X}$  is characterized by the count  $n_x$  of elements that have specific features  $x$ . These counts can also be extracted from the sum  $h(X)$  (as demonstrated in the proof). Consequently, any multiset function  $f$  transforms such counts by  $f(n_x)$ . The function  $g$  would thus first extract the counts from the sum  $h(X)$  and then apply  $f$  to the counts. If the features of a node  $v$  include its degree  $d_v$ , then mean aggregation contains the same information, as a classifier can learn to multiply  $h(X) = 1/d_v \sum_{x \in X} x$  by  $d_v$ . In contrast, max and min aggregations extract whether at least one neighbor has a specific feature property. They focus on the tails of distributions rather than full neighborhoods.

**Information loss for  $k$ -hops.** One might hope that the above theorem also applies to aggregations from hop  $k$  to  $k + 1$ . The orthogonality assumption, however, is essential and no longer met by the aggregated neighbor features  $h_n^{(k)}$  for  $k \geq 1$ . As a consequence, from  $k \geq 2$ , not all information about the distribution of features across neighbors is preserved, as Figure 4 exemplifies. In particular,  $h_1^{(2)}$  captures neither the degrees of its neighbors, nor the number and types of second-hop neighbors associated with each first-hop neighbor. Even so, aggregation still extracts useful information, and different aggregations concatenate complementary properties of neighborhoods.

- **Sum aggregation:** Sums count, for each of the  $n$  distinct orthogonal feature vectors  $x_f$ , how many nodes in the  $k$ -hop neighborhood exhibit feature  $f$ . A classifier can extract it by computing  $x_f^T h_v^{(k)}$ . Note that nodes reachable through multiple length- $k$  paths are counted multiple times.
- **Mean aggregation:** Means can partially distinguish neighbors with different degrees by considering the fraction of nodes that exhibit a specific feature vector  $x_f$ . The quantity  $x_f^T h_v^{(k)}$  weights each node  $i$  with feature  $x_f$  by  $1/d_i$ . Note that nodes reachable through multiple length- $k$  paths are again counted with multiplicity.
- **Max aggregation:** Max aggregation on one-hot encoded features returns whether at least one node within  $k$  hops has a given feature. For large neighborhoods as  $k$  increases, this indicator quickly saturates, so increasing hops adds little further information. The same reasoning applies when taking the maximum entry of the orthogonal features.
- **Min aggregation:** The same reasoning applies to the minimum as to the maximum: It indicates whether any node within  $k$  hops lacks the feature, and increasing  $k$  adds little further information.

### 4.2 LOSSLESS NEIGHBORHOOD AGGREGATION

When node features are real-valued in general, Corso et al. (2020) show that no single continuous, permutation-invariant aggregation function can be lossless for all multiset functions. This mirrors

a classical topological obstruction due to Netto (1879): There is no continuous bijection  $\mathbb{R} \rightarrow \mathbb{R}^2$  (Dauben, 1975). However, there can exist discontinuous bijections, namely space filling functions. We adopt a concrete construction based on ternary expansions and the Cantor set, adapted from a Kolmogorov-Arnold representation variant from Theorem 2 by Schmidt-Hieber (2021).

**Theorem 2** (Kolmogorov-Arnold representation from Thm. 2 of Schmidt-Hieber (2021)). *For any fixed  $d \geq 2$ , there exists a monotone function  $\phi : [0, 1] \rightarrow \mathcal{C}$  (the Cantor set) such that the map  $\Phi(x_1, \dots, x_d) = 3 \sum_{p=1}^d 3^{-p} \phi(x_p)$  is injective on  $[0, 1]^d$ . Moreover, for every continuous  $f : [0, 1]^d \rightarrow \mathbb{R}$  there exists a continuous  $g : \Phi([0, 1]^d) \rightarrow \mathbb{R}$  with  $f(x_1, \dots, x_d) = g(\Phi(x_1, \dots, x_d))$ .*

Theorem 2 isolates all required discontinuity into a *fixed* aggregation. While  $\Phi$  is not continuous, its inverse is, which makes the learnable part  $g := f \circ \Phi^{-1}$  inherit the continuity properties of  $f$ . Schmidt-Hieber (2021) has also quantified how much information is lost if  $g$  is learned instead of  $f$ . For  $f$   $\beta$ -smooth with  $\beta \leq 1$ , there is no difference in the rate of approximation. However, for higher order smoothness, the multivariate and univariate function approximation may vary. Note that this aggregation even remembers node identities. From this theorem we can learn the following insight:

A lossless, fixed, even univariate neighborhood aggregation function exists, but it has to be discontinuous for general continuous features.

### 4.3 IMPLICATIONS AND OPEN CHALLENGES

When we encode neighborhoods via the injective function  $\Phi$  and learn  $g$  so that  $f = g \circ \Phi^{-1}$ , the information content, smoothness properties, and approximation rates of the neighborhood function  $f$  transfer to  $g$ . However, this theoretical sufficiency does not guarantee strong empirical performance when  $\Phi$  is used directly as a *reducer* for FAF (see Table 9). In Appendix A.1, we visualize how  $\Phi$  maps 2D circles into the univariate Cantor set, and how  $\Phi^{-1}$  can recover them continuously. We also compare against mean and std, and observe that  $\Phi$  pushes inputs that are close together into far-apart representations, whereas mean and std bring together far-apart inputs that share commonalities. It is the case that, in practice, the simple statistics studied in § 4.1 often provide distributional summaries that downstream classifiers exploit more effectively.

An ideal aggregation function would be both injective, like  $\Phi$ , and would extract useful statistical insights, like mean. It is still an open challenge to design, or potentially learn, efficient embeddings that extract relevant information from graph neighborhoods, while easing the learning problem for the classifier (Burkholz et al., 2022). One might expect GNNs to learn such representations end-to-end without overfitting. Our experiments with FAFs (Tables 2, 6) suggest—despite some information loss at iterative hops—that simple *reducers* suffice for most standard node-classification benchmarks.

Experimentally, we find that mean aggregation alone is often among the top performers. This suggests that neighborhood feature distributions provide most task-relevant signal, and that neighbor degrees encode useful structural information, helping to distinguish the contribution of distinct neighbors. Consistent with this, the most relevant information is already provided in the immediate neighborhood ( $k = 0, 1$ , see Table 6) and the concatenation of this information is key so that it is not lost by repeated aggregations (see Table 8). Consequently, information loss at larger  $k$  is of little practical concern—except for two datasets that appear to require subtler information. Taken together, these observations motivate the following hypothesis.

Hypothesis: For most standard node-classification benchmarks, either the predictive signal is already concentrated within the first one or two hops, or current GNNs struggle to learn layer-wise aggregations that extract relevant information beyond mean or sum.

The first part of this hypothesis underscores the need for more real-world datasets where long-range interactions and richer aggregations matter, supporting prior calls to revisit benchmark design (Errica et al., 2020; Bechler-Speicher et al., 2025). Although some tasks (like Roman Empire) benefit from long-range signal (Topping et al., 2022), making deep graph models work reliably remains a challenge. Recent evidence indicates that graph models generally struggle to capture interactions beyond roughly 13 hops, irrespective of over-smoothing, over-squashing, or vanishing gradients (Zhou et al., 2025).

Table 1: Test accuracy on node classification: Best validation FAF against classic GNNs.

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
GCN	93.58 $\pm$ 0.44	95.77 $\pm$ 0.27	53.86 $\pm$ 0.48	<b>44.62 <math>\pm</math> 4.50</b>	<b>72.72 <math>\pm</math> 0.45</b>	95.73 $\pm$ 0.15	<b>97.47 <math>\pm</math> 0.08</b>
GAT	93.91 $\pm$ 0.22	96.45 $\pm$ 0.37	<b>55.51 <math>\pm</math> 0.55</b>	42.90 $\pm$ 5.47	71.82 $\pm$ 0.65	96.14 $\pm$ 0.08	97.12 $\pm$ 0.13
SAGE	93.31 $\pm$ 0.17	96.17 $\pm$ 0.44	<u>55.26 <math>\pm</math> 0.27</u>	43.11 $\pm$ 4.73	<u>71.82 <math>\pm</math> 0.81</u>	<b>96.21 <math>\pm</math> 0.10</b>	97.10 $\pm$ 0.09
MLP	87.75 $\pm$ 0.42	93.62 $\pm$ 0.36	49.04 $\pm$ 0.39	38.59 $\pm$ 3.29	57.22 $\pm$ 2.25	93.80 $\pm$ 0.19	96.02 $\pm$ 0.16
FAF <sub>bestval</sub>	<b>94.01 <math>\pm</math> 0.21</b>	<b>96.54 <math>\pm</math> 0.13</b>	55.09 $\pm$ 0.24	<u>42.96 <math>\pm</math> 2.45</u>	70.48 $\pm$ 1.24	95.37 $\pm$ 0.17	97.05 $\pm$ 0.18

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
GCN	<b>84.38 <math>\pm</math> 0.81</b>	<u>97.48 <math>\pm</math> 0.06</u>	<u>80.00 <math>\pm</math> 0.77</u>	<u>78.44 <math>\pm</math> 0.23</u>	<b>91.05 <math>\pm</math> 0.15</b>	<u>44.26 <math>\pm</math> 1.22</u>	80.06 $\pm$ 0.81
GAT	83.02 $\pm$ 1.21	97.00 $\pm$ 1.02	79.80 $\pm$ 0.94	77.72 $\pm$ 0.71	90.38 $\pm$ 0.49	39.31 $\pm$ 2.42	<b>81.01 <math>\pm</math> 0.23</b>
SAGE	<u>83.18 <math>\pm</math> 0.93</u>	<b>97.72 <math>\pm</math> 0.70</b>	77.42 $\pm$ 0.40	76.75 $\pm$ 1.07	<u>90.41 <math>\pm</math> 0.10</u>	40.22 $\pm$ 1.47	<u>80.57 <math>\pm</math> 0.42</u>
MLP	58.56 $\pm$ 1.75	51.74 $\pm$ 0.83	68.22 $\pm$ 0.96	70.40 $\pm$ 1.17	66.43 $\pm$ 0.12	39.11 $\pm$ 1.93	72.98 $\pm$ 0.49
FAF <sub>bestval</sub>	82.84 $\pm$ 0.63	90.00 $\pm$ 0.39	<b>80.96 <math>\pm</math> 1.06</b>	<b>78.69 <math>\pm</math> 0.50</b>	78.11 $\pm$ 0.38	<b>44.59 <math>\pm</math> 1.62</b>	80.25 $\pm$ 0.34

The second part of the hypothesis concerns the ability of GNNs to actually realize useful aggregations in practice. For instance, GNNs may not move far enough from their initializations. Indeed, the two datasets on which GNNs hold an advantage require linear residual transformations to realize that gap (Luo et al., 2024). Prior work also shows that GATs cannot flexibly adjust attention to shut off unhelpful neighbors (Mustafa & Burkholz, 2024b). This supports our results on rewiring the adjacency matrix before aggregation (Table 10). If GATs could learn to prune the edges that we manually drop, they would enjoy similar gains.

We see opportunity for future work along three fronts that build directly on our findings:

1. Feature/reducer engineering: FAFs highlight untapped potential for designing meaningful node features that encode graph structure, require less learning, potentially preserve more—but ideally only relevant—information, and allow for higher learning efficiency. In combination with partial feature learning, they might form the basis of a new generation of graph based learning architectures.
2. Moving beyond injectivity: As our theory and empirical results highlight, improving GNN expressiveness and thus injectivity alone is not likely to inspire practical improvements on current benchmarks, as those can be competitively solved even with simple, non-injective aggregation. We therefore call for a shift in focus from mere injectivity to other learning properties—a theoretical gap to be addressed not only for FAFs but for GNNs in general.
3. Benchmarks: Enough information to solve current benchmark tasks is already contained in early hops and can be extracted with simple, non-injective aggregation. If we really want to showcase the capabilities of GNNs to learn meaningful features, we need more difficult benchmarks that require this ability.

In theory, fixed information-preserving aggregations can reduce graph learning to tabular prediction. In practice, task relevant representations and information preservation are a challenge. Progress likely requires both more amenable reducers and better tasks for evaluation.

## 5 EXPERIMENTS

### 5.1 COMPARISON TO CLASSIC GNNs

**Performance of FAFs.** Table 1 reports test performance for classic GNNs—GCN (Kipf & Welling, 2017), GATv2 (Brody et al., 2022), and GraphSAGE (Hamilton et al., 2017)—versus our approach, which feeds Fixed Aggregation Features (FAFs) into MLPs (Figure 1 and Eq. 1). We aggregate up to the same hop depth as the GNN baselines. As a control, we include an MLP baseline with zero-hop aggregation, which performs substantially worse than all other models. We obtain the best FAF variant from validation results, shown in Table 2. FAF<sub>4</sub> uses the reducers  $\mathcal{R} = \{\text{mean, sum, max, min}\}$ , and is tuned with exactly the hyperparameter grid from Luo et al. (2024); this makes our results directly comparable to their Graph Transformers and heterophily-aware architectures, where they find that classic GNNs can also rival them. Additional FAF variants include



Table 2: Validation accuracy on node classification: FAFs against classic GNNs.

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
GCN	92.58 $\pm$ 0.10	95.42 $\pm$ 0.11	54.01 $\pm$ 0.23	48.15 $\pm$ 2.35	<b>70.36 <math>\pm</math> 0.09</b>	<u>95.32 <math>\pm</math> 0.07</u>	<b>97.16 <math>\pm</math> 0.07</b>
GAT	92.86 $\pm$ 0.06	95.93 $\pm$ 0.15	55.56 $\pm$ 0.68	46.97 $\pm$ 2.07	<u>69.52 <math>\pm</math> 0.27</u>	95.30 $\pm$ 0.08	<u>97.11 <math>\pm</math> 0.03</u>
SAGE	92.33 $\pm$ 0.17	95.60 $\pm$ 0.16	<b>55.90 <math>\pm</math> 0.54</b>	46.22 $\pm$ 2.12	68.48 $\pm$ 1.05	<b>95.51 <math>\pm</math> 0.04</b>	97.02 $\pm$ 0.09
MLP	87.89 $\pm$ 0.13	93.33 $\pm$ 0.07	48.98 $\pm$ 0.72	41.43 $\pm$ 1.77	56.80 $\pm$ 1.24	93.70 $\pm$ 0.07	95.89 $\pm$ 0.02
FAF <sub>4</sub>	<u>93.05 <math>\pm</math> 0.04</u>	<b>96.34 <math>\pm</math> 0.07</b>	55.53 $\pm$ 0.43	<b>48.51 <math>\pm</math> 2.31</b>	67.28 $\pm$ 0.64	94.93 $\pm$ 0.07	96.83 $\pm$ 0.01
FAF <sub>mean,std</sub>	93.04 $\pm$ 0.13	<u>96.23 <math>\pm</math> 0.08</u>	55.11 $\pm$ 0.40	<u>48.42 <math>\pm</math> 1.64</u>	67.20 $\pm$ 0.28	94.94 $\pm$ 0.07	96.84 $\pm$ 0.03
FAF <sub>mean</sub>	<b>93.16 <math>\pm</math> 0.04</b>	96.06 $\pm$ 0.10	53.78 $\pm$ 0.52	47.99 $\pm$ 2.02	66.92 $\pm$ 0.87	95.20 $\pm$ 0.14	97.00 $\pm$ 0.04
FAF <sub>max,std</sub>	92.32 $\pm$ 0.08	95.80 $\pm$ 0.04	<u>55.70 <math>\pm</math> 0.45</u>	48.42 $\pm$ 2.14	66.64 $\pm$ 0.54	95.04 $\pm$ 0.04	96.56 $\pm$ 0.03
FAF <sub>max</sub>	91.93 $\pm$ 0.04	95.60 $\pm$ 0.04	55.63 $\pm$ 0.29	48.06 $\pm$ 2.30	66.56 $\pm$ 0.50	95.19 $\pm$ 0.13	96.54 $\pm$ 0.01
FAF <sub>sum</sub>	90.95 $\pm$ 0.04	94.88 $\pm$ 0.04	53.48 $\pm$ 0.59	47.29 $\pm$ 1.92	67.84 $\pm$ 1.45	95.13 $\pm$ 0.09	96.65 $\pm$ 0.05
FAF <sub>std</sub>	92.50 $\pm$ 0.06	95.86 $\pm$ 0.10	55.31 $\pm$ 0.32	47.27 $\pm$ 2.15	63.44 $\pm$ 0.17	95.01 $\pm$ 0.12	96.76 $\pm$ 0.02

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
GCN	81.28 $\pm$ 0.33	<u>97.36 <math>\pm</math> 0.46</u>	79.08 $\pm$ 0.23	78.63 $\pm$ 0.23	<b>91.14 <math>\pm</math> 0.58</b>	44.88 $\pm$ 1.27	81.52 $\pm$ 0.37
GAT	81.16 $\pm$ 0.52	97.08 $\pm$ 1.16	78.84 $\pm$ 0.52	78.12 $\pm$ 1.03	<u>90.49 <math>\pm</math> 0.68</u>	43.30 $\pm$ 1.43	<b>82.38 <math>\pm</math> 0.57</b>
SAGE	81.32 $\pm$ 0.41	<b>97.68 <math>\pm</math> 0.63</b>	78.88 $\pm$ 0.91	77.35 $\pm$ 1.09	90.44 $\pm$ 0.66	40.58 $\pm$ 1.17	<u>82.27 <math>\pm</math> 0.38</u>
MLP	62.68 $\pm$ 1.15	51.12 $\pm$ 0.93	71.12 $\pm$ 0.52	71.58 $\pm$ 1.46	66.28 $\pm$ 0.27	40.57 $\pm$ 0.92	74.86 $\pm$ 0.33
FAF <sub>4</sub>	82.84 $\pm$ 0.43	89.63 $\pm$ 1.03	79.08 $\pm$ 0.36	<b>79.53 <math>\pm</math> 1.12</b>	78.68 $\pm$ 0.19	<u>47.31 <math>\pm</math> 1.39</u>	81.92 $\pm$ 0.43
FAF <sub>mean,std</sub>	<b>83.36 <math>\pm</math> 0.17</b>	89.18 $\pm$ 0.71	<b>81.28 <math>\pm</math> 0.30</b>	77.32 $\pm$ 0.36	77.59 $\pm$ 0.41	47.30 $\pm$ 1.32	81.37 $\pm$ 0.51
FAF <sub>mean</sub>	<u>83.28 <math>\pm</math> 0.30</u>	89.89 $\pm$ 0.93	81.16 $\pm$ 0.97	78.53 $\pm$ 0.87	76.67 $\pm$ 0.36	46.29 $\pm$ 1.50	81.58 $\pm$ 0.46
FAF <sub>max</sub>	81.80 $\pm$ 0.42	86.08 $\pm$ 0.77	77.48 $\pm$ 0.30	<u>79.15 <math>\pm</math> 0.86</u>	75.06 $\pm$ 0.14	46.47 $\pm$ 1.38	80.30 $\pm$ 0.56
FAF <sub>max,std</sub>	82.08 $\pm$ 0.33	87.83 $\pm$ 0.63	78.28 $\pm$ 0.30	78.86 $\pm$ 0.89	76.19 $\pm$ 0.26	<b>47.44 <math>\pm</math> 1.51</b>	80.46 $\pm$ 0.53
FAF <sub>sum</sub>	82.60 $\pm$ 0.65	89.86 $\pm$ 0.85	79.40 $\pm$ 0.57	78.12 $\pm$ 0.27	77.13 $\pm$ 0.23	46.85 $\pm$ 1.28	78.17 $\pm$ 0.23
FAF <sub>std</sub>	81.40 $\pm$ 0.51	88.20 $\pm$ 0.52	80.00 $\pm$ 0.40	76.25 $\pm$ 0.53	73.95 $\pm$ 0.49	45.91 $\pm$ 1.32	77.65 $\pm$ 0.31

mean+std, mean only, max+std, max only, sum only, and std only. The best overall result is shown in **bold**, the second best is underlined. More details on the setup are given in § B, [as well as all test accuracy counterparts in Table 11 in § D](#). We also show all training, validation, and test curves of all datasets for FAF<sub>4</sub> and GCN in Fig. 6.

Overall, we improve on 5 datasets, match within error or **1% on another 5**, and trail on **4**. On most datasets, FAF<sub>4</sub> performs comparably to mean+std. [Among the ones within 1%, we have Coauthor-CS and Coauthor-Physics \(Figures 6f, 6g\)](#), which are the largest and most feature-rich; targeted feature selection may close the gap. Among the 4 trailing datasets, two are homophilic and two heterophilic; the homophilic tasks are close to parity. Citeseer exhibits optimization instability (Fig. 6e), [and Cora has a large test-validation gap in GCNs \(Fig. 6h\)](#), not present in any other. The two heterophilic datasets, Minesweeper and Roman-Empire (Figures 6i, 6l) show larger performance drops. This behavior mirrors the decrease reported by Luo et al. (2024) when residual connections are removed. Notably, the best-performing FAFs on these two datasets use far fewer hops (4 and 2) than the GNN baselines (15 and 10), suggesting that key signal lies at longer ranges. The shallower FAFs under-aggregate relative to what those tasks require, but adding extra hops does not provide extra information, as discussed in § 4.1. We also show it in Table 6, where we concatenate up to different amount of hops. In fact, most datasets peak at  $K = 2$ , and either plateau or decrease in performance.

**Best hyperparameters.** All FAF variants benefit from normalization components (BatchNorm or LayerNorm), as aggregated features can vary widely in scale across reducers and hops. Compared with GNNs, FAFs typically favor larger learning rates, which can yield faster training, improved generalization via implicit regularization, and feature sparsity (Mohtashami et al., 2023; Sadrtadinov et al., 2024). Dropout levels, however, are broadly similar to those used for GNNs. This suggests that dropout’s gains on these node-classification tasks are driven more by dataset properties than by the specifics of training graph convolutions, which nuances prior interpretations (Luo et al., 2025b).

## 5.2 ABLATIONS

**Ablation on single reducers.** Concatenating multiple aggregations has advantages and drawbacks. On the plus side, an MLP can learn to weight each reducer, removing the need to pick one per dataset. Because our individual reducers are not lossless, different datasets may favor different ones; moreover, adding informative, correlated covariates can improve robustness and reduce variance (Reddy et al., 2025). On the downside, concatenation increases input dimensionality, with

corresponding memory and optimization costs. Table 2 reports validation results when using a single aggregation at a time. **Note that this resembles a simple feature selection over FAF<sub>4</sub>.** We keep the same hyperparameters as FAF<sub>4</sub> to isolate the effect of the aggregation choice, though the lower dimensionality could allow for alternative settings that further improve performance. Surprisingly, a single **reducer** often suffices, though the preferred choice varies by dataset. The mean is most frequently strongest, sometimes surpassing FAF<sub>4</sub> and FAF<sub>mean,std</sub>—e.g., on Pubmed, where FAF<sub>mean</sub> outperforms the GNNs while other FAF variants do not. This may reflect optimization challenges from high-dimensional inputs or increased overfitting. Still, mean is not universally best: sum or max win on some datasets (e.g., Citeseer favors sum; Amazon-Ratings favors max). Combining **reducers** therefore remains beneficial when one wishes to avoid committing to a specific one a priori.

**Comparison with one-layer classifier and last hop.** Other simplifications of GNNs (Wu et al., 2019; Micheli & Tortorella, 2024) effectively fix the aggregation and train a single linear layer on the final-hop representation. In contrast, we concatenate representations from all hops and train a well-tuned MLP classifier. This choice is crucial for matching GNN performance. As shown in Table 7, MLPs consistently outperform a single linear layer applied to the same concatenated features, indicating that their nonlinearity and increased capacity are important to learn from multi-hop features. Moreover, Table 8 shows that only using the last hop lacks important information that is not transmitted across aggregations.

**Kolmogorov-Arnold aggregation.** Our hypothesis that Roman-Empire lags due to information loss is reinforced by a FAF variant that uses the Kolmogorov–Arnold (KA) **function**  $\Phi$ , which is theoretically lossless (Theorem 2). As we also exemplified in Fig. 3, in practice, KA is hard to use for classification. We observe in Table 9 that some datasets struggle to fit (lower training accuracy), while others show mild overfitting. Nevertheless, on Roman-Empire this variant attains a validation accuracy of  $80.45 \pm 0.25$ , the highest among all FAFs, suggesting that providing full neighborhood information helps close the gap on this task. This, in turn, highlights the need for benchmarks where predictive signal genuinely arises at distant hops in complex ways.

## 6 CONCLUSIONS

We have introduced Fixed Aggregation Features (FAFs), a non-learnable tabular mapping from local neighborhoods of graph features to univariate representations that an MLP can learn to classify. Our analysis shows that fixed, injective neighborhood aggregation **functions** exist, linking multiset expressivity to Kolmogorov–Arnold factorizations; thus learned message passing is not required for expressivity in theory. But in practice, common non-injective **reducers** (mean, sum, max, min) train more reliably, underscoring a gap between what is expressive in principle and what is reliably learnable. We also highlight the practical advantages of a tabular view, such as access to the rich tabular toolkit of interpretability and tuning, and isolated representation from inference so we can attribute gains or failures to the features themselves rather than to message-passing optimization.

On node classification datasets, FAFs are a strong baseline: they match or beat classic GNNs on many benchmarks, and trail only on two datasets needing longer-range interactions, where residualized GNNs help. Two ablations explain most gains: a well-tuned MLP beats a single linear classifier on top of FAF features, and concatenating all hop beats using only the last hop. This is consistent with later hops losing detail for these practical aggregation **schemes**. Surprisingly, two hops usually suffice, suggesting either limited signal in current benchmarks, or difficulty training deep GNNs to exploit more of it. While our theory carries over other downstream tasks, other benchmarks may surface different constraints that can alter the empirical outcomes.

Our findings have immediate implications. We recommend always including a tuned FAF baseline in future studies to calibrate what fixed aggregation alone can achieve; re-evaluating—and, when appropriate, retiring—datasets on which FAFs reach state-of-the-art performance; and developing benchmarks that genuinely require long-range dependencies and inter-hop dynamics. More broadly, we advocate for simplifying models and balancing expressiveness against optimizability, rather than assuming that extra parameters or higher expressiveness extract more relevant signal than simple baselines. Notably, several phenomena that are often blamed on graph architectures—overfitting, depth-related degradation, and sensitivity to dropout—also arise in tabular settings, indicating that some limitations may stem from dataset properties rather than the graph-aware architectures alone.

## REPRODUCIBILITY STATEMENT

Experimental details are provided in Appendix B, and further experimental results are included in Appendix C. Detailed train/validation/test performance across all epochs of the main runs can be found in Appendix D. The code for the experiments is attached as supplement.

## LLM STATEMENT

To improve fluency of the text sentence level, editing has been done using large language models.

## REFERENCES

- Anonymous. Graphbench: Next-generation graph learning benchmarking. 2025. URL <https://openreview.net/forum?id=h4FOb6ak0Q>. under review.
- Gleb Bazhenov, Oleg Platonov, and Liudmila Prokhorenkova. Graphland: Evaluating graph machine learning models on diverse industrial data. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025. URL <https://openreview.net/forum?id=Gyq8lMgdk5>.
- Maya Bechler-Speicher, Ben Finkelshtein, Fabrizio Frasca, Luis Müller, Jan Tönshoff, Antoine Siraudin, Viktor Zaverkin, Michael M. Bronstein, Mathias Niepert, Bryan Perozzi, Mikhail Galkin, and Christopher Morris. Position: Graph learning will lose relevance due to poor benchmarks. In *Forty-second International Conference on Machine Learning Position Paper Track*, 2025. URL <https://openreview.net/forum?id=nDFp12lhoH>.
- Christopher Blöcker, Martin Rosvall, Ingo Scholtes, and Jevin D. West. Insights from network science can advance deep graph learning, 2025. URL <https://arxiv.org/abs/2502.01177>.
- Pietro Bongini, Niccolò Pancino, Franco Scarselli, and Monica Bianchini. Biognn: How graph neural networks can solve biological problems. In *Artificial Intelligence and Machine Learning for Healthcare*, pp. 211–231. Springer, 2023.
- Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=F72ximsx7C1>.
- Rebekka Burkholz, Nilanjana Laha, Rajarshi Mukherjee, and Alkis Gotovos. On the existence of universal lottery tickets. In *International Conference on Learning Representations*, 2022.
- Chen Cai and Yusu Wang. A simple yet effective baseline for non-attributed graph classification. In *ICLR Workshop: Representation Learning on Graphs and Manifolds*, 2019. URL <https://rlgm.github.io/papers/5.pdf>.
- Gianluca De Carlo, Andrea Mastropietro, and Aris Anagnostopoulos. Kolmogorov–arnold graph neural networks, 2025. URL <https://openreview.net/forum?id=udfjje2xXb>.
- Jinsong Chen, Kaiyuan Gao, Gaichao Li, and Kun He. NAGphormer: A tokenized graph transformer for node classification in large graphs. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=8KYeilT3Ow>.
- Eli Chien, Jianhao Peng, Pan Li, and Olga Milenkovic. Adaptive universal generalized PageRank graph neural network. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=n6jl7fLxrP>.
- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Velickovic. Principal neighbourhood aggregation for graph nets. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS ’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

- Corinna Coupette, Jeremy Wayland, Emily Simons, and Bastian Rieck. No metric to rule them all: Toward principled evaluations of graph-learning datasets. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=XbmBNwrfG5>.
- Joseph W Dauben. The invariance of dimension: Problems in the early development of set theory and topology [1]. *Historia Mathematica*, 2(3):273–288, 1975. ISSN 0315-0860. doi: [https://doi.org/10.1016/0315-0860\(75\)90066-X](https://doi.org/10.1016/0315-0860(75)90066-X). URL <https://www.sciencedirect.com/science/article/pii/031508607590066X>.
- Chenhui Deng, Zichao Yue, and Zhiru Zhang. Polynormer: Polynomial-expressive graph transformer in linear time. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hmv1LpNfXa>.
- Giovanni Donghi, Luca Pasa, Luca Oneto, Claudio Gallicchio, Alessio Micheli, Davide Anguita, Alessandro Sperduti, and Nicolò Navarin. Investigating over-parameterized randomized graph networks. *Neurocomputing*, 606:128281, 2024. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2024.128281>. URL <https://www.sciencedirect.com/science/article/pii/S092523122401052X>.
- Jon Donnelly, Srikar Katta, Cynthia Rudin, and Edward P Browne. The rashomon importance distribution: Getting RID of unstable, single model-based variable importance. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=TczT2jiPT5>.
- Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph neural networks for graph classification. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HygDF6NFPB>.
- Claudio Gallicchio and Alessio Micheli. Graph echo state networks. In *The 2010 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, 2010. doi: 10.1109/IJCNN.2010.5596796.
- Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Combining neural networks with personalized pagerank for classification on graphs. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1gL-2A9Ym>.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML’17*, pp. 1263–1272. JMLR.org, 2017.
- Yury Gorishniy, Ivan Rubachev, and Artem Babenko. On embeddings for numerical features in tabular deep learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 24991–25004. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/9e9f0ffc3d836836ca96cbf8fe14b105-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/9e9f0ffc3d836836ca96cbf8fe14b105-Paper-Conference.pdf).
- William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 1025–1035, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Adarsh Jamadandi, Celia Rubio-Madrigal, and Rebekka Burkholz. Spectral graph pruning against over-squashing and over-smoothing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=EMkrwJY2de>.
- Dimitrios Kelesis, Dimitris Fotakis, and Georgios Paliouras. Partially trained graph convolutional networks resist oversmoothing. *Machine Learning*, 114(10):211, Aug 2025. ISSN 1573-0565. doi: 10.1007/s10994-025-06865-3. URL <https://doi.org/10.1007/s10994-025-06865-3>.

- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=SJU4ayYgl>.
- Andrei Nikolaevich Kolmogorov. On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition. In *Doklady Akademii Nauk*, volume 114(5), pp. 953–956. Russian Academy of Sciences, 1957.
- Kezhi Kong, Jiuhai Chen, John Kirchenbauer, Renkun Ni, C. Bayan Bruss, and Tom Goldstein. GOAT: A global transformer on large-scale graphs. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 17375–17390. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/kong23a.html>.
- Xiang Li, Renyu Zhu, Yao Cheng, Caihua Shan, Siqiang Luo, Dongsheng Li, and Weining Qian. Finding global homophily in graph neural networks when meeting heterophily. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 13242–13256. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/li22ad.html>.
- Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle, James Halverson, Marin Soljagic, Thomas Y. Hou, and Max Tegmark. KAN: Kolmogorov–arnold networks. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Ozo7qJ5vZi>.
- Scott M. Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, pp. 4768–4777, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- Yuankai Luo, Lei Shi, and Xiao-Ming Wu. Classic GNNs are strong baselines: Reassessing GNNs for node classification. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL <https://openreview.net/forum?id=xkljKdGe4E>.
- Yuankai Luo, Lei Shi, and Xiao-Ming Wu. Can classic GNNs be strong baselines for graph-level tasks? simple architectures meet excellence. In *Forty-second International Conference on Machine Learning*, 2025a. URL <https://openreview.net/forum?id=ZH7YgIZ3DF>.
- Yuankai Luo, Xiao-Ming Wu, and Hao Zhu. Beyond random masking: When dropout meets graph convolutional networks. In *The Thirteenth International Conference on Learning Representations*, 2025b. URL <https://openreview.net/forum?id=PwxYoMvmvy>.
- Sunil Kumar Maurya, Xin Liu, and Tsuyoshi Murata. Simplifying approach to node classification in graph neural networks. *Journal of Computational Science*, 62:101695, 2022. ISSN 1877-7503. doi: <https://doi.org/10.1016/j.jocs.2022.101695>. URL <https://www.sciencedirect.com/science/article/pii/S1877750322000990>.
- Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3(2): 127–163, 2000. doi: 10.1023/A:1009953814988. URL <https://doi.org/10.1023/A:1009953814988>.
- Péter Mernyei and Cătălina Cangea. Wiki-cs: A wikipedia-based benchmark for graph neural networks. In *ICML Workshop: Graph Representation Learning and Beyond*, 2020. URL <https://arxiv.org/abs/2007.02901>.
- Alessio Micheli and Domenico Tortorella. Designs of graph echo state networks for node classification. *Neurocomputing*, 597:127965, 2024. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2024.127965>. URL <https://www.sciencedirect.com/science/article/pii/S0925231224007367>.



- Alessio Micheli and Domenico Tortorella. An empirical evaluation of rewiring approaches in graph neural networks. *Pattern Recognition Letters*, 196:134–141, 2025. ISSN 0167-8655. doi: <https://doi.org/10.1016/j.patrec.2025.05.021>. URL <https://www.sciencedirect.com/science/article/pii/S0167865525002168>.
- Amirkeivan Mohtashami, Martin Jaggi, and Sebastian U Stich. Special properties of gradient descent with large learning rates. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 25082–25104. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/mohtashami23a.html>.
- Nimrah Mustafa and Rebekka Burkholz. Dynamic rescaling for training GNNs. In *Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL <https://openreview.net/forum?id=IfZwSRpqHl>.
- Nimrah Mustafa and Rebekka Burkholz. GATE: How to keep out intrusive neighbors. In *Forty-first International Conference on Machine Learning*, 2024b. URL <https://openreview.net/forum?id=Sjv5RcqfuH>.
- Nimrah Mustafa, Aleksandar Bojchevski, and Rebekka Burkholz. Are GATs out of balance? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=qY7UqLooRa>.
- Galileo Namata, Ben London, Lise Getoor, and Bert Huang. Query-driven active surveying for collective classification. In *Proceedings of the Workshop on Mining and Learning with Graphs*, 2012.
- Oleg Platonov, Denis Kuznedelev, Artem Babenko, and Liudmila Prokhorenkova. Characterizing graph datasets for node classification: homophily-heterophily dichotomy and beyond. In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23*, Red Hook, NY, USA, 2023a. Curran Associates Inc.
- Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova. A critical look at evaluation of gnns under heterophily: Are we really making progress? In *The Eleventh International Conference on Learning Representations*, 2023b.
- Lisi Qarkaxhija, Anatol Eugen Wegner, and Ingo Scholtes. Link prediction with untrained message passing layers. In *The Third Learning on Graphs Conference*, 2024. URL <https://openreview.net/forum?id=IRFFbpbdi9>.
- Sandeep Ramachandra, Vic Degraeve, Gilles Vandewiele, Bram Steenwinkel, Sofie Van Hoecke, and Femke Ongena. Rr-gcn: Exploring untrained random embeddings for relational graphs. *International Journal of Software Engineering and Knowledge Engineering*, 35(06):809–834, 2025. doi: 10.1142/S0218194025500184. URL <https://doi.org/10.1142/S0218194025500184>.
- Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 14501–14515. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/5d4834a159f1547b267a05a4e2b7cf5e-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/5d4834a159f1547b267a05a4e2b7cf5e-Paper-Conference.pdf).
- Abbavaram Gowtham Reddy, Celia Rubio-Madrigal, Rebekka Burkholz, and Krikamol Muandet. When shift happens - confounding is to blame, 2025. URL <https://arxiv.org/abs/2505.21422>.
- Andreas Roth, Franka Bause, Nils Morten Kriege, and Thomas Liebig. Preventing representational rank collapse in mpnns by splitting the computational graph. In Guy Wolf and Smita Krishnaswamy (eds.), *Proceedings of the Third Learning on Graphs Conference*, volume 269 of *Proceedings of Machine Learning Research*, pp. 14:1–14:24. PMLR, 26–29 Nov 2025. URL <https://proceedings.mlr.press/v269/roth25a.html>.

- Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-Scale attributed node embedding. *Journal of Complex Networks*, 9(2), 05 2021. ISSN 2051-1329. doi: 10.1093/comnet/cnab014. URL <https://doi.org/10.1093/comnet/cnab014>.
- Celia Rubio-Madrigal, Adarsh Jamadandi, and Rebekka Burkholz. GNNs getting ComFy: Community and feature similarity guided rewiring. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=g6v09VxgFw>.
- Ildus Sadrtidinov, Maxim Kodryan, Eduard Pokonechny, Ekaterina Lobacheva, and Dmitry Vetrov. Where do large learning rates lead us? In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 58445–58479. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/6b7375226d4742ff910618a56ae72b7d-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/6b7375226d4742ff910618a56ae72b7d-Paper-Conference.pdf).
- Johannes Schmidt-Hieber. The Kolmogorov–Arnold representation theorem revisited. *Neural Networks*, 137:119–126, 2021. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2021.01.020>. URL <https://www.sciencedirect.com/science/article/pii/S0893608021000289>.
- Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93, Sep. 2008. doi: 10.1609/aimag.v29i3.2157. URL <https://ojs.aaai.org/index.php/aimagazine/article/view/2157>.
- Kartik Sharma, Yeon-Chang Lee, Sivagami Nambi, Aditya Salian, Shlok Shah, Sang-Wook Kim, and Srijan Kumar. A survey of graph neural networks for social recommender systems. *ACM Computing Surveys*, 56(10):1–34, June 2024. ISSN 1557-7341. doi: 10.1145/3661821. URL <http://dx.doi.org/10.1145/3661821>.
- Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls of graph neural network evaluation. *Relational Representation Learning Workshop, NeurIPS*, 2018.
- Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J. Sutherland, and Ali Kemal Sinop. Exphormer: Sparse transformers for graphs. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 31613–31632. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/shirzad23a.html>.
- Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=7UmjRGzp-A>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rJXMpikCZ>. accepted as poster.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6861–6871. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/wu19e.html>.
- Qitian Wu, Wentao Zhao, Zenan Li, David P Wipf, and Junchi Yan. Nodeformer: A scalable graph structure learning transformer for node classification. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 27387–27401. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/af790b7ae573771689438bbcf5933fe-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/af790b7ae573771689438bbcf5933fe-Paper-Conference.pdf).

- Qitian Wu, Wentao Zhao, Chenxiao Yang, Hengrui Zhang, Fan Nie, Haitian Jiang, Yatao Bian, and Junchi Yan. SGFormer: Simplifying and empowering transformers for large-graph representations. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 64753–64773. Curran Associates, Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/cc57fac10eacadb3b72a907ac48f9a98-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/cc57fac10eacadb3b72a907ac48f9a98-Paper-Conference.pdf).
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- Yilun Zheng, Sitao Luan, and Lihui Chen. What is missing for graph homophily? disentangling graph homophily for graph neural networks. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 68406–68452. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/7e810b2c75d69bel186cadd2fe3febeab-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/7e810b2c75d69bel186cadd2fe3febeab-Paper-Conference.pdf).
- Yilun Zheng, Xiang Li, Sitao Luan, Xiaojiang Peng, and Lihui Chen. Let your features tell the differences: Understanding graph convolution by feature splitting. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=I9omfcWfMp>.
- Dongzhuoran Zhou, Evgeny Kharlamov, and Egor V. Kostylev. GLora: A benchmark to evaluate the ability to learn long-range dependencies in graphs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=2jf5x5XoYk>.
- Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond homophily in graph neural networks: Current limitations and effective designs. *Advances in Neural Information Processing Systems*, 33, 2020.
- Jiong Zhu, Ryan A. Rossi, Anup Rao, Tung Mai, Nedim Lipka, Nesreen K. Ahmed, and Danai Koutra. Graph neural networks with heterophily. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):11168–11176, May 2021. doi: 10.1609/aaai.v35i12.17332. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17332>.

## APPENDIX

## A STUDY OF REDUCERS FOR NEIGHBORHOOD AGGREGATION

A.1 KOLMOGOROV-ARNOLD FUNCTION  $\Phi$  AND ITS CONTINUOUS INVERSE

Here we showcase the behavior of  $\Phi$  from the Kolmogorov-Arnold representation from Thm. 2. In Fig. 3a, we see the  $\Phi$  image of two circles colored by their angle. Colors that were close together end up in separate parts of the Cantor set; for instance, oranges and reds, or purples and blues. In contrast, in Fig. 3b we see  $\Phi^{-1}$  maps the Cantor set to the circles in such a way that all colors maintain their closeness.

If we use  $\Phi$  as a fixed neighborhood aggregation, the classifier on top needs to learn to reverse it, therefore it is advantageous to have a continuous inverse. However, this does not give information about neighborhood distributions like the commonly used mean, sum, and max. In Fig. 3c we show the behavior of mean and std; mean gives approximate location but fuses together points that are very far apart. For instance, blues and reds have an unusually large first (but different) coordinate and are mapped to the center; however, this information can be recovered with std.

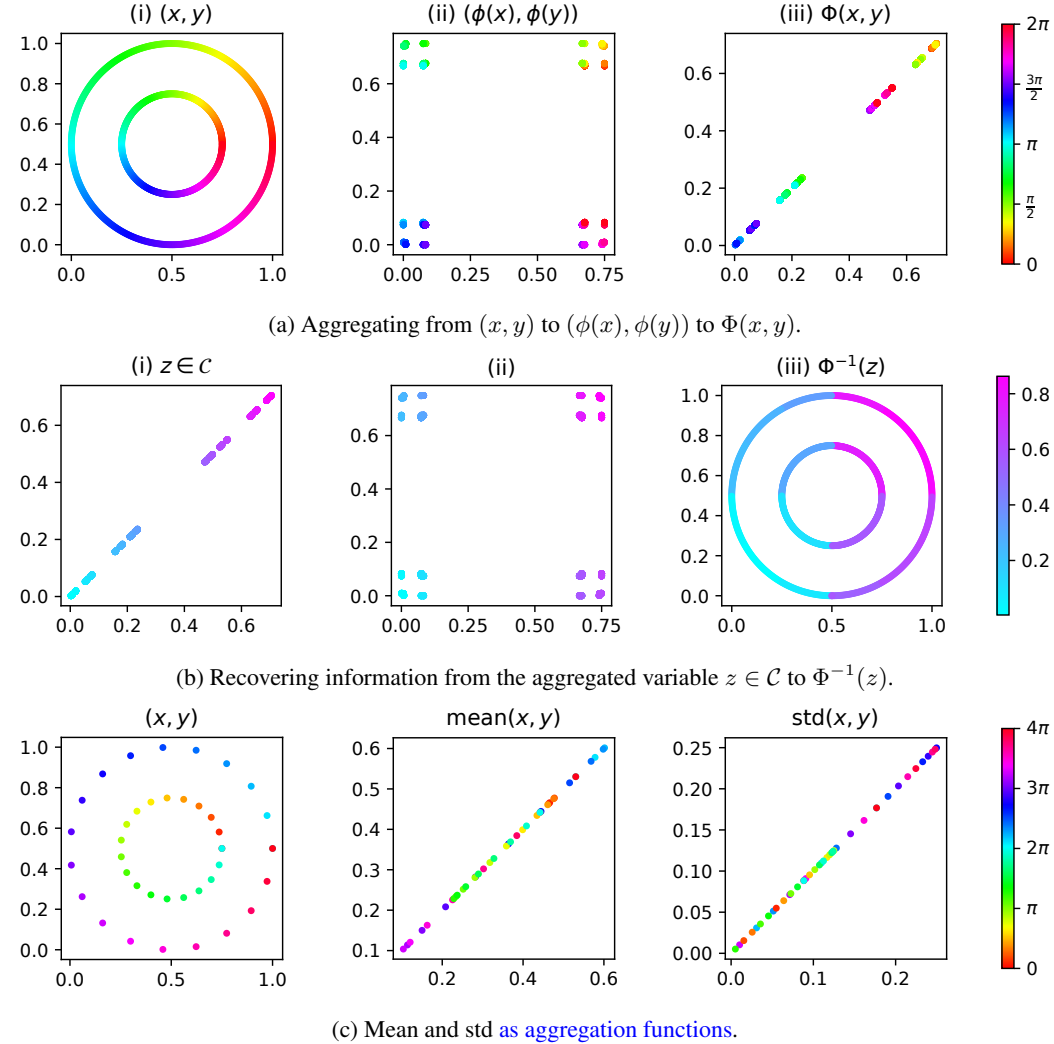


Figure 3: Functions  $\Phi$  (Thm 2)—and its inverse—, mean and std. Circles and square-like panels (a.i, a.ii, b.ii, b.iii, c.i) live in the 2D space, while segments and Cantor sets (a.iii, b.i, c.ii, c.iii) live in 1D. Colors in (a) and (c) are based on angles on 2D, while colors in (b) are based on position.

## A.2 PROOF OF MAIN THEOREM

For convenience, the following theorem restates Theorem 1 of the main paper.

**Theorem** (1-hop aggregation). *Assume the features  $\mathcal{X}$  are orthogonal. Then, the function  $h(X) = \sum_{x \in X} x$  defined on multisets  $X \subseteq \mathcal{X}$  of bounded size is *injective*. Moreover, any multiset function  $f$  can be decomposed as  $f(X) = g(\sum_{x \in X} x)$  for some function  $g$ .*

*Proof.* Note that the multiset  $X$  is fully characterized by the number  $n_f$  of nodes in the set that have a feature  $x_f$  for all possible features  $x_f$ . Our objective is to show that this information is contained in the aggregated form  $h(X) = \sum_{x \in X} x$ .

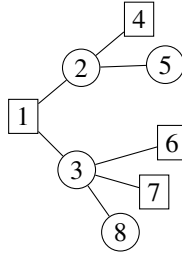
So let us assume that the features are orthogonal. Accordingly, the features  $x_v$  of each node  $v$  assume one of a finite number of possible states  $x_1, \dots, x_n \in \mathbb{R}^{n_f}$  with  $n_f \geq n$  and  $x_i^T x_j = 0$  for any pair  $i, j \in V$  with  $i \neq j$ . Note that the number of possible feature states  $n$  must be finite even in an infinitely large graph, as long as the number of features are finite, i.e.  $n_f < \infty$ . Since the feature values must be pairwise orthogonal, there can maximally exist  $n_f$  distinct feature vectors, as  $n_f$  orthogonal vectors would form a basis of  $\mathbb{R}^{n_f}$  and therefore an additional vector would become linearly dependent on the basis vectors.

So let us consider any of the possible feature states  $x_f$ . Then  $x_f^T h(X) = \sum_{x \in X} x_f^T x = \sum_{x \in X, x=x_f} 1 = n_{x_f}$  counts the number of nodes in the set  $S$  that have features  $x_f$ . Since this holds for all possible feature vectors  $x_f$ , all information about any multiset  $X$  is preserved by  $h(X)$ .

Accordingly, we can write any multiset function  $f(X) = (f(n_{x_1}), \dots, f(n_{x_n}))$  (which transforms the feature counts) into a function  $g$  that extracts first the count information from the sum  $h(X)$ . Concretely, we can define:  $g(h(X))_f := f(n_{x_f}) = f(x_f^T h(X))$ .  $\square$

## A.3 LOSS OF INFORMATION OVER SECOND HOPS

We now explore an example of a computational tree of a node with two rounds of sum aggregation, and the qualitative kind of information that is lost from the first to the second hop. As shown by Xu et al. (2019) and generalized in Thm. 1, sum is injective over one-hot encoded features, but the second aggregation round sums features that are not necessarily orthogonal, and therefore loses neighborhood information. The computational tree and calculation of hops are displayed in Figure 4.



Feature encodings:  $\square = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  and  $\circ = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$

$$H_1^{(0)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$H_2^{(0)} = H_3^{(0)} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \implies H_1^{(1)} = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$$

$$H_2^{(1)} = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \text{ and } H_3^{(1)} = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \implies H_1^{(2)} = \begin{pmatrix} 5 \\ 2 \end{pmatrix}$$

Figure 4: Example of a two-hop neighborhood with one-hot encoded features and sum aggregation.

Given the previous hops  $H_1^{(0)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  and  $H_1^{(1)} = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$ , and the second hop  $H_1^{(2)} = \begin{pmatrix} 5 \\ 2 \end{pmatrix}$ , what other combinations of two-hop neighborhoods can there be for node 1? Apart from itself, node 1's two-hop neighbors are in a  $\begin{pmatrix} 3 \\ 2 \end{pmatrix}$  feature ratio. However, we have lost the ability to recognize a) how many belong to each of its one-hop neighbor; and b) the distribution or homogeneity of each neighborhood. In reality, these 5 nodes are approximately spread out in number and distribution across the one-hop neighbors 2 and 3. But alternatively, all 5 nodes could have belonged to node 2, or all squared nodes could have belonged to node 3.

Note that, without previous hops  $H_1^{(0)}$  and  $H_1^{(1)}$ , we cannot even distinguish node 1's original features, nor distinguish its presence as a neighbor in its one-hop neighbors. Therefore, concatenating all hops is advantageous.



For completeness, we also include the calculations for mean aggregation, which are qualitatively similar to the sum in this case.

$$\begin{aligned} H_1^{(0)} &= \binom{1}{0}; & H_2^{(0)} &= H_3^{(0)} = \binom{0}{1} & \implies & H_1^{(1,m)} = \binom{0}{1} \\ H_2^{(1,m)} &= \binom{2/3}{1/3} \text{ and } H_3^{(1,m)} = \binom{3/4}{1/4} & \implies & H_1^{(2,m)} = \binom{17/24}{7/24} \end{aligned}$$

A good lossless aggregation scheme should take all possible second-hop neighborhood distributions and map them to values that would not lose information when aggregated. For instance, choosing  $a, b$  such that  $\binom{2}{1} \rightarrow a$  from node 2, and  $\binom{3}{1} \rightarrow b$  from node 3, so that,  $H_1^{(2)} = a + b$  could recover both values separately. Naturally, mapping them to one-hot encodings per distribution would suffice, but it would grow exponentially. This opens the door for better suitable fixed aggregations, or perhaps other kinds of learnable aggregation beyond current understanding of message passing.

## B EXPERIMENTAL DETAILS

### B.1 DATASET DETAILS

Datasets are taken directly from the setup of Luo et al. (2024), which includes varied node classification datasets. Here in Table 3 we include for completeness the same overview of these benchmarks.

Table 3: Details of the node classification datasets.

Dataset	Type	# Nodes	# Edges	# Features	Classes	Metric	Origin
Cora	Homophily	2,708	5,278	1,433	7	Accuracy	(McCallum et al., 2000)
CiteSeer	Homophily	3,327	4,522	3,703	6	Accuracy	(Sen et al., 2008)
PubMed	Homophily	19,717	44,324	500	3	Accuracy	(Namata et al., 2012)
Computer	Homophily	13,752	245,861	767	10	Accuracy	(Shchur et al., 2018)
Photo	Homophily	7,650	119,081	745	8	Accuracy	(Shchur et al., 2018)
CS	Homophily	18,333	81,894	6,805	15	Accuracy	(Shchur et al., 2018)
Physics	Homophily	34,493	247,962	8,415	5	Accuracy	(Shchur et al., 2018)
WikiCS	Homophily	11,701	216,123	300	10	Accuracy	(Mernyei & Cangea, 2020)
Squirrel	Heterophily	2,223	46,998	2,089	5	Accuracy	(Rozemberczki et al., 2021)
Chameleon	Heterophily	890	8,854	2,325	5	Accuracy	(Rozemberczki et al., 2021)
Roman-Empire	Heterophily	22,662	32,927	300	18	Accuracy	(Platonov et al., 2023b)
Amazon-Ratings	Heterophily	24,492	93,050	300	5	Accuracy	(Platonov et al., 2023b)
Minesweeper	Heterophily	10,000	39,402	7	2	ROC-AUC	(Platonov et al., 2023b)
Questions	Heterophily	48,921	153,540	301	2	ROC-AUC	(Platonov et al., 2023b)

### B.2 HYPERPARAMETERS

Each experiment is run on an NVIDIA A100 GPU. The setup is taken from Luo et al. (2024). That is, for a maximum of 2500 epochs, we tune the following parameters:

1. DROPOUT  $\in (0.0 \ 0.2 \ 0.3 \ 0.5 \ 0.7)$
2. LR  $\in (0.01 \ 0.005 \ 0.001 \ 0.0001)$
3. NORMALIZATION  $\in (\ln \ \text{bn} \ \text{none})$
4. HIDDEN CHANNELS  $\in (64 \ 256 \ 512)$

While Luo et al. (2024) includes weight decay as a hyperparameter, there are no concrete ranges specified for it. Therefore, we tune all the different values from the best runs of the given datasets:

5. WEIGHT DECAY  $\in (0.0 \ 1\text{e-}2 \ 1\text{e-}3 \ 5\text{e-}4 \ 5\text{e-}5)$ .

Moreover, Luo et al. (2024) tunes the local layers from 1 to 10 or 15. We instead take for each dataset the same value that they have found best for the GNNs, and use it to construct our fixed aggregation features up to that depth. In some cases where there are too many features, we restrict the depth to a smaller value, thus including a strict subset of features instead. This serves as an ad hoc feature selection to reduce overfitting.

We include as hyperparameter the MLP depth. We also include results for MLP = 1 in Table 7.

Table 4: Best hyperparameters for FAF. Classic GNNs are taken from (Luo et al., 2024).

Dataset	Model	dropout	lr	bn	ln	hidden channels	wd	hops	mlp layers	res
computer	GCN	0.5	0.001	0	1	512	5e-05	3	0	0
	GAT	0.5	0.001	0	1	64	5e-05	2	0	0
	SAGE	0.3	0.001	0	1	64	5e-05	4	0	0
	FAF	0.7	0.005	1	0	256	5e-05	2	2	0
photo	GCN	0.5	0.001	0	1	256	5e-05	6	0	1
	GAT	0.5	0.001	0	1	64	5e-05	3	0	1
	SAGE	0.2	0.001	0	1	64	5e-05	6	0	1
	FAF	0.5	0.005	1	0	256	0.0005	4	2	0
ratings	GCN	0.5	0.001	1	0	512	0	4	0	1
	GAT	0.5	0.001	1	0	512	0	4	0	1
	SAGE	0.5	0.001	1	0	512	0	9	0	1
	FAF	0.2	0.001	1	0	256	0	3	2	0
chameleon	GCN	0.2	0.005	0	0	512	0.001	5	0	0
	GAT	0.7	0.01	1	0	256	0.001	2	0	1
	SAGE	0.7	0.01	1	0	256	0.001	4	0	1
	FAF	0.3	0.001	1	0	512	0.01	5	5	0
citeseer	GCN	0.5	0.001	0	0	512	0.01	2	0	0
	GAT	0.5	0.001	0	0	256	0.01	3	0	1
	SAGE	0.2	0.001	0	0	512	0.01	3	0	0
	FAF	0	0.005	0	1	512	0.001	2	3	0
coauthor-cs	GCN	0.3	0.001	0	1	512	0.0005	2	0	1
	GAT	0.3	0.001	0	1	256	0.0005	1	0	1
	SAGE	0.5	0.001	0	1	512	0.0005	2	0	1
	FAF	0.2	0.005	1	0	64	0.01	2	2	0
coauthor-physics	GCN	0.3	0.001	0	1	64	0.0005	2	0	1
	GAT	0.7	0.001	1	0	256	0.0005	2	0	1
	SAGE	0.7	0.001	1	0	64	0.0005	2	0	1
	FAF	0	0.001	1	0	512	0.001	1	2	0
cora	GCN	0.7	0.001	0	0	512	0.0005	3	0	0
	GAT	0.2	0.001	0	0	512	0.0005	3	0	1
	SAGE	0.7	0.001	0	0	256	0.0005	3	0	0
	FAF	0.7	0.01	0	1	512	0.01	3	3	0
minesweeper	GCN	0.2	0.01	1	0	64	0	12	0	1
	GAT	0.2	0.01	1	0	64	0	15	0	1
	SAGE	0.2	0.01	1	0	64	0	15	0	1
	FAF	0.2	0.01	1	0	64	0	4	12	0
pubmed	GCN	0.7	0.005	0	0	256	0.0005	2	0	0
	GAT	0.5	0.01	0	0	512	0.0005	2	0	0
	SAGE	0.7	0.005	0	0	512	0.0005	4	0	0
	FAF	0.7	0.01	0	1	64	0	4	2	0
questions	GCN	0.3	3e-05	0	0	512	0	10	0	1
	GAT	0.2	3e-05	0	1	512	0	3	0	1
	SAGE	0.2	3e-05	0	1	512	0	6	0	0
	FAF	0.2	0.005	1	0	512	0.01	4	3	0
roman-empire	GCN	0.5	0.001	1	0	512	0	9	0	1
	GAT	0.3	0.001	1	0	512	0	10	0	1
	SAGE	0.3	0.001	1	0	256	0	9	0	0
	FAF	0.7	0.01	1	0	256	0	2	3	0
squirrel	GCN	0.7	0.01	1	0	256	0.0005	4	0	1
	GAT	0.5	0.005	1	0	512	0.0005	7	0	1
	SAGE	0.7	0.01	1	0	256	0.0005	3	0	1
	FAF	0.7	0.01	1	0	512	0.01	4	5	0
wikics	GCN	0.5	0.001	0	1	256	0	3	0	0
	GAT	0.7	0.001	0	1	512	0	2	0	1
	SAGE	0.7	0.001	0	1	256	0	2	0	0
	FAF	0.7	0.01	1	0	64	0.001	2	2	0

6. MLP LAYERS  $\in (2\ 3\ 5)$ .

On the other hand, we do not include linear residual connections, as these are used to bypass the convolutional layers in the classical GNNs. This creates a direct difference on the two datasets that most benefit from this component, Minesweeper and Roman-Empire.

Table 5: Empirical training time in seconds of FAF and GNN models, averaged over runs.

Dataset	GCN	GAT	SAGE	FAF <sub>4</sub>	FAF <sub>2</sub>	FAF <sub>1</sub>	MLP
computer	127.30 ± 0.58	29.33 ± 3.22	45.33 ± 0.58	42.67 ± 0.58	25.00 ± 0.00	17.00 ± 0.00	9.33 ± 0.58
photo	82.00 ± 0.00	26.00 ± 0.00	33.33 ± 0.58	66.67 ± 0.58	36.00 ± 0.00	21.33 ± 0.58	7.33 ± 0.58
ratings	140.30 ± 0.58	161.30 ± 0.58	330.00 ± 0.00	46.00 ± 0.00	26.00 ± 0.00	17.33 ± 0.58	10.00 ± 0.00
chameleon	22.30 ± 0.48	15.20 ± 0.42	18.00 ± 0.00	44.50 ± 0.97	22.50 ± 0.53	15.10 ± 0.32	10.00 ± 0.00
citeseer	16.20 ± 0.45	20.00 ± 0.00	26.00 ± 0.00	62.80 ± 1.79	36.20 ± 0.45	25.80 ± 1.79	13.60 ± 0.55
coauthor-cs	157.00 ± 0.00	70.33 ± 0.58	301.30 ± 20.50	296.30 ± 1.16	159.00 ± 0.00	91.00 ± 0.00	26.33 ± 0.58
coauthor-physics	65.33 ± 0.58	190.00 ± 0.00	400.30 ± 0.58	2383.00 ± 0.58	1004.00 ± 1.00	532.70 ± 0.58	183.00 ± 0.00
cora	16.40 ± 0.89	20.20 ± 0.45	11.20 ± 0.45	46.80 ± 0.45	26.00 ± 0.00	17.00 ± 0.00	8.20 ± 0.45
minesweeper	69.67 ± 0.58	100.30 ± 0.58	68.67 ± 2.08	19.00 ± 0.00	20.00 ± 3.46	19.00 ± 0.00	18.00 ± 1.73
pubmed	20.20 ± 0.45	42.00 ± 0.00	78.20 ± 0.45	34.00 ± 0.00	19.20 ± 0.45	12.00 ± 0.00	6.20 ± 0.45
questions	650.00 ± 0.00	258.00 ± 1.00	363.70 ± 0.58	180.70 ± 0.58	122.70 ± 0.58	94.00 ± 0.00	64.67 ± 1.16
roman-empire	240.30 ± 0.58	294.30 ± 0.58	93.00 ± 0.00	38.67 ± 2.89	23.33 ± 0.58	18.00 ± 0.00	14.67 ± 2.89
squirrel	29.00 ± 0.00	87.10 ± 0.32	24.70 ± 1.89	43.30 ± 0.48	25.20 ± 0.42	18.00 ± 0.00	11.70 ± 0.48
wikies	60.00 ± 3.46	97.33 ± 0.58	27.33 ± 0.58	10.00 ± 0.00	7.33 ± 0.58	7.00 ± 0.00	7.00 ± 0.00

In Table 4 we include the best hyperparameter choices for the four models: GCN, GATv2, GraphSAGE, and FAF<sub>4</sub>, the results of which are in Table 2. We run baselines directly from the setup of Luo et al. (2024), and we sweep FAF<sub>4</sub> with the same ranges included in their work. Each dataset has a specific number of splits given by their setup (from 3 to 10), which we then average.

We also include in Table 5 the training runtime of our algorithms, including FAF<sub>4</sub> and its variants, grouped by the number of reducers—as MLPs with the same input width will have the same training time. Note that all runs have the same number of epochs (2500), as in the original setup, and all datasets match the number of runs of the setup. In general, MLPs are more efficient than MPGNNs, as backpropagation over message-passing is costly. However, we increase the number of features in the data—depending on the aggregation depth and number of reducers—so for some datasets with many features the improvement is not necessarily observed. A way to reduce this overhead may be to apply dimensionality reduction to the tabular FAF representation.

### B.3 FAFs BEYOND NODE CLASSIFICATION

Our theory applies to any task that learns multiset functions over neighborhoods. In our experiments, we focus on node classification for two main reasons. First, these are the benchmarks on which GNNs have been shown to be competitive with more complex architectures in Luo et al. (2024), so they are amenable to simple models for which we have strong, well-tuned hyperparameters. Second, node-classification datasets typically provide rich features that depend on neighborhood distributions. Thus, non-injective but commonly used reducers such as mean and sum still convey highly informative distributional signals. Regarding inductive settings, they would require computing the new aggregation rounds at test time. We would not have access to the test node features when precomputing training aggregations. Otherwise, our approach is just as feasible as in the transductive case.

## C ADDITIONAL EXPERIMENTS

In Table 6 we include results for different number of hops concatenated as features for FAF<sub>4</sub>. Note that one hop already gives much of the information, and two hops often give the best performance.

In Table 7 we show the performance of two classifiers on the same FAF<sub>4</sub> features: one linear layer and a multilayer perceptron—this being our choice for other experiments.

In Table 8 we ablate on using only the last hop as features to an MLP, or using the last hop concatenated to the original features. This mimics the choice of directly freezing a GNN and using its output as features to an (often linear) classifier.

Table 9 shows the last epoch training accuracy and best epoch validation accuracy of using the Kolmogorov-Arnold function  $\Phi$  as a reducer for FAF. Following Corso et al. (2020), we make it act on multisets by sorting, which we fix by the given data order.

In Table 10 we include results on rewiring the input graph by deleting edges based on pairwise cosine similarity. REW includes hop-wise features where negative similarity neighbors are set to

0. SP includes hop-wise features where positive and negative similarity neighbors are aggregated in different features and concatenated together.

Table 6: Increasing number of concatenated hops in FAF<sub>4</sub>, compared to the best GNN (Classic).

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
Classic	92.86 ± 0.06	95.93 ± 0.15	55.90 ± 0.54	48.15 ± 2.35	70.36 ± 0.09	95.51 ± 0.04	97.16 ± 0.07
FAF+0	87.89 ± 0.13	93.33 ± 0.07	48.98 ± 0.72	41.43 ± 1.77	53.80 ± 0.82	93.70 ± 0.07	95.89 ± 0.02
FAF+1	92.53 ± 0.08	96.14 ± 0.07	54.17 ± 0.14	46.91 ± 1.43	65.52 ± 0.64	94.84 ± 0.08	96.83 ± 0.01
FAF+2	93.05 ± 0.04	96.23 ± 0.08	55.02 ± 0.67	47.30 ± 1.76	67.28 ± 0.64	94.93 ± 0.07	96.63 ± 0.02
FAF+4	93.04 ± 0.10	96.34 ± 0.07	55.08 ± 0.19	48.42 ± 2.22	50.52 ± 3.84	94.88 ± 0.07	96.63 ± 0.04
FAF+8	92.96 ± 0.13	96.21 ± 0.07	55.26 ± 0.30	48.74 ± 1.69	40.72 ± 2.26	94.92 ± 0.09	-

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
Classic	81.32 ± 0.41	97.68 ± 0.63	79.08 ± 0.23	78.63 ± 0.23	91.14 ± 0.58	44.88 ± 1.27	82.38 ± 0.57
FAF+0	62.68 ± 1.15	51.12 ± 0.93	71.12 ± 0.52	71.58 ± 1.46	66.28 ± 0.27	40.57 ± 0.92	74.86 ± 0.33
FAF+1	82.16 ± 0.33	87.65 ± 0.47	78.24 ± 0.36	77.44 ± 1.07	77.36 ± 0.55	47.26 ± 1.31	81.37 ± 0.51
FAF+2	82.84 ± 0.38	89.48 ± 1.08	78.52 ± 0.41	79.71 ± 0.86	78.68 ± 0.19	47.18 ± 1.44	81.92 ± 0.43
FAF+4	81.80 ± 0.35	89.63 ± 1.02	79.08 ± 0.36	79.67 ± 0.89	77.48 ± 0.17	47.61 ± 1.43	81.73 ± 0.53
FAF+8	74.28 ± 0.23	89.10 ± 1.03	76.80 ± 0.20	79.94 ± 0.88	75.08 ± 0.23	47.95 ± 1.36	81.58 ± 0.63

Table 7: Comparison of 1 linear layer (1L) versus multiple layers (MLP) as the classifier over FAF<sub>4</sub>.

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
FAF+MLP	93.05 ± 0.04	96.34 ± 0.07	55.53 ± 0.43	48.51 ± 2.31	67.28 ± 0.64	94.93 ± 0.07	96.83 ± 0.01
FAF+1L	91.50 ± 0.08	96.01 ± 0.00	47.65 ± 0.75	47.11 ± 2.68	66.76 ± 0.99	93.32 ± 0.05	96.62 ± 0.01

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
FAF+MLP	82.84 ± 0.43	89.63 ± 1.03	79.08 ± 0.36	79.53 ± 1.12	78.68 ± 0.19	47.31 ± 1.39	81.92 ± 0.43
FAF+1L	81.08 ± 0.41	88.99 ± 0.76	78.00 ± 0.51	77.22 ± 1.04	75.56 ± 0.28	45.89 ± 1.75	81.25 ± 0.96

Table 8: Using only the last hop ( $H^{(K)}$ ), that and original features ( $H^{(0) \oplus (K)}$ ), and all hops (FAF<sub>4</sub>).

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
FAF <sub>4</sub>	93.05 ± 0.04	96.34 ± 0.07	55.53 ± 0.43	48.51 ± 2.31	67.28 ± 0.64	94.93 ± 0.07	96.83 ± 0.01
$H^{(0) \oplus (K)}$	92.35 ± 0.08	95.53 ± 0.04	54.95 ± 0.23	48.16 ± 2.61	66.80 ± 0.42	95.13 ± 0.02	96.65 ± 0.01
$H^{(K)}$	91.67 ± 0.04	92.18 ± 0.04	49.64 ± 0.40	49.18 ± 2.14	66.68 ± 0.46	92.83 ± 0.10	96.22 ± 0.06

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
FAF <sub>4</sub>	82.84 ± 0.43	89.63 ± 1.03	79.08 ± 0.36	79.53 ± 1.12	78.68 ± 0.19	47.31 ± 1.39	81.92 ± 0.43
$H^L$	82.12 ± 0.52	69.40 ± 1.50	79.24 ± 0.36	78.14 ± 0.93	55.43 ± 0.51	46.10 ± 1.46	80.14 ± 0.67
$H^{(0) \oplus (K)}$	81.80 ± 0.28	73.30 ± 1.26	79.64 ± 0.38	79.29 ± 0.27	76.42 ± 0.13	46.41 ± 1.84	81.22 ± 0.52
$H^{(K)}$	82.12 ± 0.52	69.40 ± 1.50	79.24 ± 0.36	78.14 ± 0.93	55.43 ± 0.51	46.10 ± 1.46	80.14 ± 0.67

Table 9: (Last epoch) training and (best) validation accuracy of the KA function  $\Phi$  from Thm. 2.

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
FAF <sub>KA</sub> (train)	94.73 ± 0.12	99.65 ± 0.09	99.92 ± 0.02	96.35 ± 3.44	70.17 ± 41.23	99.49 ± 0.22	99.84 ± 0.28
FAF <sub>KA</sub> (val)	87.88 ± 0.19	93.40 ± 0.07	51.97 ± 0.12	41.65 ± 1.91	55.76 ± 1.07	93.66 ± 0.08	95.90 ± 0.05

Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikics
FAF <sub>KA</sub> (train)	14.29 ± 0.00	51.68 ± 0.44	33.33 ± 0.00	99.85 ± 0.09	86.41 ± 0.23	39.30 ± 2.30	97.64 ± 1.09
FAF <sub>KA</sub> (val)	29.56 ± 0.79	51.61 ± 0.42	42.80 ± 0.76	74.28 ± 1.95	80.45 ± 0.25	40.43 ± 1.10	77.38 ± 0.88

Table 10: Feature augmentations based on similarity-based rewiring (REW) and computational graph splitting (SP) on a subset of datasets that benefit from mean aggregation.

Dataset	computer	photo	chameleon	citeseer	cora	pubmed	wikis
FAF <sub>mean</sub>	93.16 $\pm$ 0.04	<b>96.06 <math>\pm</math> 0.10</b>	47.99 $\pm$ 2.02	66.92 $\pm$ 0.87	83.28 $\pm$ 0.30	81.16 $\pm$ 0.97	81.58 $\pm$ 0.46
REW <sub>mean</sub>	93.25 $\pm$ 0.08	95.90 $\pm$ 0.04	43.94 $\pm$ 2.39	66.92 $\pm$ 0.78	82.36 $\pm$ 0.17	80.80 $\pm$ 0.42	82.44 $\pm$ 0.56
SP <sub>mean</sub>	93.20 $\pm$ 0.11	95.97 $\pm$ 0.10	43.71 $\pm$ 1.71	67.32 $\pm$ 0.99	81.80 $\pm$ 0.24	80.84 $\pm$ 0.62	82.38 $\pm$ 0.46
FAF <sub>mean</sub> +REW <sub>mean</sub>	<b>93.33 <math>\pm</math> 0.12</b>	96.03 $\pm$ 0.04	<b>48.26 <math>\pm</math> 1.66</b>	67.48 $\pm$ 0.36	83.20 $\pm$ 0.20	80.84 $\pm$ 0.52	<b>82.46 <math>\pm</math> 0.48</b>
FAF <sub>mean</sub> +SP <sub>mean</sub>	93.31 $\pm$ 0.10	<b>96.06 <math>\pm</math> 0.04</b>	47.99 $\pm$ 1.83	<b>67.88 <math>\pm</math> 0.33</b>	<b>83.52 <math>\pm</math> 0.36</b>	<b>81.24 <math>\pm</math> 0.43</b>	82.43 $\pm$ 0.51

In Figure 5 we show two more plots of feature importance using SHAP (Lundberg & Lee, 2017) for Pubmed and Amazon-Ratings, on the MLP over single-reducer FAFs. Features are sorted by global importance and broken down over the different hops by color. While the implementation of SHAP on MLPs used (GradientExplainer) relies on local linearization and often assumes input feature independence, the explanations still reveal informative qualitative patterns. In Pubmed, feature 346 is most important at hops 1 and 2, and remains second at hops 0 and 4, whereas the most important base feature (205) contributes little at other hops. By contrast, in Amazon-Ratings, importance is more evenly distributed across features and hops.

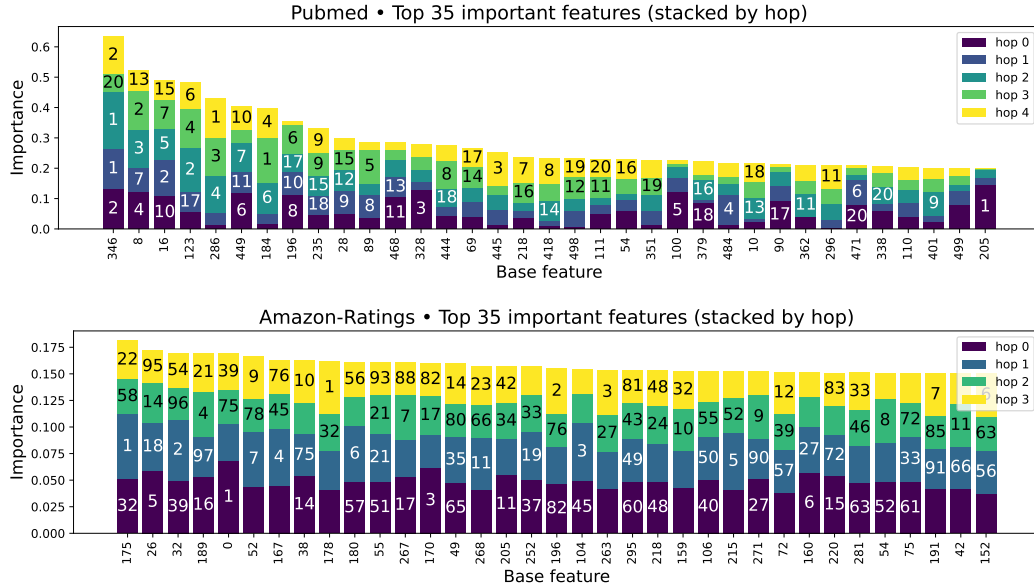


Figure 5: SHAP feature importance for Pubmed and Amazon-Ratings. The base features are ranked according to the sum of their importance values across hops. Numbers on the stacked bars indicate the ranking of that particular feature on that particular hop.



## D TRAINING, VALIDATION, AND TEST ACCURACY CURVES

In this section we compare training MLPs on FAF features to training GCNs, by tracking train/validation/test accuracy over epochs (Figure 6). Below we summarize the behaviors on datasets where differences arise between the two methods:

- Amazon-Computer (6a) and Amazon-Photo (6b) behave similarly, but GCNs are more unstable.
- FAF for Chameleon (6d) has much better training accuracy but similar generalization; in contrast, GCN for Squirrel (6m) has much better training, but slightly worse generalization than FAF.
- Citeseer (6e) with FAF breaks at the end of training, which indicates instability. However, this could be overcome with standard learning rate schedules.
- Coauthor-CS (6f) and Coauthor-Physics (6g) have dips in all metrics for both models.
- Questions (6f) with FAF is more (locally) unstable but also more stationary and does not degrade performance later on.
- As mentioned in the main text, Minesweeper (6i) and Roman-Empire (6l) are the two datasets that seem to truly lose neighborhood information with FAF.

### D.1 TEST ACCURACY OF MAIN RESULTS

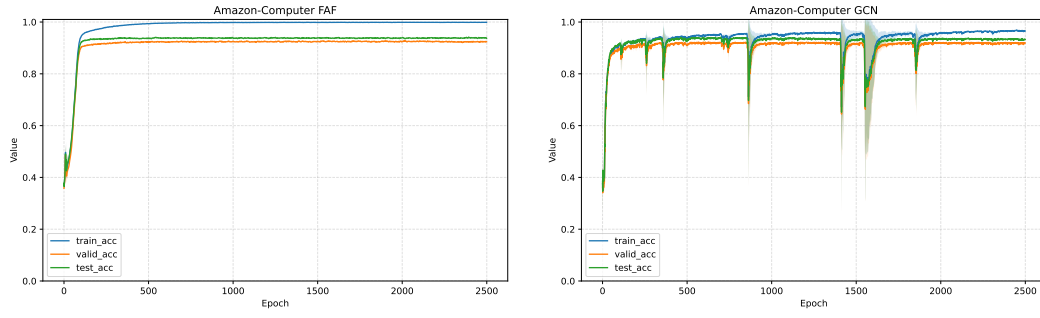
Here in Table 11 we report the test accuracy of the main FAF variants of our experimental results (§ 5) where in Table 1 we only have the best validation FAF’s test results, and in Table 2 we show validation results. FAFs in all datasets are  $\pm 1\%$  away from the best classic GNN, except for those already mentioned in the main text (Citeseer, Cora, Roman-Empire, and Minesweeper).

Table 11: Test accuracy on 14 node classification benchmarks: FAFs+MLP against classic GNNs. Validation accuracy is shown in Table 2.

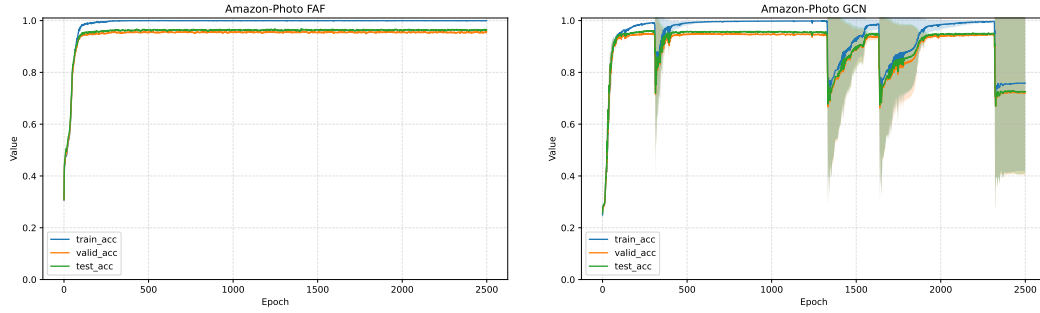
Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs	coauthor-physics
GCN	93.58 $\pm$ 0.44	95.77 $\pm$ 0.27	53.86 $\pm$ 0.48	44.62 $\pm$ 4.50	<b>72.72 <math>\pm</math> 0.45</b>	95.73 $\pm$ 0.15	<b>97.47 <math>\pm</math> 0.08</b>
GAT	93.91 $\pm$ 0.22	96.45 $\pm$ 0.37	<b>55.51 <math>\pm</math> 0.55</b>	42.90 $\pm$ 5.47	71.82 $\pm$ 0.65	96.14 $\pm$ 0.08	97.12 $\pm$ 0.13
SAGE	93.31 $\pm$ 0.17	96.17 $\pm$ 0.44	55.26 $\pm$ 0.27	43.11 $\pm$ 4.73	71.82 $\pm$ 0.81	<b>96.21 <math>\pm</math> 0.10</b>	97.10 $\pm$ 0.09
MLP	87.75 $\pm$ 0.42	93.62 $\pm$ 0.36	49.04 $\pm$ 0.39	38.59 $\pm$ 3.29	57.22 $\pm$ 2.25	93.80 $\pm$ 0.19	96.02 $\pm$ 0.16
FAF <sub>bestval</sub>	<b>94.01 <math>\pm</math> 0.21</b>	96.54 $\pm$ 0.13	55.09 $\pm$ 0.24	42.96 $\pm$ 2.45	70.48 $\pm$ 1.24	95.37 $\pm$ 0.17	97.05 $\pm$ 0.18
FAF <sub>4</sub>	93.75 $\pm$ 0.04	96.54 $\pm$ 0.13	54.42 $\pm$ 0.45	42.96 $\pm$ 2.45	69.42 $\pm$ 1.32	95.33 $\pm$ 0.20	96.96 $\pm$ 0.09
FAF <sub>mean,std</sub>	94.00 $\pm$ 0.25	96.30 $\pm$ 0.23	54.73 $\pm$ 0.22	<b>45.13 <math>\pm</math> 3.42</b>	67.90 $\pm$ 0.95	95.34 $\pm$ 0.14	96.93 $\pm$ 0.04
FAF <sub>mean</sub>	<b>94.01 <math>\pm</math> 0.21</b>	<b>96.71 <math>\pm</math> 0.16</b>	53.12 $\pm$ 0.44	43.21 $\pm$ 2.24	66.82 $\pm$ 1.74	95.37 $\pm$ 0.17	97.05 $\pm$ 0.18
FAF <sub>max,std</sub>	93.60 $\pm$ 0.25	96.01 $\pm$ 0.41	55.09 $\pm$ 0.24	43.20 $\pm$ 2.42	67.18 $\pm$ 0.88	95.53 $\pm$ 0.10	96.61 $\pm$ 0.04
FAF <sub>max</sub>	92.98 $\pm$ 0.22	96.12 $\pm$ 0.10	54.79 $\pm$ 0.15	42.15 $\pm$ 3.19	67.52 $\pm$ 0.40	95.55 $\pm$ 0.08	96.84 $\pm$ 0.13
FAF <sub>sum</sub>	91.77 $\pm$ 0.24	95.08 $\pm$ 0.61	53.44 $\pm$ 0.18	39.63 $\pm$ 2.90	70.48 $\pm$ 1.24	95.08 $\pm$ 0.12	96.86 $\pm$ 0.06
FAF <sub>std</sub>	93.54 $\pm$ 0.26	96.17 $\pm$ 0.10	54.77 $\pm$ 0.14	42.68 $\pm$ 2.75	62.70 $\pm$ 1.18	95.77 $\pm$ 0.12	96.97 $\pm$ 0.09

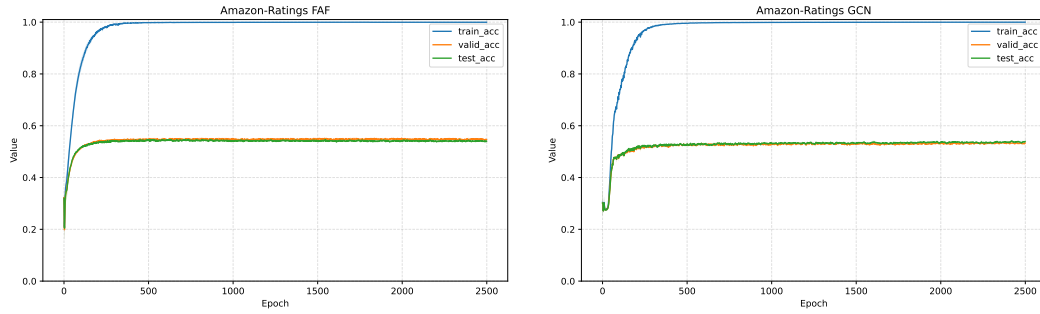
Dataset	cora	minesweeper	pubmed	questions	roman-empire	squirrel	wikis
GCN	<b>84.38 <math>\pm</math> 0.81</b>	97.48 $\pm$ 0.06	80.00 $\pm$ 0.77	78.44 $\pm$ 0.23	<b>91.05 <math>\pm</math> 0.15</b>	44.26 $\pm$ 1.22	80.06 $\pm$ 0.81
GAT	83.02 $\pm$ 1.21	97.00 $\pm$ 1.02	79.80 $\pm$ 0.94	77.72 $\pm$ 0.71	90.38 $\pm$ 0.49	39.31 $\pm$ 2.42	<b>81.01 <math>\pm</math> 0.23</b>
SAGE	83.18 $\pm$ 0.93	<b>97.72 <math>\pm</math> 0.70</b>	77.42 $\pm$ 0.40	76.75 $\pm$ 1.07	90.41 $\pm$ 0.10	40.22 $\pm$ 1.47	80.57 $\pm$ 0.42
MLP	58.56 $\pm$ 1.75	51.74 $\pm$ 0.83	68.22 $\pm$ 0.96	70.40 $\pm$ 1.17	66.43 $\pm$ 0.12	39.11 $\pm$ 1.93	72.98 $\pm$ 0.49
FAF <sub>bestval</sub>	82.84 $\pm$ 0.63	90.00 $\pm$ 0.39	<b>80.96 <math>\pm</math> 1.06</b>	<b>78.69 <math>\pm</math> 0.50</b>	78.11 $\pm$ 0.38	<b>44.59 <math>\pm</math> 1.62</b>	80.25 $\pm$ 0.34
FAF <sub>4</sub>	81.44 $\pm$ 0.38	90.01 $\pm$ 0.51	77.20 $\pm$ 0.45	<b>78.69 <math>\pm</math> 0.50</b>	78.11 $\pm$ 0.38	44.02 $\pm$ 2.08	80.25 $\pm$ 0.34
FAF <sub>mean,std</sub>	82.84 $\pm$ 0.63	90.17 $\pm$ 0.51	<b>80.96 <math>\pm</math> 1.06</b>	75.82 $\pm$ 1.27	77.14 $\pm$ 0.52	43.83 $\pm$ 2.34	79.48 $\pm$ 0.81
FAF <sub>mean</sub>	82.80 $\pm$ 0.70	90.00 $\pm$ 0.39	79.88 $\pm$ 0.92	76.83 $\pm$ 1.19	76.36 $\pm$ 0.55	42.44 $\pm$ 1.73	79.61 $\pm$ 0.56
FAF <sub>max,std</sub>	79.34 $\pm$ 0.95	88.36 $\pm$ 0.74	77.52 $\pm$ 0.77	76.62 $\pm$ 0.79	75.89 $\pm$ 0.30	<b>44.59 <math>\pm</math> 1.62</b>	78.44 $\pm$ 0.67
FAF <sub>max</sub>	79.34 $\pm$ 0.67	86.39 $\pm$ 1.22	77.18 $\pm$ 0.13	77.59 $\pm$ 1.67	75.01 $\pm$ 0.43	43.03 $\pm$ 1.90	78.63 $\pm$ 0.35
FAF <sub>sum</sub>	81.46 $\pm$ 0.62	89.96 $\pm$ 0.45	77.46 $\pm$ 0.43	76.12 $\pm$ 1.08	76.90 $\pm$ 0.28	44.07 $\pm$ 1.98	76.59 $\pm$ 0.36
FAF <sub>std</sub>	79.50 $\pm$ 0.39	88.93 $\pm$ 0.68	79.06 $\pm$ 1.09	73.99 $\pm$ 1.67	73.80 $\pm$ 0.21	43.63 $\pm$ 1.43	76.09 $\pm$ 0.26



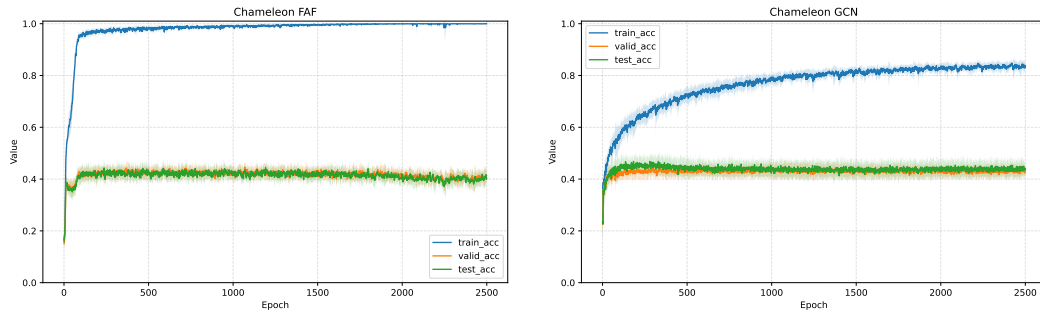
(a) Amazon-Computer: FAF vs. GCN



(b) Amazon-Photo: FAF vs. GCN

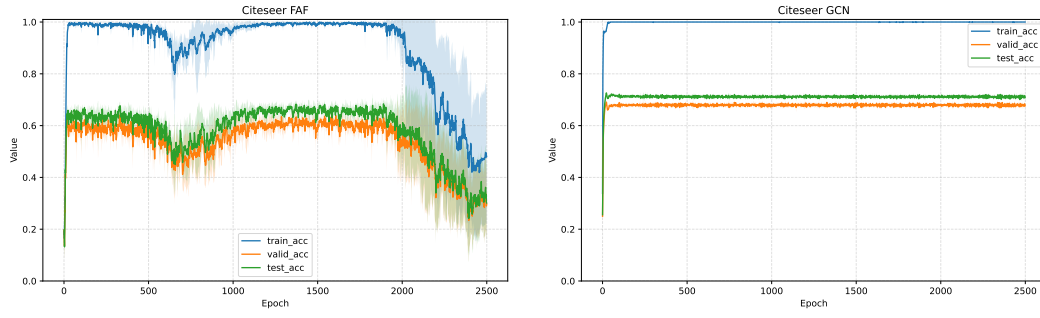


(c) Amazon-Ratings: FAF vs. GCN

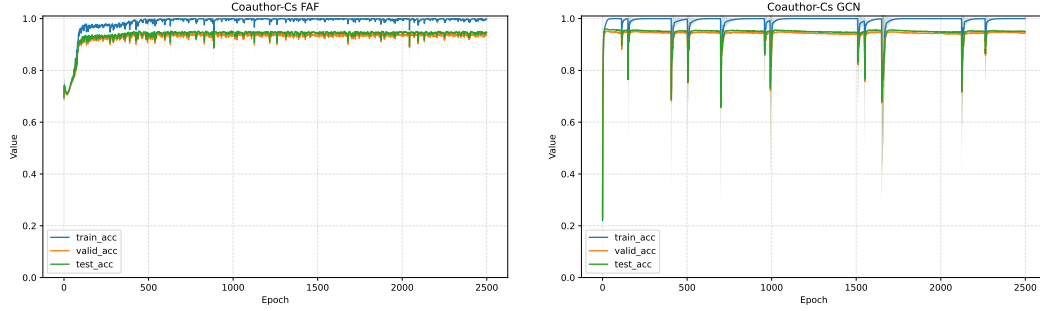


(d) Chameleon: FAF vs. GCN

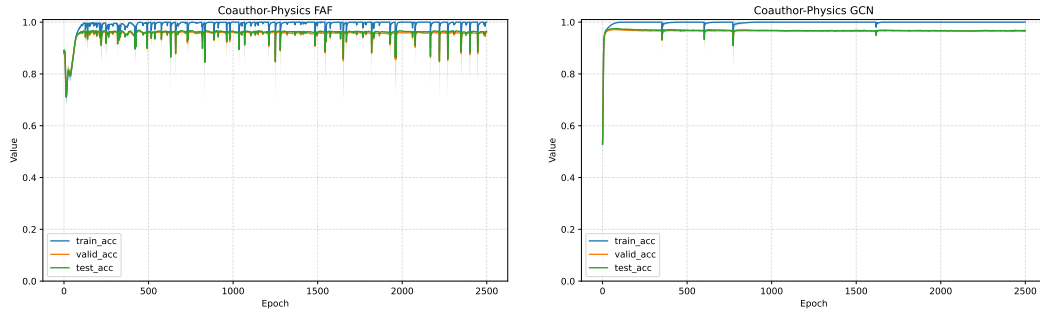
Figure 6: Train, validation, and test accuracy of FAF+MLP versus GCN. (i)



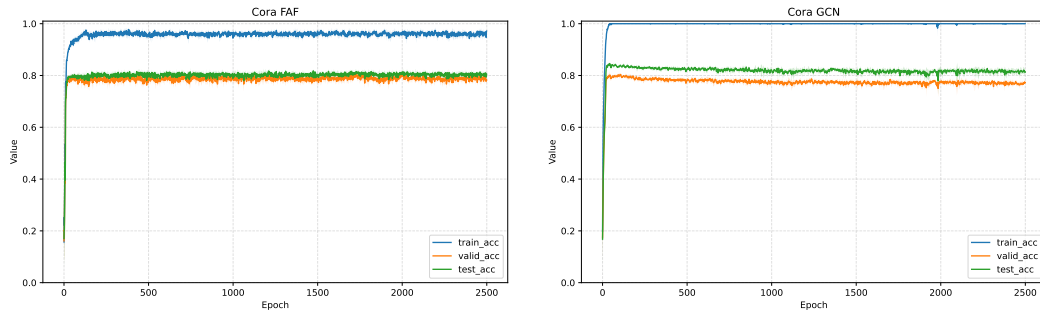
(e) Citeseer: FAF vs. GCN



(f) Coauthor-CS: FAF vs. GCN

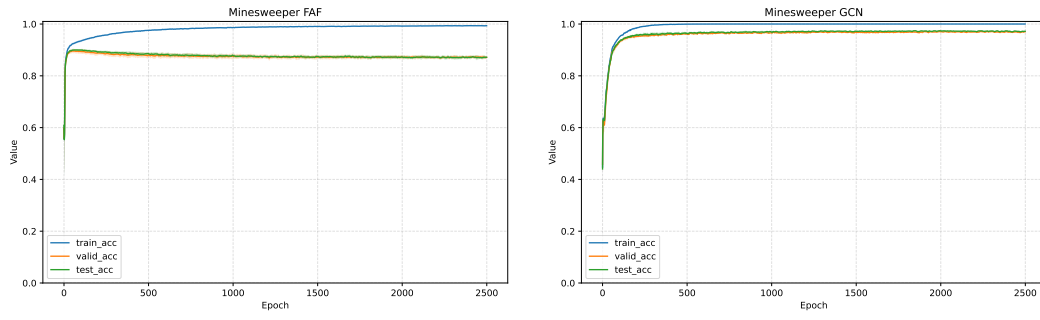


(g) Coauthor-Physics: FAF vs. GCN

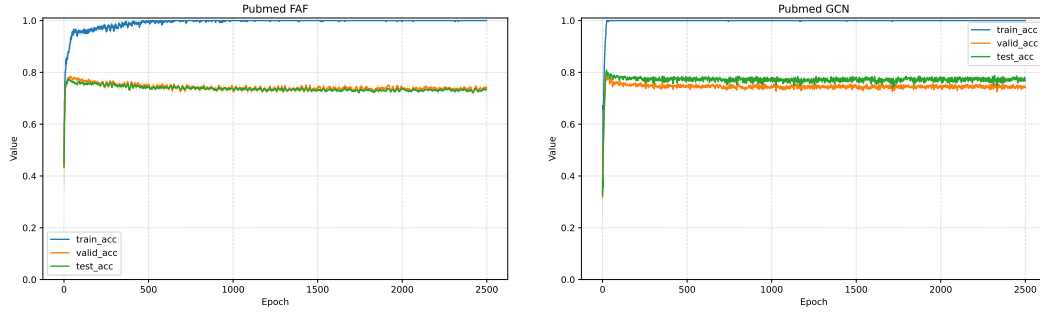


(h) Cora: FAF vs. GCN

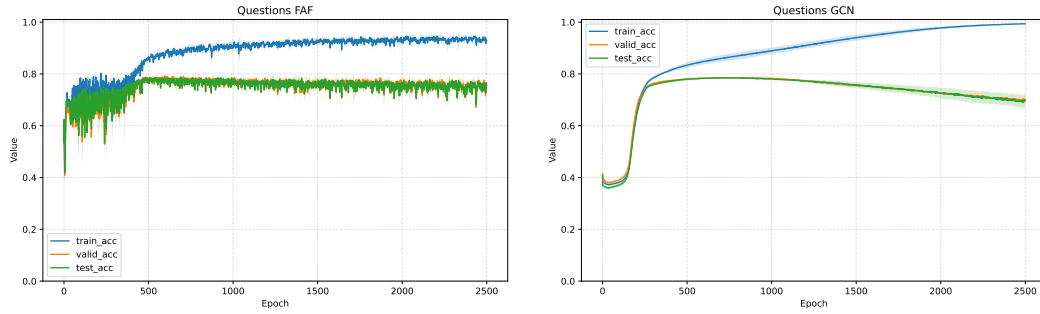
Figure 6: Train, validation, and test accuracy of FAF+MLP versus GCN. (ii)



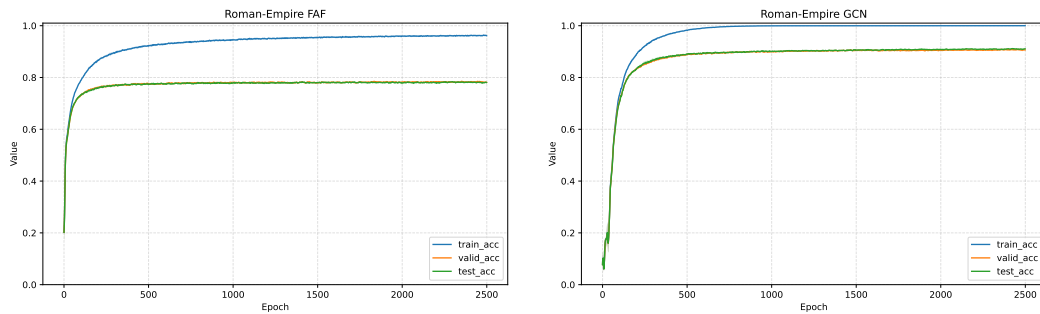
(i) Minesweeper: FAF vs. GCN



(j) Pubmed: FAF vs. GCN

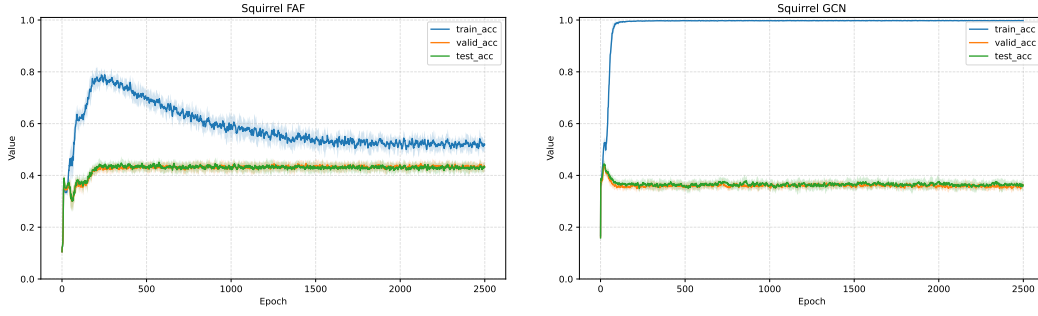


(k) Questions: FAF vs. GCN

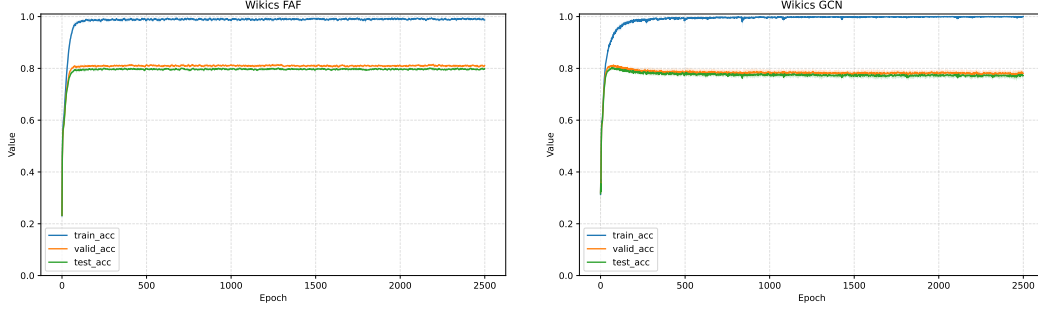


(l) Roman-Empire: FAF vs. GCN

Figure 6: Train, validation, and test accuracy of FAF+MLP versus GCN. (iii)



(m) Squirrel: FAF vs. GCN



(n) Wikics: FAF vs. GCN

Figure 6: Train, validation, and test accuracy of FAF+MLP versus GCN. (iv)

## E PRELIMINARY RESULTS ON OTHER BENCHMARKS

We include preliminary results for FAFs on the GraphLand benchmark (Bazhenov et al., 2025). We do *not* perform the full hyperparameter sweep, therefore we indicate FAFs with an asterisk (\*), as there could be better performing versions. We copy baselines from the original paper: MLP, MLP-NFA (one-hop FAFs), GCN and GAT. For FAFs, we only report results for mean+std and mean aggregations. We choose the best validation hyperparameters we have been able to find so far (shown in Table 13) and report the resulting test accuracy in Table 12. While we do not yet achieve the performance of GATs, we approach that of GCNs, and we improve upon MLPs and NFA.

Table 12: Test accuracy of 4 GraphLand datasets (averaged over 10 runs).

	artnet-exp	hm-categories	tolokers-2	pokec-regions
ResMLP	$35.07 \pm 2.34$	$37.72 \pm 0.18$	$41.16 \pm 1.13$	$4.88 \pm 0.01$
ResMLP-NFA	$38.25 \pm 0.56$	$48.72 \pm 0.38$	$48.14 \pm 1.40$	$8.05 \pm 0.03$
GCN	$43.09 \pm 0.38$	$61.70 \pm 0.35$	$51.32 \pm 0.96$	$34.96 \pm 0.38$
GAT	$46.62 \pm 0.32$	$67.96 \pm 0.33$	$53.78 \pm 1.34$	$46.17 \pm 0.32$
FAF* <sub>mean,std</sub>	$41.56 \pm 0.26$	$59.50 \pm 0.15$	$52.74 \pm 0.53$	$28.44 \pm 0.22$
FAF* <sub>mean</sub>	$39.25 \pm 0.38$	$54.50 \pm 0.13$	$50.72 \pm 0.38$	$31.23 \pm 0.18$

Table 13: Best hyperparameters found so far for FAFs on GraphLand datasets.

	dropout	lr	bn	hidden channels	weight decay	local layers	mlp layers
artnet-exp	0.7	0.01	1	256	0.01	3	2
hm-categories	0.5	0.001	1	512	0.0005	3	3
tolokers-2	0.3	0.0001	1	512	5e-05	3	5
pokec-regions	0	0.001	1	512	0.001	3	3



## F COMPARISON TO GRAPH ECHO STATE NETWORKS

Graph Echo State Networks (GESN) (Gallicchio & Micheli, 2010) compute label-independent node embeddings via fixed “reservoir” layers, one per hop, followed by a linear readout. We compare FAFs to this approach in Table 14 as a representative previously proposed simplification of GNNs. We do not include the *coauthor-physics* or *questions* datasets, as GESN exceeds memory capacity on them.

We use the public implementation at <https://github.com/dtortorella/graph-esn>, keeping most GESN-specific hyperparameters as in their example. We set the “depth” (local layers) and “hidden units” (hidden channels) as in the best GCN. Each layer is given 10 minutes to compute its embedding.

As the classifier, we replace the original linear readout with a MLP of the same architecture as for FAFs to provide a fair comparison of the role of the embeddings. For the MLP hyperparameters, we try both the best GCN and the best FAF configurations. Since our different FAFs could share MLP hyperparameters, we expected them to transfer well here, but in this case the GCN hyperparameters perform slightly better. As shown below, FAFs seem to be more suitable for the tested benchmark tasks.

Table 14: Test accuracy of GESN (Gallicchio & Micheli, 2010) against FAFs.

Dataset	computer	photo	ratings	chameleon	citeseer	coauthor-cs
FAF <sub>bestval</sub>	94.01 $\pm$ 0.21	96.54 $\pm$ 0.13	55.09 $\pm$ 0.24	42.96 $\pm$ 2.45	70.48 $\pm$ 1.24	95.37 $\pm$ 0.17
GESN+MLP	90.80 $\pm$ 0.10	92.72 $\pm$ 0.27	50.34 $\pm$ 0.26	41.64 $\pm$ 3.63	41.44 $\pm$ 0.34	89.99 $\pm$ 0.09

Dataset	cora	minesweeper	pubmed	roman-empire	squirrel	wikics
FAF <sub>bestval</sub>	82.84 $\pm$ 0.63	90.00 $\pm$ 0.39	80.96 $\pm$ 1.06	78.11 $\pm$ 0.38	44.59 $\pm$ 1.62	80.25 $\pm$ 0.34
GESN+MLP	65.78 $\pm$ 0.26	50.93 $\pm$ 1.28	64.98 $\pm$ 1.57	11.76 $\pm$ 0.38	36.58 $\pm$ 1.18	73.98 $\pm$ 0.85