# WHICH EIGENVECTORS DO GRAPH TRANSFORMERS NEED FOR NODE CLASSIFICATION?

**Anonymous authors**Paper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

031

033

034

036

037

038

040

041

042

043

044

046

047

048

049

051

052

## **ABSTRACT**

Graph transformers have emerged as powerful tools for modeling complex graphstructured data, offering the ability to capture long-range dependencies beyond the graph adjacency. Yet their performance on node classification often lags behind that of message passing and spectral graph networks. Unlike these methods, graph transformers require explicit positional encodings to inject structural information, which are most commonly derived from the eigenvectors of the graph Laplacian. Existing methods select eigenvectors using data-agnostic heuristics, assuming onesize-fits-all rules suffice. In contrast, we show that the spectral distribution of class information is graph-specific. To address this, we introduce *Broaden the Spectrum* (BTS), a novel, intuitive, and data-driven algorithm for selecting subsets of Laplacian eigenvectors for node classification. Our method is grounded in theory: we characterize the structure of optimal attention matrices for classification and show, in a simplified setting, how BTS naturally emerges as the eigenvector selection rule for achieving such attention matrices. When evaluated with standard graph transformer architectures, it delivers substantial performance gains across a wide range of node classification benchmarks. Our work shows that the performance of graph transformers on node classification has been held back by the choice of positional encodings and can be improved by employing a broader, well-chosen set of Laplacian eigenvectors.

# 1 Introduction

Graph transformers provide a flexible framework for modeling graph-structured data, with global receptive fields that can capture interactions beyond the reach of local message passing (Hoang et al., 2024). This flexibility has made them appealing for graph-level tasks such as molecular property prediction and long-range dependency modeling, where message-passing Graph Neural Networks (GNNs) often struggle with phenomena like oversquashing and oversmoothing (Topping et al., 2022; Rusch et al., 2023). Yet their performance on node classification has often lagged behind both message-passing and spectral methods (Luo et al., 2024; Bo et al., 2023).

A central reason for this lack of performance in node-classification tasks lies in how transformers incorporate information about the graph topology. Unlike message-passing networks, which propagate information directly along edges, transformers rely on *positional encodings* (PEs) to inject structural information. In practice, these encodings are commonly derived from the graph's Laplacian eigenvectors (Hoang et al., 2024; Dwivedi & Bresson, 2020), which most models truncate to a few lowest-frequency components (see Table 7 in Appendix B). While this low-frequency truncation heuristic works well for some homophilic graphs, it does not account for the fact that class-relevant information can appear in very different parts of the spectrum depending on the dataset. Heuristics such as fixed low-high splits in Kim et al. (2022) suggest that including a broader spectrum can help, but their effectiveness varies on the nature of the graphs themselves.

Crucially, these heuristics are data-agnostic. In contrast, the spectral GNN literature has demonstrated the advantages of *adaptive* frequency selection, combining high- and low-frequency operations based on the task (Sun et al., 2022; Bo et al., 2021; Dong et al., 2021). This motivates the need for adaptive methods that select graph transformer positional encodings in a task-aware manner, rather than relying on one-size-fits-all rules.

To address this, we introduce *Broaden the Spectrum (BTS)*, a simple, data-driven, and theoretically grounded approach for selecting Laplacian eigenvectors in graph transformers. BTS identifies the parts of the spectrum that are most aligned with class information and uses them as positional encodings. Empirically, BTS yields consistent gains across homophilic, heterophilic, and long-range benchmarks. On challenging heterophilic datasets such as Chameleon and Squirrel, even a simple transformer backbone equipped with BTS improves accuracy by more than 20%. More advanced models such as NAGphormer (Chen et al., 2023) and GraphGPS (Rampášek et al., 2022) also see substantial boosts when augmented with BTS, revealing the existence of performance bottlenecks due to under-utilization of graph topology.

## Our contributions are as follows:

- We introduce Broaden the Spectrum (BTS), a lightweight, architecture-agnostic, and datadriven algorithm for selecting Laplacian eigenvectors as positional encodings specifically for node classification.
- We provide a theoretical analysis of attention-based node classification in an illustrative linear model, showing that the optimal attention matrix has a class-wise block structure. We also derive the eigenvector selection criteria to achieve such optimal attention matrices.
- We demonstrate that including a broad and well-chosen spectrum of eigenvectors leads to significant gains in node classification performance of *existing* graph transformers across a wide range of benchmarks.

# 2 METHOD

Transformers are inherently permutation-invariant, which makes positional information essential to break symmetry and provide meaningful structure to the model. For graphs, Laplacian position encodings (LPE) have been identified by prior work (Dwivedi & Bresson, 2020; Hoang et al., 2024) to be effective, and are a natural extension of the Fourier basis used in other sequence modeling domains (Vaswani et al., 2017; Dosovitskiy et al., 2021; Nie et al., 2023). More formally, we define the LPE as follows.

**Definition 2.1** (Laplacian Position Encodings). Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a graph with  $|\mathcal{V}| = n$  nodes, adjacency matrix A, and degree matrix D. The normalized graph Laplacian is  $L = I - D^{-1/2}AD^{-1/2}$ . Let  $L = V\Lambda V^{\top}$  be its eigendecomposition with eigenvalues ordered as  $\lambda_1 \leq \cdots \leq \lambda_n$ . The Laplacian positional encodings are defined as  $X_{\text{pos}} \in \mathbb{R}^{n \times k}$  formed by selecting any k Laplacian eigenvectors.

# 2.1 CLASS-LABEL ENERGY SPECTRAL DENSITY

When using LPEs, there is a critical design choice to be made: which subset of eigenvectors should be used? The prevailing practice is to use fixed heuristics, such as using the first k eigenvectors (Dwivedi & Bresson, 2020; Chen et al., 2023; Rampášek et al., 2022; Kreuzer et al., 2021; Hoang et al., 2024), or using an equal number of low- and high-frequency<sup>1</sup> components (Kim et al., 2022). However, Spectral GNN literature has shown the benefit of adaptively performing high- and low-frequency operations (Sun et al., 2022; Bo et al., 2021; Dong et al., 2021). Motivated by these works, we have developed a theoretically grounded method for adaptively selecting frequency components in a data-driven manner. Intuitively, our method involves finding eigenvectors that are most aligned with the class labels, which we characterize in terms of *energy spectral density*, defined as follows:

**Definition 2.2** (Energy Spectral Density (ESD) of class-labels). Given the orthonormal Laplacian eigenvectors  $V \in \mathbb{R}^{n \times n}$ , a one-hot class-indicator matrix  $Y \in \{0, 1\}^{n \times c}$ , and let  $V_i$  denote the  $i^{\text{th}}$  column of V. We define the class label ESD of the class-labels:

$$ESD_{i} = \frac{\|V_{i}^{\top}Y\|_{2}^{2}}{\sum_{j=1}^{n} \|V_{j}^{\top}Y\|_{2}^{2}}, \quad i = 1, \dots, n$$

<sup>&</sup>lt;sup>1</sup>Eigenvectors associated with small eigenvalues vary slowly across the graph, capturing "low-frequency" global variations, while eigenvectors corresponding to large eigenvalues vary rapidly, encoding "high-frequency" signals that change significantly across neighboring nodes (Shuman et al., 2013; Ortega et al., 2018).

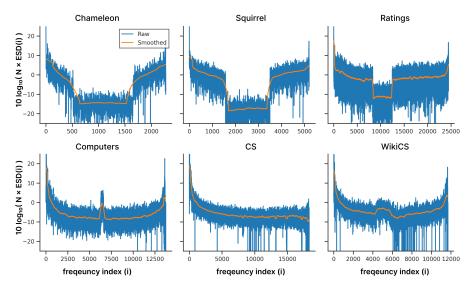


Figure 1: Energy spectral density (ESD) of the class labels across the Laplacian spectrum for real-world graphs. Peaks in mid- and high-frequency regions indicate that class-relevant signals are not confined to the low end of the spectrum.

Here,  $ESD_i$  measures the proportion of label energy aligned with eigenvector  $V_i$ . In Section 3, we show using our theoretical framework that choosing the eigenvectors with the highest ESD leads to desirable attention matrices with a class-aligned structure. Thus we rank eigenvectors by their ESD, providing a data-driven approach for efficiently utilizing the graph spectrum as positional encodings.

**ESD reveals limitations of data-agnostic selection heuristics.** Figure 1 reveals that label energy can appear in different regions of the spectrum: sometimes concentrated at low frequencies, sometimes at high frequencies, and often distributed heterogeneously. A key limitation of current eigenvector selection strategies is that they are data-agnostic, relying on fixed rules such as truncating to the lowest modes, enforcing symmetric low–high splits, or sampling at random. Because no single band is universally optimal for all graphs, such heuristics yield positional encodings that are misaligned with the task, limiting the effectiveness of graph transformers.

## 2.2 Broaden the Spectrum (BTS)

To overcome this limitation, we introduce *Broaden the Spectrum (BTS)*, a principled algorithm for selecting eigenvectors that are most informative for node classification. Rather than assuming that useful signal lies in a particular band of frequencies, we measure the alignment between eigenvectors and class labels, and select those with the highest *label ESD*. This simple yet powerful idea reframes positional encodings as a *learned spectral alignment problem*, bridging the gap between graph signal processing and transformer-based architectures.

Given a budget of k eigenvectors, BTS selects the top-k eigenvectors with the highest label ESD. This broadens the usable spectrum to include whichever eigenvectors carry discriminative signal, rather than assuming they reside at the bottom/top of the spectrum. Importantly, this selection is computed only from training labels, and we employ boxcar smoothing to mitigate noise. (Algorithm 1)

## 2.3 ARCHITECTURAL MODIFICATIONS

When expanding the spectrum of LPEs used in transformer models, we found it critical to incorporate a slight modification to the input encoder design. The selected eigenvectors  $X_{\rm pos}$ , along with the node features  $X_{\rm node}$  are passed through two independent MLPs, followed by concatenation to form the transformer's input tokens:

$$Y_{\text{GT}^*} = \text{Transformer}\left(\left[\text{MLP}_{\text{n}}(X_{\text{node}}); \text{MLP}_{\text{p}}(\text{norm}(X_{\text{pos}}))\right]\right) W_{\text{out}} \tag{1}$$

# Algorithm 1 Pseudocode for BTS eigenvector selection

**Input:** Laplacian eigenvectors  $V \in \mathbb{R}^{n \times n}$ ,

Training node indices  $\mathcal{I}_{\text{train}}$ , training labels  $Y_{\text{train}} \in \{0,1\}^{n_{\text{train}} \times c}$ , Number of eigenvectors to choose k, Smoothing window size w

**Output:** Indices  $\mathcal{I}_k$  of selected eigenvectors

```
1: V_{\text{train}} \leftarrow V[\mathcal{I}_{\text{train}}]  
ightharpoonup \text{Restrict to training nodes}

2: \tilde{Y}_{\text{train}} \leftarrow V_{\text{train}}^{\top} Y_{\text{train}}  
ightharpoonup \text{Graph Fourier transform}

3: \mathbf{for} \ i = 1 \text{ to } N \text{ do}  
ightharpoonup \text{Compute energy spectrum}

4: E_i \leftarrow \|\tilde{Y}_{\text{train},i}\|_2^2

5: ESD \leftarrow E/\sum_{i=1}^N E_i  
ightharpoonup \text{Normalize energy}

6: \overline{ESD} \leftarrow \text{BoxcarSmooth}(ESD, w)  
ightharpoonup \text{Smooth with window-size } w

7: \mathbf{return} \ \text{TopK}(\overline{ESD}, k)  
ightharpoonup \text{Return top-} k \ \text{ESD indices}
```

Here,  $\operatorname{norm}(\cdot)$  denotes row-wise  $\ell_2$  normalization. As we show in Section 4.3, this simple modification significantly improves performance when we increase the amount of selected eigenvectors. Normalization is critical here because the scale of Laplacian eigenvector elements is  $\sim 1/n$ , which quickly vanishes for reasonably sized graphs. Meanwhile, independent MLPs provide a more expressive mapping from raw inputs to transformer-compatible tokens.

# 3 THEORETICAL ANALYSIS

Having introduced BTS, we now present its theoretical foundations. Positional encodings play a crucial role in shaping the attention matrices of a transformer. This motivates us to break the problem of finding a good LPE subset into two steps: (i) we first show, through an illustrative linear model, that the optimal attention matrix for node classification has a *class-wise block structure* (Section 3.1), and (ii) we find that the Laplacian eigenvectors needed to be able to approximate such a block structure are exactly the ones as described by our label-ESD-based criteria. (Section 3.2)

## 3.1 Understanding Optimal Attention Matrices for node classification

Given a data matrix  $X \in \mathbb{R}^{n \times d}$ , the attention operation is defined as:

$$\operatorname{Attn}(X) = \operatorname{softmax}(XW_Q^\top W_K X^\top) X W_V^\top, \tag{2}$$

for some learnable weights  $W_Q$ ,  $W_K$ , and  $W_V$ . Here, the attention score matrix, softmax $(XW_QW_K^{\top}X)$ , is constrained by the softmax operator and the dependence on X. We lift these constraints, simplifying the softmax-attention operation into a linear one, and ask:

**Q1.** What form should a general attention matrix  $A \in \mathbb{R}^{n \times n}$  take so that the resulting latents Z := AX are most easily classifiable?

We assume a single-layer setting with a linear classifier. Given c classes, let  $Y \in \{0,1\}^{n \times c}$  be the one-hot class assignment matrix with  $Y_{ij} = 1$  if node i belongs to class j and zero otherwise. The classification objective is:

$$\mathcal{L}_{\text{class}}(A, W_C) = \text{CrossEntropy}(AXW_C^{\top} \mid Y), \tag{3}$$

where  $W_C \in \mathbb{R}^{c \times d}$  are classifier weights. We assume a mixture model  $X = YM_X + \sigma N$  for classmeans  $M_X \in \mathbb{R}^{c \times d}$ , noise variance  $\sigma > 0$ , and isotropic zero-mean noise N with  $\mathbb{E}[NN^{\top}] = dI_n$  We also assume balanced classes.

We first simplify the loss formulation through the following structural lemma:

**Lemma 3.1.** There exists a global minimizer  $(A^*, W_C^*)$  of Equation (3) such that all samples from the same class are mapped to the same latent, i.e.

$$A^*X = YM_Z$$
 for some  $M_Z \in \mathbb{R}^{c \times d}$ .

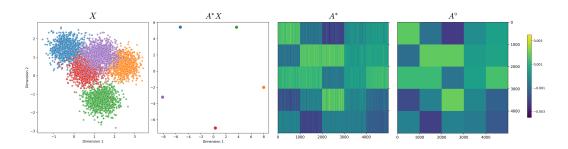


Figure 2: **Optimal Attention Matrices for Node Classification**. We find the minimizer for  $\mathcal{L}_{\text{class}}$  (Equation (3)) on 2 dimensional Guassian mixture model data with 5 classes. The resulting latents and attention matrix are plotted. Upon arranging the nodes based on their class index, we observe that the corresponding optimal attention matrix,  $A^*$  has an approximate class-wise block-structure of a form as predicted in  $A^{\circ}$ .

Intuitively, Lemma 3.1 (see Appendix C.1 for proof) shows that the optimal attention  $A^*$  clusters latents by class. Our goal is to probe the structure of such an  $A^*$ , and therefore, we study the surrogate objective:

$$A^{\circ} = \underset{A}{\operatorname{argmin}} \ \mathbb{E}_{N} \left[ \|AX - YM_{Z}\|_{F}^{2} \right], \tag{4}$$

which leads to the following:

**Theorem 3.1** (Optimal attention has class-wise block structure). Suppose  $X = YM_X + \sigma N$  as above with balanced classes. Then every minimizer  $A^{\circ}$  of Equation (4) admits the representation

$$A^{\circ} = Y M_A^{\circ} Y^{\top}, \quad \text{for some } M_A^{\circ} \in \mathbb{R}^{c \times c}.$$

That is,  $A_{ij}^{\circ}$  depends only on class memberships of nodes i and j. Theorem 3.1 (proof in Appendix C.2) formalizes the intuition that the *best attention matrix acts as a block matrix over classes*, ignoring within-class differences. Simulations confirm that empirical minimizers of Equation (3) indeed exhibit this block structure (Figure 2).

## 3.2 EIGENVECTOR SUBSETS FOR APPROXIMATING BLOCK ATTENTION

We now connect the block-structured optimal attention to spectral encodings. Consider a simplified linear attention formulation using only positional encodings:

$$A = X_{\text{pos}} W X_{\text{pos}}^{\top} \tag{5}$$

where W is a learnable full-rank matrix. Our next goal is to understand:

**Q2.** Which eigenvectors allow the best approximation of the block-structured optimum  $A^{\circ} = YMY^{\top}$ , uniformly for any  $M \in \mathbb{R}^{c \times c}$ ?

Essentially, our goal is to find a k-sized subset of eigenvectors so as to best approximate  $A \approx Y M Y^{\top}$  uniformly for any  $M \in \mathbb{R}^{c \times c}$ . The following theorem gives a sufficient criterion as a corollary:

**Theorem 3.2** (Uniform error bound for attention approximation). Let V denote the Laplacian eigenvectors and  $H \in \{0,1\}^{n \times k}$  denote an eigenvector selection matrix s.t.  $H_{ij} = 1$  iff eigenvector i is selected at position j. Define the diagonal 0/1 projector  $\tilde{H} := HH^{\top}$ , with  $\tilde{H}_{ii} = 1$  iff eigenvector i was selected. Set  $X_{pos} = VH$  in Equation (5). Then the uniform (in M) error functional:

$$\Phi(\tilde{H}) \; := \; \sup_{\|M\|_2 \leq 1} \; \min_{W} \, \left\| \, X_{pos} W X_{pos}^\top - Y M Y^\top \right\|_F$$

is bounded by the residual  $E := (I - \tilde{H})V^{\top}Y$ :

$$\Phi(\tilde{H}) \le 2\sqrt{n} \|E\|_F + \|E\|_F^2.$$

Theorem 3.2 (see Appendix C.3 for proof) shows that the quality of approximation is controlled entirely by the residual  $\|E\|_F$ , i.e., how well the selected eigenvectors capture the class-indicator matrix Y. In other words, the smaller the projection error of Y, the closer the resulting attention matrix is to the block-structured optimum. This is exactly what is minimized by the class-label ESD based ranking described in our method Section 2.2, leading to the following corollary.

**Corollary 3.2.1** (Class-label ESD based eigenvector selection). Among all k-sized eigenvector subsets, the selector that chooses the k eigenvectors with the largest label spectral-energy  $\|V_i^{\mathsf{T}}Y\|_2^2$  minimizes the bound in Theorem 3.2.

Corollary 3.2.1 (see Appendix C.4 for proof) tells us that our label-ESD based eigenvector selection criteria would lead to class-wise block-structured attention matrices, which we have identified to be desirable in Section 3.1. Moreover, in Section 4.2, we empirically validate this prediction and show that models trained with BTS produce attention matrices with stronger class-wise block structure.

## 4 RESULTS

In this section, we evaluate our approach on a diverse set of node classification benchmarks, analyzing its effectiveness across three established graph transformer architectures. Our evaluation focuses on measuring improvements in classification performance and understanding how spectral information is utilized.

**Experimental setup.** We evaluate our approach on homophilic, heterophilic, and long-range datasets (see Appendix D for details), using three standard graph transformer backbones: GT, NAGphormer, and GraphGPS. GT (Dwivedi & Bresson, 2020) is a direct application of transformers to graphs. NAGphormer (Chen et al., 2023) restricts attention to *K*-hop neighborhoods using a normalized adjacency matrix. GraphGPS (Rampášek et al., 2022) combines message passing with transformer-based global attention and uses LPE for positional encoding.

We use the subscript BTS to denote models tuned and trained with our ESD-based eigenvector selection approach (Section 2), as well as the input encoder modifications (Section 2.3), with  $k \leq 8192$ . Models without the subscript follow standard truncation to the  $k \leq 16$  lowest eigenvectors, (Table 7). For fairness, we also expand the GT baseline to broad spectrum setting (Section 4.3), and show that expansion alone provides little benefit without encoder modifications. Full training and hyperparameter details are given in Appendix E.

Table 1: Node classification performance on heterophilic benchmarks. Performance numbers for GT/GT<sub>BTS</sub>, NAGphormer/NAGphormer<sub>BTS</sub>, and GraphGPS/GraphGPS<sub>BTS</sub> were (re-)produced with our consistent experimental setup. Performance for other models are reported from existing literature. "-" indicates absence of a particular evaluation in existing literature. The top-1<sup>st</sup>, top-2<sup>nd</sup>, and top-3<sup>rd</sup> results are highlighted.

Model	Chameleon Accuracy ↑	<b>Squirrel</b> Accuracy ↑	<b>Tolokers</b> AU-ROC↑	Ratings Accuracy ↑
GCN	38.44 ± 1.92	$31.52 \pm 0.71$	$83.64 \pm 0.67$	$48.70 \pm 0.63$
GraphSAGE	$58.73 \pm \textbf{1.68}$	$41.61 \pm 0.74$	$82.43 \pm 0.44$	$\textbf{53.63} \pm \textbf{0.39}$
GAT	$48.36 \pm \textbf{1.58}$	$36.77 \pm 1.68$	$83.70 \pm 0.47$	$49.09 \pm \textbf{0.63}$
NodeFormer	$36.38 \pm 3.85$	$38.89 \pm 2.67$	$78.10 \pm 1.03$	43.79 ± 0.57
SGFormer	$45.21 \pm 3.72$	$42.65 \pm 4.21$	-	$\textbf{54.14} \pm \textbf{0.62}$
Exphormer	-	-	$83.53 \pm 0.28$	$50.48 \pm 0.34$
SpExphormer	-	-	$83.34 \pm 0.31$	$50.48 \pm \textbf{0.34}$
GT	$50.48 \pm 2.08$	$34.70 \pm 1.77$	$80.30 \pm 0.91$	$49.02 \pm 0.61$
GT <sub>BTS</sub>	<b>73.09</b> $\pm$ 1.00 +22.61 $\uparrow$	<b>65.06</b> ± 1.93 +30.36↑	$84.45 \pm 0.66$ +4.15 $\uparrow$	$50.37 \pm 0.48 \textbf{+1.35} \uparrow$
NAGphormer	$52.41 \pm 2.21$	$40.21 \pm 1.77$	$83.69 \pm 0.86$	$50.16 \pm 0.69$
$NAGphormer_{BTS}$	$\textbf{73.90} \pm \textbf{1.68} + \textbf{21.49} \uparrow$	<b>65.04</b> ± 1.69 +24.83↑	$\textbf{85.47} \pm \textbf{0.72} + \textbf{1.78} \uparrow$	$49.65 \pm 0.65  \textbf{-0.51} \!\!\downarrow$
GraphGPS	$60.92 \pm 2.54$	$43.43 \pm 1.46$	$\textbf{86.29} \pm \textbf{0.68}$	$50.19 \pm 0.51$
<b>GraphGPS</b> <sub>BTS</sub>	<b>73.16</b> $\pm$ 1.70 +12.24 $\uparrow$	$\textbf{65.87} \pm \textbf{1.30}  + \textbf{22.44} \!\!\uparrow$	$\pmb{86.31} \pm \textbf{0.63} + \pmb{0.02} \uparrow$	$\textbf{51.33} \pm \textbf{0.58} + \textbf{1.14} \uparrow$

<sup>&</sup>lt;sup>2</sup>For graphs larger than 8192 nodes, we only compute the low and high 4096 eigenvectors.

Table 2: Node classification accuracy (%) on homophilic benchmarks. Results for GraphGPS/GraphGPS $_{BTS}$ , NAGphormer/NAGphormer $_{BTS}$ , and GT/GT $_{BTS}$  were (re-)produced with our consistent experimental setup. Performance for other models are reported from existing literature. The top-1<sup>st</sup>, top-2<sup>nd</sup>, and top-3<sup>rd</sup> results are highlighted.

Model	Physics Accuracy ↑	<b>CS</b> Accuracy ↑	<b>Photo</b> Accuracy ↑	<b>Computers</b> Accuracy ↑	<b>WikiCS</b> Accuracy ↑	ogbn-arXiv Accuracy ↑
NodeFormer	$96.45 \pm 0.28$	$95.64 \pm 0.22$	$93.46 \pm 0.35$	$86.98 \pm 0.62$	$74.73 \pm 0.94$	59.90 ± 0.42
SGFormer	$96.60 \pm 0.18$	$94.78 \pm 0.34$	$95.10 \pm 0.47$	$91.99 \pm 0.70$	$73.46 \pm 0.56$	$\textbf{72.63} \pm 0.13$
Exphormer	$96.89 \pm 0.09$	$94.93 \pm 0.01$	$95.35 \pm 0.22$	$91.47 \pm 0.17$	$78.19 \pm 0.29$	$\textbf{71.27} \pm 0.27$
SpExphormer	$96.70 \pm 0.05$	$95.00 \pm 0.15$	$95.33 \pm 0.49$	$91.09 \pm 0.08$	$78.20 \pm 0.14$	$70.82 \pm 0.24$
GT	$96.02 \pm 0.20$	94.66 ± 0.44	$91.59 \pm 0.68$	$85.65 \pm 0.59$	$72.91 \pm 0.59$	55.68 ± 0.39
$GT_{BTS}$	$96.90 \pm 0.18  \textbf{+0.88} \uparrow$	$95.44 \pm 0.33  \textbf{+0.78} \uparrow$	$\textbf{95.95} \pm \textbf{0.48}  + \textbf{4.36} \uparrow$	$91.46 \pm 0.51  \textbf{+5.81} \uparrow$	<b>78.94</b> ± 0.26 +6.03↑	$70.30 \pm 0.12 + 14.62 \uparrow$
NAGphormer	$96.98 \pm 0.13$	$\textbf{95.71} \pm \textbf{0.26}$	$95.51 \pm 0.41$	$91.39 \pm 0.41$	$78.73 \pm 0.66$	69.43 ±0.32
$NAGphormer_{BTS}$	97.05 ± 0.18 +0.07↑	$95.42 \pm 0.39  \textbf{-0.29} \! \downarrow$	<b>95.90</b> ± 0.37 +0.39↑	$91.85 \pm 0.44 + 0.46 \uparrow$	$\textbf{79.42} \pm \textbf{0.55} + \textbf{0.69} \uparrow$	<b>71.29</b> ± 0.13 + <b>1.86</b> ↑
GraphGPS	$97.13 \pm 0.17$	$95.70 \pm 0.38$	$95.35 \pm 0.45$	$91.64 \pm 0.46$	$77.67 \pm 0.73$	65.16 ± 1.45
<b>GraphGPS</b> <sub>BTS</sub>	$\textbf{97.21} \pm \textbf{0.14} + \textbf{0.08} \uparrow$	$\textbf{95.72} \pm \textbf{0.37} + \textbf{0.02} \uparrow$	$\textbf{95.87} \pm \textbf{0.42} + \textbf{0.52} \uparrow$	$\textbf{91.87} \pm \textbf{0.45} + \textbf{0.23} \uparrow$	$\textbf{79.47} \pm \textbf{0.48} + \textbf{1.80} \uparrow$	$70.92 \pm 0.33$ +5.76 $\uparrow$

## 4.1 Main Results

**Results on heterophilic benchmarks.** We find substantial improvements when using BTS on heterophilic graphs (Table 1). For example, on *Chameleon*, performance improves by over 22%, and on *Squirrel*, by over 30% when BTS-filtered eigenvectors are used with a simple transformer architecture (GT). Remarkably, this brings the vanilla transformer architecture into close competition with, and in some cases even surpassing, more complex graph transformer models proposed in recent literature. To the best of our knowledge, the performance reported here for *Chameleon*, *Squirrel*, and *Tolokers* represents the strongest results achieved by any graph transformer model to date.

These improvements can be understood through the lens of the class-label ESD. As shown in Figure 1, graphs like *Chameleon* and *Squirrel* exhibit significant class energy in a broad set of low and high frequency regions. These findings highlight that the historical reliance on low-frequency truncation was a critical bottleneck, masking the true representational and generalization potential of graph transformers.

**Results on homophilic benchmarks.** As shown in Table 2, BTS improves performance even on homophilic graphs.  $GT_{BTS}$  achieves +5.8% on *Computers*, +6.0% on *WikiCS*, and +14.4% on *ogbn-arXiv*. Gains for NAGphormer and GraphGPS are smaller but consistent. These results are expectedly more

modest than in heterophilic settings, as the spectral energy of class signals in many homophilic graphs is concentrated in low frequencies. Still, BTS captures the broader spectrum whenever useful, yielding robust improvements.

Results on long-range benchmarks. Graph transformers are naturally suited for capturing long-range dependencies due to their global attention mechanism. However, a number of recent papers have shown that graph transformers suffer from overglobalization (Xing et al., 2024) and perform poorly on long-range tasks. As shown in Table 3, our method achieves substantial improvements on the long-range benchmark (Liang et al., 2025) over baseline graph transformer architectures. These datasets exhibit particularly strong gains when using BTS-selected features. For instance, performance for GT on *Paris* improves by over 38%, and on *Shanghai* by 31%, bringing the vanilla transformer model on par with strong baseline methods.

Table 3: Node classification accuracy (%) on Long Range Benchmarks The top-1<sup>st</sup>, top-2<sup>nd</sup>, and top-3<sup>rd</sup> results are highlighted.

Model	Paris	Shanghai
GCN	$47.30 \pm 0.20$	$52.40 \pm 0.30$
GraphSAGE	$49.10 \pm 0.60$	$60.40 \pm 0.30$
SGFormer	$45.00 \pm 0.20$	$53.5 \pm 0.30$
GT	$15.46 \pm 3.93$	$21.05 \pm 0.51$
$GT_{BTS}$	$\textbf{53.79} \pm \textbf{0.17}$	$52.66 \pm 0.83$
Δ	+38.33 ↑	+31.61↑
NAGphormer	$25.26 \pm 0.34$	$24.94 \pm 0.28$
$NAGphormer_{BTS}$	$\textbf{53.68} \pm 0.23$	$\textbf{57.26} \pm \textbf{0.34}$
$\Delta$	+28.42 ↑	+32.32 ↑
GraphGPS	$28.99 \pm 0.31$	$28.46 \pm 0.31$
GraphGPS <sub>BTS</sub>	$\textbf{54.12} \pm \textbf{0.23}$	$\textbf{55.31} \pm \textbf{0.33}$
Δ	+25.13 ↑	+26.85↑

## 4.2 ATTENTION MATRICES

Our theoretical analysis (Section 3) predicts that the optimal attention matrix for node classification should have a class-wise block structure and that using BTS-selected eigenvectors would lead to such attention matrices. To validate this, we examine the attention matrices obtained after training. As shown in Figure 3, models trained with BTS yield attention matrices that indeed display a substantially

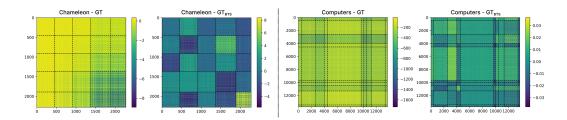


Figure 3: Attention matrices learned by vanilla GT and GT<sub>BTS</sub> on Chameleon (left) and Computers (right). Incorporating BTS leads to clearer class-wise block structures, consistent with our theoretical prediction that such attention patterns are optimal for node classification (Section 3.1). These attention matrices correspond to the first head in the last layer, and no)es are sorted according to their classes.

stronger block-wise organization aligned with class partitions. Since BTS models perform better, this also suggests that class-wise block structures are indeed desirable attention patterns.

## 4.3 ABLATIONS

We conduct ablation studies to isolate the contributions of two key factors: (i) the design of the encoder that processes positional encodings, and (ii) the strategy used to select eigenvectors. Our results show that architectural support is necessary to benefit from larger k, while principled selection is essential for avoiding overfitting and consistently improving performance.

**Spectral expansion and encoder design.** The results in Table 4 indicate that expanding k (num. of eigenvectors) beyond this range yields only marginal improvements without our input encoder modifications. For example, with GT, Chameleon's performance increases from 50.48% to 52.28%. This partly explains why most approaches have restricted positional encodings to a small set of low-frequency eigenvectors. However, we see substantial gains when the positional encodings are normalized and processed with a MLP. With these modifications, performance improves to 67.83% on Chameleon, and from 34.70% to 62.91% on Squirrel. This indicates that, while a broader spectrum of eigenvectors provides useful signal, careful design of the encoder is needed to fully utilize them.

Table 4: Ablation of our encoder modifications. Increasing number of eigenvectors (k) beyond low frequencies alone offers only marginal gains without our architectural modifications.

<b>Position Encoder</b>	k	Chameleon	Squirrel	WikiCS	Computers
Baseline (linear)	$k \le 16$	$50.48 \pm 2.08$	34.70 ± 1.77	$72.91 \pm 0.59$	$85.65 \pm 0.59$
Baseline (linear)	no-bound	$52.28 \pm 2.88$	$37.71 \pm 1.79$	$73.75 \pm 0.63$	$87.21 \pm 0.55$
Norm + Linear	no-bound	$65.07 \pm 1.08$	$56.96 \pm 1.23$	$77.67 \pm 0.48$	$91.39 \pm 0.40$
Norm + MLP	no-bound	$\textbf{67.83} \pm \textbf{1.82}$	$\textbf{62.91} \pm \textbf{1.04}$	$\textbf{78.46} \pm \textbf{0.56}$	$\textbf{91.66} \pm \textbf{0.41}$

**Eigenvector selection strategies.** Next, we explore how different eigenvector selection heuristics influence performance. Table 5 shows that simply using the full spectrum severely hurts generalization due to overfitting, demonstrating that indiscriminate inclusion of all frequencies is detrimental. We also tested four different selection heuristics: low-only, high-only, low+high, and low+medium+high, each corresponding to retaining different bands of the spectrum. While all four of these heuristics show improvements over baseline performance, we find that their effectiveness is not consistent across datasets. For example, the best out of the four variants is high-only on both heterophilic datasets (Chameleon, Squirrel) and low-only on homophilic datasets (WikiCS, Computers). In contrast, our data-driven selection (BTS) consistently achieves the best results.

## 5 RELATED WORK

**Graph Transformers and Positional Encodings.** Graph Transformers (GTs) allow graph nodes to interact globally via self-attention (Dwivedi & Bresson, 2020). Because self-attention is permutation-invariant, GTs require positional encodings (PEs) to inject structural information. Laplacian eigenvectors have become a common choice (Hoang et al., 2024), providing a spectral basis that reflects

Table 5: Impact of different eigenvector selection strategies on the performance of GT\*. The top-1<sup>st</sup>, top-2<sup>nd</sup>, and top-3<sup>rd</sup> results are highlighted.

Eigenvector Selection	Chameleon	Squirrel	WikiCS	Computers
Baseline	$50.48 \pm 2.08$	$34.70 \pm 1.77$	$72.91 \pm 0.59$	$85.65 \pm 0.59$
Full spectrum	$48.14 \pm 3.05$	$35.30 \pm 1.61$	$72.35 \pm 0.72$	$85.20 \pm 0.30$
Low-only	$67.83 \pm 1.82$	$62.91 \pm 1.04$	$\textbf{78.46} \pm \textbf{0.56}$	$91.26 \pm 0.41$
High-only	$72.79 \pm 1.37$	$63.22 \pm 1.75$	$73.45\ \pm0.40$	$89.80\ \pm0.55$
Low + High (equal)	<b>72.50</b> $\pm$ 0.93	$63.19 \pm 1.49$	$\textbf{78.37} \ \pm \textbf{0.52}$	$90.99 \pm 0.50$
Low + Medium + High (equal)	$66.51 \pm 1.88$	$44.66 \pm 1.59$	$75.38 \pm 0.56$	$89.79 \pm 0.33$
BTS	$\textbf{73.09} \pm 1.68$	$\textbf{65.06} \pm 1.93$	$\textbf{78.94} \pm \textbf{0.26}$	$91.46 \pm 0.51$

graph topology. However, nearly all existing GTs adopt a heuristic truncation: retaining only the first k low-frequency eigenvectors. The low-frequency bias is evident in models such as GT (Dwivedi & Bresson, 2020), NAGphormer (Chen et al., 2023), GraphTrans (Wu et al., 2021), GraphGPS (Rampášek et al., 2022), UGT (Lee et al., 2024), SAN (Kreuzer et al., 2021), and SAT (Chen et al., 2022). Scalability-oriented variants such as ANS-GT (Zhang et al., 2022), Gapformer (Liu et al., 2023), and Exphormer (Shirzad et al., 2023), which focus on efficient attention also retain the same positional encoding heuristic. TokenGT (Kreuzer et al., 2021) extends beyond low frequencies by including both low- and high-frequency eigenvectors, but relies on fixed splits rather than task-driven selection.

The Role of High-Frequency Signals in Node Classification. In contrast, the MPNN literature has increasingly emphasized the importance of high-frequency information for node classification, especially in heterophilic graphs. Spectral methods explicitly operate in the frequency domain and modulate both low- and high-frequency signals (Wu et al., 2019; Dong et al., 2020). Early spectral models such as Spectral CNN (Bruna et al., 2013) and ChebNet (Defferrard et al., 2016) introduced learnable filters over eigenvalues, and subsequent works demonstrated that high-frequency information is critical for node-level expressiveness, such as adaptive propagation (Chien et al., 2020), complete spectral filtering (Luan et al., 2020), and frequency-based feature selection (Bo et al., 2021). These techniques explicitly exploit high-frequency modes, while adaptive approaches such as AdaGNN (Dong et al., 2021), and spectral attention (Chang et al., 2021) dynamically balance contributions from different parts of the spectrum.

## 6 Conclusion

In this paper, we introduced Broaden the Spectrum (BTS), a lightweight and architecture-agnostic algorithm for selecting Laplacian eigenvectors as positional encodings. We show that spectral distribution of class information is unique for each graph, making fixed heuristics for eigenvector selection (such as low-frequency truncation) inherently limited. BTS addresses this by selecting eigenvectors in a data-driven way, consistently improving the performance of graph transformers across homophilic, heterophilic, and long-range benchmarks. In several cases, BTS elevates even simple transformer architectures to state-of-the-art levels.

To ground our method, we also developed a theoretical framework for attention-based node classification, showing that the optimal attention matrix exhibits a class-wise block structure and that eigenvectors most correlated with labels are best suited to approximate it. While we focus on the Laplacian spectrum, our framework is general and can be extended to other orthonormal bases. However, both BTS and our theoretical analysis are tailored to the supervised node classification setting, and it remains an open direction to adapt similar principles to other tasks such as link prediction, graph classification, or self-supervised pretraining.

Overall, our work highlights that the performance bottleneck of graph transformers on node classification lies in how positional encodings are chosen. By broadening the spectrum with a data-driven, theoretically grounded selection, we show that this bottleneck can be overcome and that even simple graph transformers are competitive and reach near state-of-the-art performance.

# REFERENCES

- Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *The 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2623–2631, 2019.
- Mehdi Azabou, Venkataramana Ganesh, Shantanu Thakoor, Chi-Heng Lin, Lakshmi Sathidevi, Ran Liu, Michal Valko, Petar Veličković, and Eva L Dyer. Half-hop: A graph upsampling approach for slowing down message passing. In *International Conference on Machine Learning*, pp. 1341–1360. PMLR, 2023.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K.Q. Weinberger (eds.), Advances in Neural Information Processing Systems, volume 24. Curran Associates, Inc., 2011. URL https://proceedings.neurips.cc/paper\_files/paper/2011/file/86e8f7ab32cfd12577bc2619bc635690-Paper.pdf.
- Deyu Bo, Xiao Wang, Chuan Shi, and Huawei Shen. Beyond low-frequency information in graph convolutional networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 3950–3957, 2021.
- A survey on spectral graph neural networks, 2023. URL https://arxiv.org/abs/2302.05631.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- Heng Chang, Yu Rong, Tingyang Xu, Wenbing Huang, Somayeh Sojoudi, Junzhou Huang, and Wenwu Zhu. Spectral graph attention network with fast eigen-approximation. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pp. 2905–2909, 2021.
- Dexiong Chen, Leslie O'Bray, and Karsten Borgwardt. Structure-aware transformer for graph representation learning. In *International Conference on Machine Learning*, pp. 3469–3489. PMLR, 2022.
- Jinsong Chen, Kaiyuan Gao, Gaichao Li, and Kun He. Nagphormer: A tokenized graph transformer for node classification in large graphs. 2023. URL https://arxiv.org/abs/2206.04910.
- Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network. *arXiv preprint arXiv:2006.07988*, 2020.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- Xiaowen Dong, Dorina Thanou, Laura Toni, Michael Bronstein, and Pascal Frossard. Graph signal processing for machine learning: A review and new perspectives. *IEEE Signal processing magazine*, 37(6):117–127, 2020.
- Yushun Dong, Kaize Ding, Brian Jalaian, Shuiwang Ji, and Jundong Li. Adagnn: Graph neural networks with adaptive frequency response filter. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pp. 392–401, 2021.
- An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL https://arxiv.org/abs/2010.11929.
- Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs. *arXiv preprint arXiv:2012.09699*, 2020.
- Van Thuy Hoang, O Lee, et al. A survey on structure-preserving graph transformers. *arXiv preprint* arXiv:2401.16176, 2024.
  - Pure transformers are powerful graph learners, 2022. URL https://arxiv.org/abs/2207.02505.
    - Devin Kreuzer, Dominique Beaini, William L. Hamilton, Vincent Létourneau, and Prudencio Tossou. Rethinking graph transformers with spectral attention. 2021. URL https://arxiv.org/abs/2106.03893.
    - O-Joun Lee et al. Transitivity-preserving graph representation learning for bridging local connectivity and role-based similarity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 12456–12465, 2024.
    - Huidong Liang, Haitz Sáez de Ocáriz Borde, Baskaran Sripathmanathan, Michael Bronstein, and Xiaowen Dong. Towards quantifying long-range interactions in graph machine learning: a large graph dataset and a measurement. *arXiv preprint arXiv:2503.09008*, 2025.

- Chuang Liu, Yibing Zhan, Xueqi Ma, Liang Ding, Dapeng Tao, Jia Wu, and Wenbin Hu. Gapformer: Graph transformer with graph pooling for node classification. In *Proceedings of the 32nd International Joint Conference on Artificial Intelligence (IJCAI-23)*, pp. 2196–2205, 2023.
  - Decoupled weight decay regularization, 2019. URL https://arxiv.org/abs/1711.05101.
  - Sitao Luan, Mingde Zhao, Chenqing Hua, Xiao-Wen Chang, and Doina Precup. Complete the missing half: Augmenting aggregation filtering with diversification for graph convolutional networks. *arXiv preprint arXiv:2008.08844*, 2020.
  - Classic gnns are strong baselines: Reassessing gnns for node classification, 2024. URL https://arxiv.org/abs/2406.08993.
  - Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, pp. 43–52, 2015.
  - Wiki-cs: A wikipedia-based benchmark for graph neural networks, 2022. URL https://arxiv.org/abs/2007.02901.
  - A time series is worth 64 words: Long-term forecasting with transformers, 2023. URL https://arxiv.org/abs/2211.14730.
  - Antonio Ortega, Pascal Frossard, Jelena Kovačević, José M. F. Moura, and Pierre Vandergheynst. Graph signal processing: Overview, challenges, and applications. *Proceedings of the IEEE*, 106(5):808–828, 2018. doi: 10.1109/JPROC.2018.2820126.
  - Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geometric graph convolutional networks. *arXiv preprint arXiv:2002.05287*, 2020.
  - Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova. A critical look at the evaluation of gnns under heterophily: Are we really making progress? *arXiv preprint arXiv:2302.11640*, 2023.
  - Ladislav Rampášek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. *Advances in Neural Information Processing Systems*, 35:14501–14515, 2022.
  - Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. *Journal of Complex Networks*, 9(2):cnab014, 2021.
  - A survey on oversmoothing in graph neural networks, 2023. URL https://arxiv.org/abs/2303.10993.
  - Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J Sutherland, and Ali Kemal Sinop. Exphormer: Sparse transformers for graphs. In *International Conference on Machine Learning*, pp. 31613–31632. PMLR, 2023.
  - Hamed Shirzad, Honghao Lin, Balaji Venkatachalam, Ameya Velingker, David Woodruff, and Danica Sutherland. Even sparser graph transformers. *arXiv preprint arXiv:2411.16278*, 2024.
  - David I Shuman, Sunil K. Narang, Pascal Frossard, Antonio Ortega, and Pierre Vandergheynst. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine*, 30(3):83–98, 2013. doi: 10.1109/MSP.2012.2235192.
  - Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June Hsu, and Kuansan Wang. An overview of microsoft academic service (mas) and applications. In *Proceedings of the 24th international conference on world wide web*, pp. 243–246, 2015.
  - Jiaqi Sun, Lin Zhang, Shenglin Zhao, and Yujiu Yang. Improving your graph neural networks: A high-frequency booster. In 2022 IEEE International Conference on Data Mining Workshops (ICDMW), pp. 748–756. IEEE, 2022.
  - Understanding over-squashing and bottlenecks on graphs via curvature, 2022. URL https://arxiv.org/abs/2111.14522.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. 2017. URL https://arxiv.org/pdf/1706.03762.pdf.

convolutional networks. In International conference on machine learning, pp. 6861-6871. PMLR, 2019. Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E Gonzalez, and Ion Stoica. Represent-ing long-range context for graph neural networks with global attention. Advances in Neural Information Processing Systems, 34:13266-13279, 2021. Yujie Xing, Xiao Wang, Yibo Li, Hai Huang, and Chuan Shi. Less is more: on the over-globalizing problem in graph transformers. arXiv preprint arXiv:2405.01102, 2024. Zaixi Zhang, Qi Liu, Qingyong Hu, and Chee-Kong Lee. Hierarchical graph transformer with adaptive node sampling. Advances in Neural Information Processing Systems, 35:21171–21183, 2022. 

Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph

# **APPENDIX**

# A ABLATIONS FOR ALL ARCHITECTURES

We conducted an ablation study in Section 4.3 to evaluate the impact of higher-order spectral components, encoder design, and label-aware selection on the Graph Transformer (GT). Here, we extend this analysis to all baseline transformer architectures. As shown in Table 6, the same trends hold across models, showing that access to more eigenvectors and a better encoder design can improve performance.

Table 6: Node classification performance with full eigenvector spectrum vs with left-truncated spectrum but tuned K.

		Heterophilic		Home	ophilic
Model	Eigenvectors	Chameleon	Squirrel	WikiCS	Computers
GT	$K \in [4, 16]$ (tuned)	$50.48 \pm 2.08$	$34.70 \pm 1.77$	$72.91 \pm 0.59$	$85.65 \pm 0.59$
GT	$K \in [4, N]$ (tuned)	$52.28 \pm 2.88$	$37.71 \pm 1.79$	$73.75 \pm 0.63$	$87.21 \pm 0.55$
GT* Linear	$K \in [4, N]$ (tuned)	$65.07 \pm 1.08$	$56.96 \pm 1.23$	$77.67 \pm 0.48$	$91.39 \pm 0.40$
GT* MLP	$K \in [4, N]$ (tuned)	$67.83 \pm 1.82$	$62.91 \pm 1.04$	$78.46 \pm 0.56$	$91.66 \pm 0.41$
GT*	full spectrum (fixed)	$48.14 \pm 3.05$	$35.30 \pm 1.61$	$72.35 \pm 0.72$	$85.20 \pm 0.30$
$GT_{BTS}$	BTS, $K \in [4, N]$ (tuned)	$73.09 \pm 1.68$	$65.06 \pm 1.93$	$78.94 \pm 0.26$	$91.46 \pm 0.51$
NAGphormer	$K \in [4, 16]$ (tuned)	$52.41 \pm 2.21$	40.21 ± 1.77	$78.73 \pm 0.66$	$91.39 \pm 0.41$
NAGphormer	$K \in [4, N]$ (tuned)	$57.06 \pm 1.96$	$40.89 \pm 2.06$	$79.70 \pm 0.50$	$91.61 \pm 0.42$
NAGphormer* Linear	$K \in [4, N]$ (tuned)	- ± -	- ± -	$79.33 \pm 0.44$	$92.16 \pm 0.48$
NAGphormer* MLP	$K \in [4, N]$ (tuned)	$70.07 \pm 2.33$	$63.87 \pm 1.51$	$79.83 \pm 0.63$	$91.96 \pm 0.37$
NAGphormer*	full spectrum (fixed)	$59.63 \pm 2.06$	$52.27 \pm 1.28$	$79.63 \pm 0.63$	$91.53 \pm 0.47$
NAGphormer <sub>BTS</sub>	BTS, $K \in [4, N]$ (tuned)	$73.90 \pm 1.68$	$65.04 \pm 1.69$	$79.42 \pm 1.55$	$91.85 \pm 0.44$
GraphGPS	$K \in [4, 16]$ (tuned)	$60.92 \pm 2.54$	$43.43 \pm 1.46$	$77.67 \pm 0.73$	$91.64 \pm 0.46$
GraphGPS	$K \in [4, N]$ (tuned)	$64.67 \pm 2.98$	$47.12 \pm 4.21$	$77.40 \pm 0.45$	$91.60 \pm 0.45$
GraphGPS* Linear	$K \in [4, N]$ (tuned)	$68.8 \pm 2.13$	$57.67 \pm 0.97$	$78.92 \pm 0.53$	$91.96 \pm 0.22$
GraphGPS* MLP	$K \in [4, N]$ (tuned)	$70.24 \pm 2.08$	$63.40 \pm 1.16$	$78.68 \pm 0.46$	$91.80 \pm 0.40$
GraphGPS*	full spectrum (fixed)	$59.14 \pm 1.95$	$42.88 \pm 1.77$	$77.42 \pm 0.98$	$91.24 \pm 0.40$
$GPS_{BTS}$	BTS, $K \in [4, N]$ (tuned)	$73.16 \pm 1.70$	$65.87 \pm 1.30$	$79.47 \pm 0.48$	$91.87 \pm 0.45$

# B HIGHEST MAXIMUM FREQUENCIES USED IN THE LITERATURE

Table 7: Maximum number of eigenvectors ( $K_{\text{max}}$ ) used by recent graph transformer models, based on publicly available code.

Model	$K_{max}$
NAGphormer (Chen et al., 2023)	15
GraphGPS (Rampášek et al., 2022)	10
SAN (Kreuzer et al., 2021)	10
GT (Dwivedi & Bresson, 2020)	10
UGT (Lee et al., 2024)	10
Exphormer (Shirzad et al., 2023)	10

To better understand the typical frequency truncation choices in existing graph transformer models, we compile representative values of the maximum number of Laplacian eigenvectors (k) used across a range of published works. As shown in Table 7, most methods restrict k to a small number—often below 16—reinforcing the low-pass inductive bias observed in current practice.

# C PROOFS

#### C.1 Proof for Lemma 3.1

*Proof.* Given a data matrix  $X \in \mathbb{R}^{n \times d}$  and the minimization problem:

$$\min_{A,W_C} \mathcal{L}_{\text{class}}(A, W_C), \tag{6}$$

our goal is to characterize the associated latents  $Z^* := A^*X$  of an optimal solution  $(A^*, W_C^*)$ .

Let  $(A, W_C)$  be an arbitrary candidate solution to the minimization problem (6) and define the resulting latents Z := AX consisting of rows  $\mathbf{z}_i \in \mathbb{R}^d$ . Let  $Y \in \{0,1\}^{n \times c}$  denote the one-hot class indicator matrix,  $y_i$  denote the class index of node i, and  $n_j$  is the number of points in class j. First we write the cross entropy classification objective:

$$\mathcal{L}_{\text{class}}(A, W_C) = \sum_{i=1}^{n} \ell(W_C \mathbf{z}_i, y_i), \qquad (7)$$

with  $\ell(\mathbf{u}, y_i) = -\mathbf{u}_{y_i} + \log \sum_{j=1}^c e^{\mathbf{u}_j}$  for any predicted logits  $\mathbf{u} \in \mathbb{R}^c$ . We note that  $\ell(\cdot, y_i)$  is convex in  $\mathbf{u}$ .

Now define the  $j^{th}$ -class mean-latent:

$$\overline{\mathbf{z}}_j := \frac{1}{n_j} \sum_{i: y_i = j} \mathbf{z}_i. \quad \forall j \in \{1 \dots c\}$$
 (8)

By Jensen's inequality, since  $\ell$  is convex and  $\mathbf{z} \mapsto W_C \mathbf{z}$  is linear, we have

$$\frac{1}{n_{j}} \sum_{i:y_{i}=j} \ell\left(W_{C} \mathbf{z}_{i}, y_{i}\right) \geq \ell\left(W_{C} \left(\frac{1}{n_{j}} \sum_{i:y_{i}=j} \mathbf{z}_{i}\right), y_{i}\right) \quad \forall j \in \{1 \dots c\}$$

$$\triangleq \ell\left(W_{C} \overline{\mathbf{z}}_{i}, y_{i}\right), \tag{9}$$

i.e.

$$\sum_{i:y_i=j} \ell(W_C \mathbf{z}_i, y_i) \ge n_j \ell(W_C \overline{\mathbf{z}}_j, y_i). \quad \forall j \in \{1 \dots c\}$$
(10)

Hence, if we define  $Z^*$  as the latent matrix with every row  $\mathbf{z}_i^*$  set to  $\overline{\mathbf{z}}_{y_i}$ , then we have for any  $W_C$ :

$$\mathcal{L}_{\text{class}}(Z^*, W_C) \le \mathcal{L}_{\text{class}}(Z, W_C). \tag{11}$$

In particular, if we minimize over all  $W_C$ , we have:

$$\min_{W_C} \mathcal{L}_{\text{class}}(Z^*, W_C) \le \min_{W_C} \mathcal{L}_{\text{class}}(Z, W_C). \tag{12}$$

Thus, from any (approximately) optimal Z we can construct a collapsed  $Z^*$  (with  $\mathbf{z}_i^* := \overline{\mathbf{z}}_{y_i}$ ) that is at least as good as Z. This relationship holds for any  $W_C$ , therefore after solving the minimization problem in Equation (6) we obtain a global minimizer  $(A^*, W_C^*)$  with:

$$A^*X \triangleq Z^* = YM_Z \tag{13}$$

where  $M_Z \in \mathbb{R}^{c \times d}$  is the matrix of class mean-latents with rows  $\overline{\mathbf{z}}_i \in \mathbb{R}^d$ .

## C.2 Proof for Theorem 3.1

*Proof.* Given X and Y defined as in the proof from Appendix C.1, assume a mixture model for the data matrix:

$$X = YM_X + \sigma N \tag{14}$$

where  $M_X \in \mathbb{R}^{c \times d}$  is a matrix of class-means,  $\sigma > 0$  denotes noise variance, and N denotes isotropic zero-mean noise with  $\mathbb{E}[NN^T] = dI_n$ .

Consider the noise-averaged surrogate objective defined in Equation (4):

$$\min_{A} \mathcal{L}_{\text{attn}}(A) = \mathbb{E}_{N} \left[ \|AX - YM_{Z}\|_{F}^{2} \right]$$
 (15)

where  $M_Z \in \mathbb{R}^{c \times d}$  is the class mean-latent matrix as defined in Lemma 3.1 and Appendix C.1. We shall omit the N subscript from  $\mathbb{E}_N$  for notational convenience. The purpose of the surrogate loss is to probe the structure of an attention matrix that attempts to solves the classification problem defined in Equation (12).

Assume Y is defined such that the classes are balanced, i.e.  $Y^{\top}Y = \frac{n}{c}I_c$ . Then define the orthogonal projector onto col(Y):

$$P := Y(Y^{\top}Y)^{-1}Y^{\top} = \frac{c}{n}YY^{\top}$$
 (16)

Note that P is symmetric and idempotent. Our goal is to reveal class-wise block structure of attention matrices that minimize the surrogate objective.

First, let the attention  $A \in \mathbb{R}^{n \times n}$  be arbitrary. Since  $YM_Z \in \text{col}(Y) = \text{range}(P)$ , we have the orthogonal decomposition:

$$||AX - YM_Z||_F^2 = ||P(AX - YM_Z) + (I_n - P)(AX - YM_Z)||_F^2$$
  
=  $||P(AX - YM_Z)||_F^2 + ||(I_n - P)AX||_F^2$  (17)

Now define  $\tilde{A} := PA$  and notice that, by idempotency of P, the first term in the decomposition stays the same while the second term vanishes. Hence, we have:

$$\|\tilde{A}X - YM_Z\|_F^2 = \|P(AX - YM_Z)\|_F^2 \le \|AX - YM_Z\|_F^2$$
(18)

Thus, since the inequality is preserved under expectation, we have  $\mathcal{L}_{attn}(PA) \leq \mathcal{L}_{attn}(A)$  for all A. In particular, given a minimizer  $A^{\circ}$  of  $\mathcal{L}_{attn}$ , we have  $\mathcal{L}_{attn}(PA^{\circ}) = \mathcal{L}_{attn}(A^{\circ})$ . But then by the decomposition, we know that

$$0 = \mathbb{E}\left[ \|(I_n - P)A^{\circ}X\|_F^2 \right] \ge \sigma^2 d \|(I_n - P)A^{\circ}\|_F^2$$
(19)

using the data model assumption  $X = YM_X + \sigma N$ . Hence, we have  $PA^{\circ} = A^{\circ}$ .

Similarly, for any A, we can again using the data model assumption to expand the surrogate loss:

$$\mathcal{L}_{\text{attn}}(A) = \|AYM_X - YM_Z\|_F^2 + \sigma^2 d \|A\|_F^2$$
 (20)

By definition of P, since  $YM_X \in \operatorname{col}(Y)$ , we have  $PYM_X = YM_X$ , i.e. the signal term stays the same if we make the substitution  $\tilde{A} := AP$ . Hence, by idempotency of P (and consequently idempotency of  $I_n - P$ ),

$$\mathcal{L}_{\operatorname{attn}}(A) - \mathcal{L}_{\operatorname{attn}}(AP) = \sigma^{2} d \left( \|A\|_{F}^{2} - \|AP\|_{F}^{2} \right)$$

$$= \sigma^{2} d \left( \operatorname{tr}(AA^{\top}) - \operatorname{tr}(APA^{\top}) \right)$$

$$= \sigma^{2} d \left( \operatorname{tr}(A(I_{n} - P)A^{\top}) \right)$$

$$= \sigma^{2} d \|A(I_{n} - P)\|_{F}^{2} \ge 0$$
(21)

Again, if we take a minimizer  $A^{\circ}$  of  $\mathcal{L}_{\text{attn}}$ , then this tells us that  $\mathcal{L}_{\text{attn}}(A^{\circ}) = \mathcal{L}_{\text{attn}}(A^{\circ}P)$ , i.e.  $\|A^{\circ}(I_n - P)\|_F = 0$ , i.e.  $A^{\circ} = A^{\circ}P$ .

Combining the results so far, we can conclude that any minimizer  $A^{\circ}$  admits the representation

$$A^{\circ} = PA^{\circ}P$$
$$= YM_{A}^{\circ}Y^{\top} \tag{22}$$

with  $M_A^{\circ} := \frac{c^2}{n^2}(Y^{\top}AY)$ . Since we know a minimizer will admit such a representation, we can restrict the surrogate loss in terms of the kernel matrix  $M_A$ :

$$\mathcal{L}_{\text{attn}}(YM_{A}Y^{\top}) = \mathbb{E}\left[\left\|YM_{A}Y^{\top}(YM_{X} + \sigma N) - YM_{Z}\right\|_{F}^{2}\right]$$

$$= \left\|Y\left(\frac{n}{c}M_{A}M_{X} - M_{Z}\right)\right\|_{F}^{2} + \sigma^{2}\mathbb{E}\left[\left\|YM_{A}Y^{\top}N\right\|_{F}^{2}\right]$$

$$= \frac{n}{c}\left\|\frac{n}{c}M_{A}M_{X} - M_{Z}\right\|_{F}^{2} + \frac{\sigma^{2}n^{2}d}{c^{2}}\left\|M_{A}\right\|_{F}^{2}$$
(23)

Lastly, we show that  $M_A^\circ := \frac{c}{n} M_Z M_X^\top (M_X M_X^\top + \frac{c}{n} \sigma^2 dI_c)^{-1}$  yields a closed-form minimizer for the surrogate loss, and it is in fact the *unique* minimizer. To do so, consider another candidate solution specified by  $M_A := M_A^\circ + \Delta$ . We expand the loss:

$$\mathcal{L}_{\text{attn}} (Y M_A Y^{\top}) = \frac{n}{c} \left\| \frac{n}{c} M_A^{\circ} M_X - M_Z + \frac{n}{c} \Delta M_X \right\|_F^2 + \frac{\sigma^2 n^2 d}{c^2} \left\| M_A^{\circ} + \Delta \right\|_F^2$$

$$= \mathcal{L}_{\text{attn}} (Y M_A^{\circ} Y^{\top}) + \frac{n}{c} \left\| \frac{n}{c} \Delta M_X \right\|_F^2 + \frac{\sigma^2 n^2 d}{c^2} \left\| \Delta \right\|_F^2 + 2\mathcal{R}(\Delta)$$
(24)

where  $\mathcal{R}(\Delta)$  is given by:

$$\mathcal{R}(\Delta) = \frac{n}{c} \left\langle \frac{n}{c} \Delta M_X, \frac{n}{c} M_A^{\circ} M_X - M_Z \right\rangle_F + \frac{\sigma^2 n^2 d}{c^2} \left\langle \Delta, M_A^{\circ} \right\rangle_F$$
$$= \frac{n^2}{c^2} \left[ \left\langle \Delta, \left( \frac{n}{c} M_A^{\circ} M_X - M_Z \right) M_X^{\top} + \sigma^2 d M_A^{\circ} \right\rangle_F \right] \tag{25}$$

We can expand the right term in the Frobenius inner product using the proposed solution  $M_A^{\circ}$ :

$$\left(\frac{n}{c}M_{A}^{\circ}M_{X} - M_{Z}\right)M_{X}^{\top} + \sigma^{2}dM_{A}^{\circ}$$

$$= M_{Z}M_{X}^{\top}\left(M_{X}M_{X}^{\top} + \frac{c}{n}\sigma^{2}dI_{c}\right)^{-1}M_{X}M_{X}^{\top} + \frac{c}{n}\sigma^{2}dM_{Z}M_{X}^{\top}\left(M_{X}M_{X}^{\top} + \frac{c}{n}\sigma^{2}dI_{c}\right)^{-1} - M_{Z}M_{X}^{\top}$$

$$= M_{Z}M_{X}^{\top}\left(M_{X}M_{X}^{\top} + \frac{c}{n}\sigma^{2}dI_{c}\right)^{-1}\left(M_{X}M_{X}^{\top} + \frac{c}{n}\sigma^{2}dI_{c}\right) - M_{X}M_{X}^{\top}$$

$$= M_{Z}(M_{X}^{\top} - M_{X}^{\top})$$

$$= 0$$
(26)

Hence,  $\mathcal{R}(\Delta) = 0$ , i.e.  $\mathcal{L}_{\text{attn}}(YM_AY^\top) \leq \mathcal{L}_{\text{attn}}(YM_A^{\circ}Y^\top)$  with equality iff  $\Delta = 0$ . Indeed,  $M_A^{\circ}$  is the unique solution.

#### C.3 Proof for Theorem 3.2

*Proof.* To reiterate the setup, we consider the simplified linear attention formulation using only positional encodings, as in Equation (5). Let V denote the Laplacian eigenvectors and  $H \in \{0,1\}^{n \times k}$  denote an eigenvector selection matrix s.t.  $H_{ij} = 1$  iff eigenvector i is selected at position j. Define the diagonal 0/1 projector  $\tilde{H} := HH^{\top}$ , with  $\tilde{H}_{ii} = 1$  iff eigenvector i was selected, and set  $X_{\text{pos}} := VH$  in the formulation. Then define the error functional over all  $M \in \mathbb{R}^{c \times c}$  within the unit ball  $\|M\|_2 \le 1$ :

$$\Phi(\tilde{H}) := \sup_{\|M\|_{2} \le 1} \min_{W} \|X_{\text{pos}} W X_{\text{pos}}^{\top} - Y M Y^{\top}\|_{F} 
= \sup_{\|M\|_{2} \le 1} \min_{W} \|V H W H^{\top} V^{\top} - Y M Y^{\top}\|_{F}.$$
(27)

The error functional defines the quality of an approximation of the block-structured representation  $YMY^{\top}$ , uniformly across all kernels M. Our goal is to provide a uniform error bound for this functional

We first observe that since eigenvectors are chosen without replacement, VH will have full column rank. Hence, the inner optimization problem is a well-defined classical linear least squares problem that can be solved exactly for W:

$$W^*(\tilde{H}) := H^\top V^\top Y M Y^\top V H. \tag{28}$$

Hence, the functional becomes:

$$\Phi(\tilde{H}) = \sup_{\|M\|_{2} \le 1} \|VHH^{\top}V^{\top}YMY^{\top}VHH^{\top}V^{\top} - YMY^{\top}\|_{F} 
= \sup_{\|M\|_{2} \le 1} \|\tilde{H}V^{\top}YMY^{\top}V\tilde{H} - V^{\top}YMY^{\top}V\|_{F} 
= \sup_{\|M\|_{2} \le 1} \|\tilde{H}V^{\top}YM(\tilde{H}V^{\top}Y)^{\top} - V^{\top}YMY^{\top}V\|_{F}$$
(29)

since  $\tilde{H}$  is symmetric. Define the residual  $E := (\mathbb{I}_n - \tilde{H})V^\top Y$ . Then we can rewrite the inner norm and bound it:

$$\|\tilde{H}V^{\top}YM(\tilde{H}V^{\top}Y)^{\top} - V^{\top}YMY^{\top}V\|_{F} = \|(V^{\top}Y - E)M(V^{\top}Y - E)^{\top} - V^{\top}YMY^{\top}V\|_{F}$$

$$= \|EMY^{\top}V + V^{\top}YME^{\top} - EME^{\top}\|_{F}$$

$$\leq \|M\|_{2} \left(2\|E\|_{F}\|V^{\top}Y\|_{F} + \|E\|_{F}^{2}\right)$$

$$\leq \|M\|_{2} \left(2\sqrt{n}\|E\|_{F} + \|E\|_{F}^{2}\right) \tag{30}$$

since V is orthonormal. Hence, within the unit ball  $||M||_2 \le 1$ , we have the upper bound:

$$\|\tilde{H}V^{\top}YM(\tilde{H}V^{\top}Y)^{\top} - V^{\top}YMY^{\top}V\|_{F} \le 2\sqrt{n}\|E\|_{F} + \|E\|_{F}^{2}.$$
 (31)

By definition, we have:

$$\Phi(\tilde{H}) \le 2\sqrt{n} \|E\|_F + \|E\|_F^2. \tag{32}$$

Thus,  $||E||_F$  gives us uniform control over the upper error bound of the given  $\tilde{H}$  over all  $||M||_2 \leq 1$ .

# C.4 Proof for Corollary 3.2.1

*Proof.* Following the conclusion of Theorem 3.2, we showed that the uniform error bound for attention approximation is controlled by the residual  $||E||_F$ , where  $E := (I - \tilde{H})V^\top Y$ . To minimize  $||E||_F$ , we first observe that:

$$||E||_F^2 \triangleq \sum_{i=1}^n \left(1 - \tilde{H}_{ii}\right) ||V_i^\top Y||_2^2$$
 (33)

where  $\tilde{H}_{ii}$  indicates whether or not eigenvector i was included, and  $s_i := ||V_i^\top Y||_2^2$  is the  $\ell_2$ -norm of the i<sup>th</sup> row of  $V^\top Y$ . Thus, to minimize  $||E||_F$ , we should include eigenvectors corresponding to the k largest  $s_i$ 's (in terms of *spectral-energy*  $||s_i||_2^2$ ), i.e. the optimal selector  $\tilde{H}^*$  that minimizes the upper error bound is exactly the one based on the top k indices of the  $s_i$ 's.

## D BASELINE SOURCES

## D.1 HETEROPHILIC DATASETS

We use five real-world datasets with graphs that have a homophily level  $\leq 0.30$ : Actor (Pei et al., 2020), Chameleon and Squirrel (Rozemberczki et al., 2021), as well as Ratings and Tolokers (Platonov et al., 2023). Key statistics for these datasets are listed in Table 8. We follow the experimental setup in (Pei et al., 2020) for Actor, Chameleon, and Squirrel, and for Ratings and Tolokers, we adopt the setup described in (Platonov et al., 2023), using the 10 train/validation/test splits provided.

The results for GCN-based methods and heterophily-based methods in Table 8 for Actor, Chameleon, and Squirrel have been sourced from (Azabou et al., 2023). Similarly, results for Ratings and Tolokers are sourced from (Platonov et al., 2023), while results for transformer-based methods across all datasets are obtained from (Shirzad et al., 2024).

918 919

Table 8: Statistics of heterophilic datasets used in our experiments.

924 925

926 927 928

929

930 931 932

934 935 936

933

937 938

939 940 941

946 947

948 949 950

951

956

957

958 959 960

961

962 963 964 965

966 967

968 969

970

971

DATASET **NODES EDGES** CLASSES HOMOPHILY RATIO 2,277 5 CHAMELEON 31,421 0.23 SQUIRREL 5,201 198,493 5 0.22 **TOLOKERS** 11,758 519,000 2 0.09 244,92 5 0.14 RATINGS 39,402

Table 9: Statistics of homophilic datasets used in our experiments.

DATASET	Nodes	EDGES	CLASSES	HOMOPHILY RATIO
PHYSICS	34,493	495,924	5	0.92
CS	18,333	81,894	15	0.83
Рното	7,650	238,162	8	0.84
COMPUTERS	13,752	491,722	10	0.79
WIKICS	11,701	216,123	10	0.66
OGBN-ARXIV	169,343	1,166,243	40	0.65

#### D.2 HOMOPHILIC DATASETS

We use five real-world datasets: Amazon Computers and Amazon Photos (McAuley et al., 2015), Coauthor CS and Coauthor Physics (Sinha et al., 2015), and WikiCS (Mernyei & Cangea, 2022). Key statistics for these datasets are listed in Table 9. The experimental setup follows that of (Shirzad et al., 2024), where the datasets are split into development and test sets. All hyperparameter tuning is performed on the development set, and the best models are subsequently evaluated on the test set.

We use a 60:20:20 train/validation/test split for the Amazon and Coauthor datasets. The results reported for all datasets in Table 2 are sourced from (Shirzad et al., 2024).

## D.3 Long Range Benchmark Datasets

To evaluate the ability of models to capture long-range dependencies, we use the City-Networks benchmark (Liang et al., 2025), which consists of large-scale road network graphs derived from OpenStreetMap data. We focus on two representative cities—Paris and Shanghai—which feature grid-like topology, low clustering coefficients, and large diameters. These characteristics make them particularly well-suited for studying long-range signal propagation. Key statistics for these datasets are provided in Table 10.

Following the experimental protocol in (Liang et al., 2025), we perform transductive node classification using a 10:10:80 train/validation/test split. The node labels are defined by eccentricity-based quantiles, ensuring that the task inherently depends on information from distant nodes.

Table 10: Statistics of City-Networks datasets used in our experiments.

DATASET	Nodes	EDGES	CLASSES	HOMOPHILY RATIO
Paris	114,127	- ,-	10	0.70
Shanghai	183,917		10	0.75

## D.4 BASELINE MODEL PERFORMANCE ACROSS DATASETS FROM EXISTING LITERATURE

The previously reported performance of baseline models (GT, GraphGPS, and NAGphormer) on multiple graph datasets is summarized in Table 11. The reported values, sourced from existing literature.

Table 11: Performance across datasets for GT, GraphGPS, and NAGphormer models previously reported in existing literature.

Dataset	GT	GraphGPS	NAGphormer
Chameleon	-	$40.79 \pm 4.03$	-
Squirrel	-	$39.67 \pm 2.84$	-
Tolokers	-	$83.71 \pm 0.48$	$78.32 \pm 0.95$
Ratings	-	$53.10 \pm 0.42$	$51.26 \pm 0.72$
Physics	$97.05 \pm 0.05$	$97.12 \pm 0.19$	$97.34 \pm 0.03$
CS	$94.64 \pm 0.13$	$93.93 \pm 0.12$	$95.75 \pm 0.09$
Photo	$94.74 \pm 0.13$	$95.06 \pm 0.13$	$95.49 \pm 0.11$
Computers	$91.18 \pm 0.17$	$91.19 \pm 0.54$	$91.22 \pm 0.14$
WikiCS	-	$78.66 \pm 0.49$	$77.16 \pm 0.72$
Arxiv	-	$70.97\pm0.41$	$70.13\pm0.55$

## E ADDITIONAL TRAINING DETAILS

**Optimizer.** We use the AdamW optimizer (Loshchilov & Hutter, 2019) for all runs, and fixed the number of epochs to 200. We additionally employ the linear-warmup-cosine-decay learning rate schedule. Linear rate warmup happens over 10 epochs (fixed), and cosine decay happens over the remaining 190 epochs (also fixed). All other hyperparameters are chosen by the tuning algorithm explained below.

**Hyperparameter tuning.** We optimize hyperparameters using the Tree-structured Parzen Estimator (TPE) algorithm (Bergstra et al., 2011), as implemented in Optuna (Akiba et al., 2019). The complete hyperparameter space used in our experiments is detailed in Table 12. The number of tuning trials is adjusted based on the size of each dataset: for graphs with up to 7,500 nodes, we perform 300 tuning trials; for graphs with up to 15,000 nodes, we allow 200 trials; and for larger graphs, we limit the number of trials to 100. Performing complete hyperparameter tuning, on a machine with  $4 \times NVidia$  L40S GPUs takes 2-4 hours depending on the size of the dataset. Hyperparameters are selected based on validation-set performance, and all reported results correspond to test-set performance using the best configurations found. The scripts used for hyperparameter tuning are also included in our codebase.

# F LLM USAGE DISCLOSURE

We used LLMs solely for the purpose of editing and polishing the paper.

Table 12: Complete hyperparameter search space for all model variants presented in this paper.

Hyperparameter	SEARCH SPACE	SAMPLING TYPE	
COMMON PARAMETERS			
Learning rate	$[10^{-4}, 10^{-1}]$	Logarithmic	
WEIGHT DECAY	$[10^{-7}, 10^{-2}]$	Logarithmi	
Dropout	[0, 0.5]	LINEAR	
ATTENTION DROPOUT	[0, 0.5]	LINEAR	
WINDOW LENGTH	$\{256, 512, 1024, 2048, 4096\}$		
Transformer Depth	$\{1,2,\ldots,8\}$	LINEAR	
Number of attention heads	$\{0, 1, 2, 4, 8\}$		
COMMON FOR GT <sub>BTS</sub> /NAGPHORMER <sub>BTS</sub> /GRAPHGPS <sub>BTS</sub>			
Number of eigenvectors $(K)$	$\{4, 8, 16, \dots, 1024\}$		
Pos. feature encoder - output dimension	$\{8, 16, 32, 64, 128\}$		
Pos. feature encoder - hidden dimension	$\{16, 32, 64, \dots, 2048\}$		
Pos. feature encoder - # hidden layers	$\{1, 2, 3, 4\}$		
Node feature encoder - output dimension	$\{8, 16, 32, 64, 128\}$		
NODE FEATURE ENCODER - HIDDEN DIMENSION	$\{16, 32, 64, 128\}$		
Node feature encoder - # hidden layers	{1}		
SPECIFIC FOR GT			
Number of eigenvectors $(K)$	$\{4, 8, 16\}$		
TOKEN DIMENSION	$\{64, 128, 256, 512\}$		
SPECIFIC FOR NAGPHORMER			
Number of eigenvectors $(K)$	$\{4, 8, 16\}$		
Token dimension	$\{64, 128, 256, 512\}$		
NUMBER OF HOPS (SAME FOR NAGPHORMER <sub>BTS</sub> )	$\{1, 2, 3, \dots, 20\}$		
SPECIFIC FOR GRAPHGPS			
Number of eigenvectors $(K)$	$\{4, 8, 16\}$		
LPE - NUMBER OF LAYERS	$\{1, 2, 3, \dots, 8\}$		
LPE - NUMBER OF POST-LAYERS	$\{0, 1, 2, 3, 4\}$		