
FIGURE: Simple and Efficient Unsupervised Node Representations with Filter Augmentations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Unsupervised node representations learnt using contrastive learning-based methods
2 have shown good performance on downstream tasks. However, these methods rely
3 on augmentations that mimic low-pass filters, limiting their performance on tasks
4 requiring different eigen-spectrum parts. This paper presents a simple filter-based
5 augmentation method to capture different parts of the eigen-spectrum. We show
6 significant improvements using these augmentations. Further, we show that sharing
7 the same weights across these different filter augmentations is possible, reducing the
8 computational load. In addition, previous works have shown that good performance
9 on downstream tasks requires high dimensional representations. Working with high
10 dimensions increases the computations, especially when multiple augmentations
11 are involved. We mitigate this problem and recover good performance through
12 lower dimensional embeddings using simple random Fourier feature projections.
13 Our method, FiGURE, achieves an average gain of up to 4.4%, compared to the
14 state-of-the-art unsupervised models, across all datasets in consideration, both
15 homophilic and heterophilic.

16 1 Introduction

17 Contrastive learning is a powerful method for unsupervised graph representation learning, achieving
18 notable success in various applications [35, 8]. However, these evaluations typically focus on tasks
19 exhibiting homophily, where task labels strongly correlate with the graph’s structure. An existing
20 edge suggests the connected nodes likely share similar labels in these scenarios. However, these
21 representations often struggle when dealing with heterophilic tasks, where edges tend to connect
22 nodes with different labels.

23 Several papers [4, 10, 3, 20] have tackled the problem of heterophily by leveraging information from
24 both low and high-frequency components. However, these methods operate in the semi-supervised
25 setting, and the extension of these ideas in unsupervised learning still needs to be explored. Inspired
26 by the insights in these papers, we propose a simple method incorporating these principles. Our
27 approach introduces filter banks as additional views and learns separate representations for each filter
28 bank. However, this approach faces two main challenges: Firstly, storing representations from each
29 view can become prohibitively expensive for large graphs; secondly, contrastive learning methods
30 typically demand high-dimensional representations, which increase both the computational cost of
31 training and the storage burden.

32 We employ a shared encoder for all filter banks to tackle the first challenge. Our results confirm
33 that a shared encoder performs on par with independent encoders for each filter bank. This strategy
34 enables us to reconstruct filter-specific representations as needed, drastically reducing the storage
35 requirement.

36 For the second challenge, we train our models with low-dimensional embeddings. Then, we use
 37 random Fourier feature projection [31] to lift these low-dimensional embeddings into a higher-
 38 dimensional space. Kernel tricks [15] were typically used in classical machine learning to project
 39 low-dimensional representation to high dimensions where the labels can become linearly separable.
 40 However, constructing and leveraging the kernels in large dataset scenarios could be expensive. To
 41 avoid this issue, several papers [31, 32, 13, 28, 18] proposed to approximate the map associated with
 42 the kernel. For our scenario, we use the map associated with Gaussian kernel [31]. We empirically
 43 demonstrate that using such a simple approach preserves high performance for downstream tasks,
 44 even in the contrastive learning setting. Consequently, our solution offers a more efficient approach
 45 to unsupervised graph representation learning in computation and storage, especially concerning
 46 heterophilic tasks.

47 Our contributions in this work are, 1] We propose a simple scheme of using filter banks for learning
 48 representations that can cater to both heterophily and homophily tasks, 2] We address the computa-
 49 tional and storage burden associated with this simple strategy by sharing the encoder across these
 50 various filter views, 3] By learning a low-dimensional representation and later projecting it to high
 51 dimensions using random Fourier Features, we further reduce the burden, 4] We study the perfor-
 52 mance of our approach on four homophilic and seven heterophilic datasets. Our method achieves
 53 new SOTA performance in unsupervised representation learning on heterophilic datasets, achieving
 54 18% gains over prior methods, and in homophilic datasets [22], achieving 1.5% gains over prior
 55 methods. Our method also performs better than supervised methods such as GCN [17] on several
 56 heterophilic datasets and is competitive on homophilic datasets.

57 2 Related Work

58 Several unsupervised representation learning methods have been proposed in prior literature. Random
 59 walk-based methods like Node2Vec [7] and DeepWalk [29] preserve node proximity but tend to ne-
 60 glect structural information and node features. Contrastive methods, such as DEEP GRAPH INFOMAX
 61 (DGI) [35], maximize the mutual information (MI) between local and global representations while
 62 minimizing the MI between corrupted representations. Methods like MVGRL [8] and GRACE [36]
 63 expand on this, by integrating additional views into the MI maximization objective.

64 However, most of these methods focus on the low frequency components, overlooking critical insights
 65 from other parts. Semi-supervised methods like GPRGNN [4], BERNNET [10], and PPGNN [20]
 66 address this by exploring the entire eigenspectrum, but these concepts are yet to be applied in the
 67 unsupervised domain.

68 This work proposes the use of a filter bank to capture information across the full eigenspectrum
 69 while sharing an encoder across filters. Given the high-dimensional representation demand of
 70 contrastive learning methods, we propose using Random Fourier Features (RFF) to project lower-
 71 dimensional embeddings into higher-dimensional spaces, reducing computational load without
 72 sacrificing performance. The ensuing sections define our problem, describe filter banks and random
 73 feature maps, and explain our model and experimental results.

74 3 Problem Setting

75 In the domain of unsupervised representation learning, our focus lies on graph data, denoted as
 76 $G = (V; E)$, where V is the set of vertices and E the set of edges ($E \subseteq V \times V$). We associate an
 77 adjacency matrix with G , referred to as $\mathbf{A} : \mathbf{A} \in \mathbb{R}^{n \times n}$, where $n = |V|$ corresponds to the
 78 number of nodes. Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be the feature matrix. We use \mathbf{A}_I to represent $\mathbf{A} + \mathbf{I}$ with \mathbf{I} is the
 79 identity matrix, while $\mathbf{D}_{\mathbf{A}_I}$ signifies the degree matrix of \mathbf{A}_I . We also define \mathbf{A}_n as $\mathbf{D}_{\mathbf{A}_I}^{-1/2} \mathbf{A}_I \mathbf{D}_{\mathbf{A}_I}^{-1/2}$.
 80 No additional information is provided during training. The goal is to learn a parameterized encoder,
 81 $E : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^{n \times d}$, where $d \ll n$. This encoder produces a set of node representations
 82 $E(\mathbf{X}; \mathbf{A}_n) = \{h_1; h_2; \dots; h_n\}$ where each $h_i \in \mathbb{R}^d$ represents a rich representation for node i . The
 83 subsequent section will provide preliminary details about filter banks and random feature maps before
 84 we discuss the specifics of the proposed approach.

85 4 Preliminaries

86 Our proposed approach hinges on the critical components of filter banks and random feature maps. In
87 this section, we delve into brief details about these two facets, setting the stage for a comprehensive
88 description of our approach.

89 4.1 Filter Banks

90 Graph Fourier Transform (GFT) forms the basis of Graph Neural Networks (GNNs). A GFT is defined
91 using a reference operator \mathcal{R} which admits a spectral decomposition. Traditionally, in the case of
92 GNNs, this reference operator has been the symmetric normalized Laplacian \mathcal{L}_{sym} or the \mathcal{A}_n
93 as simplified in [17]. A graph filter is an operator that acts independently on the entire eigenspace of
94 a diagonalisable and symmetric reference operator \mathcal{R} by modulating their corresponding eigenvalues.
95 [34, 33]. Thus, a graph filter \mathcal{H} is defined via the graph filter function $g(\cdot)$ operating on the reference
96 operator as $\mathcal{H} = g(\mathcal{R}) = U g(\Lambda) U^T$. Here, $\Lambda = \text{diag}([\lambda_1; \lambda_2; \dots; \lambda_n])$, where λ_i denotes the
97 eigenvalues of the reference operator.

98 We describe a filter bank as a set of filters, denoted as $\mathcal{F} = \{F_1; F_2; \dots; F_K\}$. Both GPRGNN [4]
99 and BERNNET [10] employ filter banks, comprising polynomial filters, and amalgamate the represen-
100 tations from each filter bank to enhance the performance across heterophilic datasets. BERNNET uses
101 a filter bank defined as $\mathcal{F}_{\text{GPRGNN}} = \{f_1; A_n; \dots; A_n^{K-1}\}$, while $\mathcal{F}_{\text{BERNNET}} = \{B_0; B_1; \dots; B_K\}$
102 characterizes the filter bank utilized by BERNNET. Here, $B_i = \frac{1}{2^{K-i}} \mathbf{1}^{\otimes K-i} (2I - L_n)^{K-i} (L_n)^i$.

103 Each filter in these banks highlights different parts of the eigenspectrum. By tuning the combination
104 on downstream tasks, it offers the choice to select and leverage the right spectrum to enhance perfor-
105 mance. Notably, unlike traditional GNNs, which primarily emphasize low-frequency components,
106 higher frequency components have proved useful for heterophilic [10, 20]. Consequently, a
107 vital takeaway is that for comprehensive representations, we must aggregate information from
108 different parts of the eigenspectrum and re-tune it for specific downstream tasks

109 4.2 Random Feature Maps for Kernel Approximations

110 Before the emergence of deep learning models, the kernel trick was instrumental in learning non-linear
111 models. A kernel function $k : \mathcal{R}^d \times \mathcal{R}^d \rightarrow \mathbb{R}$, accepts two input features and returns a real-valued
112 score. Given a positive-definite kernel, Mercer's Theorem [26] assures the existence of a feature map
113 $\phi(\cdot)$, such that $k(x; y) = \langle \phi(x); \phi(y) \rangle$. Leveraging the kernel trick, researchers combined Mercer's
114 theorem with the representer theorem [5], enabling the construction of non-linear models that remain
115 linear in k . These models created directly using ϕ instead of the potentially complex k , outperformed
116 traditional linear models. The implicit maps linked with these kernels projected the features into a
117 significantly high-dimensional space, where targets were presumed to be linearly separable. However,
118 computational challenges arose when dealing with large datasets.

119 Addressing these issues, subsequent works [3, 28, 31] introduced approximations of the map
120 associated with individual kernels through random projections into higher-dimensional spaces ϕ
121 This approach ensures that $\langle \phi(x); \phi(y) \rangle \approx k(x; y)$. These random feature maps are inexpensive to
122 compute and affirm that simple projections to higher-dimensional spaces can achieve linear separa-
123 bility. The critical insight is that computationally efficient random feature maps exist, capable
124 of projecting lower-dimensional representations into higher dimensions. These projections
125 enhance the adaptability of these representations for downstream tasks. Random Fourier
126 features (RFF) [31] provide a prime example of such techniques

127 5 Proposed Approach

128 The following section delineates the process of unsupervised representation learning. Post that, we
129 give details on how the representations learned from each filter bank is used in downstream tasks
130 using random feature maps.

131 5.1 Unsupervised Representation Learning

132 Our method **FIGURE** (Filter-based Graph Unsupervised Representation Learning) builds on concepts
 133 introduced in [1, 35], extending the maximization of mutual information between node and global
 134 iter representations for each iter in the iter bank $\mathbb{R} = \{F_1; F_2; \dots; F_K\}$. We construct an encoder
 135 for each iter to maximize the mutual information between the input data and encoder output. For
 136 the i^{th} iter, we learn an encoder $E_i : X_i \rightarrow X_i^0$, denoted by learnable parameters. In this
 137 context, X_i represents a set of examples, where each example $\mathcal{X}_{ij} \in X_i$ consists of a iter
 138 F_i , its corresponding nodes and node features drawn from an empirical probability distribution
 139 which captures the joint distribution of features and node representations \mathcal{X}_{ij} . X_i denotes the
 140 set of representations learnt by the encoder on utilizing feature information as well as topological
 141 information from the samples, sampled from the joint distribution $P_{\mathcal{X}_{ij}}$. The goal, aligned with
 142 [21, 11, 35], is to identify E_i that maximizes mutual information between \mathcal{X}_{ij} and $E_i(\mathcal{X}_{ij}; F_i)$,
 143 or $I(\mathcal{X}_{ij}; F_i; E_i(\mathcal{X}_{ij}; F_i))$. While exact mutual information (MI) computation is unfeasible due to
 144 unavailable exact data and learned representations distributions, we can estimate using the
 145 Jensen-Shannon MI estimator [5, 25], defined as:

$$I_{i,j}^{\text{JSD}}(\mathcal{X}_{ij}; F_i; E_i(\mathcal{X}_{ij}; F_i)) := E_{P_i}[\text{sp}(T_{i,j}([\mathcal{X}_{ij}; \mathcal{F}_{ij}]; E_i(\mathcal{X}_{ij}; \mathcal{F}_{ij})))] \\ E_{P_i}[\text{sp}(T_{i,j}([\mathcal{X}_{ij}; \mathcal{F}_{ij}]; E_i(\mathcal{X}_{ij}; \mathcal{F}_{ij})))] \quad (1)$$

146 Here, $T_{i,j} : X_i \times X_i^0 \rightarrow \mathbb{R}$ represents a discriminator function with learnable parameters. Note that
 147 $[\mathcal{X}_{ij}; \mathcal{F}_{ij}]$ is an input sampled from P_i , which is a marginal of the joint distribution of the input data
 148 and the learned node representations. The function $\text{sp}(\cdot)$ corresponds to the softplus function [6].
 149 Additionally, $T_{i,j} = D_w(R(E_i(\mathcal{X}_{ij}; \mathcal{F}_{ij})); E_i(\mathcal{X}_{ij}; \mathcal{F}_{ij}))$, where R denotes the readout function
 150 responsible for summarizing all node representations by aggregating and distilling information into a
 151 global iter representation.

152 In our approach, we first obtain node
 153 representations by feeding the iter-speci c topology and associated node
 154 features into the encoder $H_i =$
 155 $E_i(\mathcal{X}_{ij}; F_i) = [h_1^{F_i}; h_2^{F_i}; \dots; h_n^{F_i}]$.
 156 To obtain global representations, we
 157 employ a readout function $R :$
 158 $\mathbb{R}^{N \times d^0} \rightarrow \mathbb{R}^{d^0}$, which combines
 159 and distills information into a global
 160 representation $h_g^{F_i} = R(H_i) =$
 161 $R(E_i(\mathcal{X}_{ij}; F_i))$. Instead of directly
 162 maximizing the mutual information
 163 between the local and global repre-
 164 sentations, we introduce a learnable
 165 discriminator $D_i : \mathbb{R}^{d^0} \times \mathbb{R}^{d^0} \rightarrow \mathbb{R}$,
 166 where $D_i(\cdot; \cdot)$ represents the joint
 167 probability score between the global
 168 representation and the node-speci c
 169 patch representation. This joint prob-
 170 ability score should be higher when
 171 considering global and local repre-
 172 sentations obtained from the same iter,
 173 as opposed to the joint probability
 174 score between the global representa-
 175 tion from one iter and the local rep-
 176 resentation from an arbitrary iter.

Figure 1: Unsupervised learning of node embeddings by maximizing mutual information between node and graph representations over the graphs from the iter bank. Note that the parameter is shared across all the iters.

178 To generate negative samples for contrastive learning, we employ a corruption function $C : \mathbb{R}^{N \times d^0} \rightarrow \mathbb{R}^{M \times d^0}$
 179 $\mathbb{R}^{N \times d^0} \rightarrow \mathbb{R}^{M \times d^0}$, which yields corrupted samples denoted as $[\mathcal{X}_{ij}^c; \mathcal{F}_{ij}^c] = C(\mathcal{X}_{ij}; F_i)$. The
 180 designed corruption function generates data decorrelated with the input data.

181 In order to learn representations across all lters in the lter bank, we aim to maximise the average
 182 estimate of mutual information (MI) across all lters, considering lters.

$$I_F = \frac{1}{K} \sum_{i=1}^K I_{i;F_i}^{\text{JSD}}([X; F_i]; E(X; F_i)) \quad (2)$$

183 Maximising the Jensen-Shannon estimator is equivalent to reducing the binary cross entropy loss
 184 defined between positive samples (sampled from the joint) and the negative samples (sampled from
 185 the product of marginals). Therefore, for each lter, we minimise the following objective:

$$L_{F_i} = \frac{1}{N + M} \sum_{j=1}^N E_{(X; F_i)}[\log(D_j(h_j^{F_i}; h_g^{F_i}))] + \sum_{j=1}^M E_{(X; F_i)}[\log(D_j(h_j^{F_i}; h_g^{F_i}))] \quad (3)$$

186 Therefore to learn meaningful representations across all lters the following objective is minimised:

$$L = \frac{1}{K} \sum_{i=1}^K L_{F_i} \quad (4)$$

187 However, managing the computational cost of training and storage for large graphs with separate node
 188 representations for each lter presents a significant challenge, exacerbated by the high dimensional
 189 requirements of contrastive learning methods. We implement parameter sharing to mitigate the
 190 issue, borrowing the concept from studies such as [4, 10], thereby sharing the encoder's parameters
 191 and the discriminator's parameters across all lters. Instead of storing dense lter-specific node
 192 representations, we only store the parameters of the shared encoder and the first-hop neighbourhood
 193 information of each node per lter, which has a lower storage cost. For downstream tasks, we
 194 retrieve the embeddings by reconstructing lter-specific representations. To ensure quick and efficient
 195 reconstruction, we use a simple one-layer GNN. This on-demand reconstruction of lter-specific
 196 representations significantly reduces the computational and storage requirements associated with
 197 individual node representations. Fig 1 illustrates such a simple encoder's mutual information-based
 198 learning process.

199 Addressing the second issue, we initially train our models to generate low-dimensional embeddings.
 200 These encapsulate latent classes, as discussed in a superset of classes pertinent to downstream
 201 tasks. Although the low-dimensional embeddings harbour latent class information, they lack linear
 202 separability. Hence, we project these embeddings into a higher-dimensional space using random
 203 Fourier feature (RFF) projections, a strategy inspired by kernel methods (Section 4.2). Using this
 204 approach allows for improved linear separability of the latent classes. Our experimental findings (Sec-
 205 tion 6.2) affirm the effectiveness of projecting lower-dimensional embeddings into higher dimensions,
 206 confirming the retention of latent class information in these embeddings.

207 5.2 Supervised Representation Learning

208 After obtaining representations for each lter post the reconstruction of the node representations,
 209 learning an aggregation mechanism to combine information from representations that capture different
 210 parts of the eigenspectrum for the given task is necessary. We adopt learning schemes proposed in
 211 [4, 10, 20], where we learn a weighted combination of lter-specific representations. Therefore, the
 212 combined representations we learn for the downstream task are as follows (considering lters
 213 from the lter bank F):

$$Z = \sum_{i=1}^K \alpha_i \phi(E(X; F_i)) \quad (5)$$

214 The parameters α_i 's are learnable. Additionally, the function $\phi(\cdot)$ represents either the RFF pro-
 215 jection or an identity transformation, depending on whether $E(X; F_i)$ is low-dimensional or not.
 216 A classifier model (e.g. logistic regression) consumes these embeddings, where we train both the
 217 α_i 's and the weights of the classifier. Fig 2 illustrates this process. The main distinction between
 218 semi-supervised methods such as [4, 10] and our method is that the semi-supervised methods
 219 learn both the encoder and the combination coefficients based on labelled data. However, we pre-
 220 train the encoder in our method and subsequently learn a task-specific combination of lter-specific
 221 representations.

Figure 2: Supervised Learning: Using the trained parameters we generate the node embeddings by encoding the iterated graphs that get consumed in the classification task.

6 Experimental Results

Training Details: We define a single-layer graph convolutional network (GCN) with shared weights (θ) across all layers in the layer bank (\mathcal{F}) as our encoder. Therefore, the encoder can be expressed as follows: $E(X; F_i) = \sigma(F_i X)$. It is important to note that F_i represents a normalized layer with self-loops, which ensures that its eigenvalues are within the range of $[0, 2]$. The non-linearity function σ refers to the parametric rectified linear unit (PReLU). As we work with a single graph, we obtain the positive samples by sampling nodes from the graph. Using these sampled nodes, we construct a new adjacency list that only includes the edges between these sampled nodes in layer F_i . On the other hand, the corruption function operates on the same sampled nodes. However, it randomly shuffles the node features instead of perturbing the adjacency list. Similarly, we employ a straightforward readout function that involves averaging the representations across all nodes for a specific layer F_i : $R(H_i) = \frac{1}{N} \sum_{j=0}^N h_j^{F_i}$ where σ denotes the sigmoid non-linearity. We utilize a bilinear scoring function, whose parameters are also shared across all layers:

$$D_l(h_j^{F_i}; h_g^{F_i}) = (h_j^{F_i T} W h_g^{F_i}) \quad (6)$$

We learn the encoder and discriminator parameters by optimizing Eq. 4. While we could use various layer banks, we specifically employ the layer bank corresponding to F_{GPRGNN} for all our experiments. However, we also conduct an ablation study (see 6.5) to compare the performance when using F_{GPRGNN} versus F_{BERNNET} . For more detailed training information, please refer to the supplementary material.

We conducted a series of comprehensive experiments to evaluate the effectiveness and competitiveness of our proposed model compared to SOTA models and methods. These experiments address the following research questions: [RQ1] How does FIGURE perform compared to SOTA unsupervised models? [RQ2] Can we perform satisfactorily even with lower dimensional representations using projections such as RFF? [RQ3] Does shared encoder decrease performance? [RQ4] What is the computational efficiency gained by using lower dimensional representations compared to methods that rely on higher dimensional representations? [RQ5] Can alternative layer banks be employed to recover good quality representations?

Datasets and Setup: We evaluated our model on a diverse set of real-world datasets, which include both heterophilic and homophilic networks, to assess its effectiveness. Similar to previous works, we utilized the node classification task as a proxy to evaluate the quality of the learned representations. Please refer to the supplementary material for detailed information about the benchmark datasets.

The heterophilic datasets used in our evaluation include CHAMELEON, SQUIRREL, ROMAN-EMPIRE, and MINESWEEPER. For CHAMELEON and SQUIRREL, we adopted the ten random splits (with 48%, 32%, and 20% of nodes allocated for the train, validation, and test sets, respectively) from [27]. For ROMAN-EMPIRE and MINESWEEPER, we used the ten random splits provided in [30]. Additionally, we evaluated our model on four homophilic datasets: CORA, CITESEER, and PUBMED, as borrowed from [14]. We report the mean and standard deviation of the test accuracy across different splits. Please refer to the supplementary material for detailed statistics of each dataset.

259 Baselines: In our comparison against baselines, we considered common unsupervised approaches,
 260 such as DEEPWALK and NODE2VEC, and state-of-the-art mutual information-based methods, namely
 261 DGI, MVGRL, GRACE, and SUGRL. We also include the performance numbers of the widely used
 262 GCN for reference. It is important to note that unless explicitly mentioned, we set the representation
 263 size to 512 dimensions for all reported results, consistent with previous work. Please refer to the
 264 supplementary material for detailed comparisons with other supervised methods and the link to our
 265 codebase.

266 6.1 RQ1: FiGURe versus SOTA Methods

Table 1: Contains node classification accuracy percentages on homophilic and heterophilic datasets. FiGURe₃₂ and FiGURe₁₂₈ refer to FiGURe trained with 32 and 128 dimensional representations, respectively, and then projected using RFF. The remaining models are trained at 512 dimensions. Higher numbers indicate better performance. It is worth noting that FiGURe achieves superior performance or remains competitive with the baseline methods in all cases. The rightmost column Av. gain represents the average accuracy gain of FiGURe over the model in that row, averaged across the different datasets. Blue, Red and Green represent 1st, 2nd and 3rd best performing models, for a particular dataset.

	HETEROPHILIC DATASETS				HOMOPHILIC DATASETS			Av. gain
	SQUIRREL	CHAMELEON	ROMAN-EMPIRE	MINESWEEPER	CORA	CITeseer	PUBMED	
DEEPWALK	38.66 (1.44)	53.42 (1.73)	13.08 (0.59)	79.96 (0.08)	83.64 (1.85)	63.66 (3.36)	80.85 (0.44)	0.35
NODE2VEC	42.60 (1.15)	54.23 (2.30)	12.12 (0.30)	80.00 (0.00)	78.19 (1.14)	57.45 (6.44)	73.24 (0.59)	0.56
DGI	39.61 (1.81)	59.28 (1.23)	47.54 (0.76)	82.51 (0.47)	84.57 (1.22)	73.96 (1.61)	86.57 (0.52)	0.67
MVGRL	39.90 (1.39)	54.61 (2.29)	68.50 (0.38)	85.60 (0.35)	86.22 (1.30)	75.02 (1.72)	87.12 (0.35)	0.39
GRACE	53.15 (1.10)	68.25 (1.77)	47.83 (0.53)	80.22 (0.45)	84.79 (1.51)	67.60 (2.01)	87.04 (0.45)	0.54
SUGRL	43.13 (1.36)	58.60 (2.04)	39.40 (0.49)	82.40 (0.58)	81.21 (2.07)	67.50 (1.62)	86.90 (0.59)	0.80
FiGURe ₃₂	48.89 (1.55)	65.66 (2.52)	67.67 (0.77)	85.28 (0.71)	82.56 (0.87)	71.25 (2.20)	84.18 (0.53)	0.18
FiGURe ₁₂₈	48.78 (2.48)	66.03 (2.19)	68.10 (1.09)	85.16 (0.58)	86.14 (1.13)	73.34 (1.91)	85.41 (0.52)	0.11
FiGURe	52.23 (1.19)	68.55 (1.87)	70.99 (0.52)	85.58 (0.49)	87.00 (1.24)	74.77 (2.00)	88.60 (0.40)	0.00

Table 2: Comparison of Node classification accuracy percentages with the widely used supervised model GCN. Despite not having access to task specific labels, FiGURe learns good quality representations.

	SQUIRREL	CHAMELEON	ROMAN-EMPIRE	MINESWEEPER	CORA	CITeseer	PUBMED
GCN	47.78 (2.13)	61.43 (2.70)	73.69 (0.74)	89.75 (0.52)	87.36 (0.91)	76.47 (1.34)	88.41 (0.46)
FiGURe	52.23 (1.19)	68.55 (1.87)	70.99 (0.52)	85.58 (0.49)	87.00 (1.24)	74.77 (2.00)	88.60 (0.44)

267 We analyzed the results in Table 1 and made important observations. Across homophilic and
 268 heterophilic datasets, FiGURe consistently outperforms several SOTA unsupervised models, except
 269 in a few cases where it achieves comparable performance. We want to emphasize the rightmost
 270 column of the table, which shows the average percentage gain in performance across all datasets.
 271 This metric compares the improvement that FiGURe provides over each baseline model for each
 272 dataset and averages these improvements. This metric highlights the performance consistency of
 273 FiGURe across diverse datasets. No other baseline model achieves the same consistent performance
 274 across all datasets as FiGURe. Even the recent state-of-the-art contrastive models GRACE and
 275 SUGRL experience average performance drops of approximately 5% and 10%, respectively. This
 276 result indicates that FiGURe learns representations that exhibit high generalization and task-agnostic
 277 capabilities. Another important observation is the effectiveness of RFF projections in improving lower
 278 dimensional representations. We compared FiGURe at different dimensions, including FiGURe₃₂
 279 and FiGURe₁₂₈, corresponding to learning 32 and 128-dimensional embeddings, respectively, in
 280 addition to the baseline representation size of 512 dimensions. Remarkably, even at lower dimensions,
 281 FiGURe with RFF projections demonstrates competitive performance across datasets, surpassing the
 282 512-dimensional baselines in several cases. This result highlights the effectiveness of RFF projections
 283 in enhancing the quality of lower dimensional representations. Section 6.2 discusses more insights
 284 about the effectiveness of RFF projections. Furthermore, we include the widely used supervised
 285 model, GCN, in Table 2 as a benchmark for comparison. Notably, FiGURe outperforms GCN on
 286 heterophilic datasets, except for ROMAN-EMPIRE and MINESWEEPER, while achieving competitive
 287 performance on homophilic datasets. Please refer to supplementary material for detailed comparisons
 288 with supervised methods.

Table 3: Node classification accuracy percentages with and without using Random Fourier Feature projections (on 32 dimensions). A higher number means better performance. The performance is improved by using RFF in almost all cases, indicating the usefulness of this transformation

	RFF	CORA	CITeseer	SQUIRREL	CHAMELEON
DGI		81.65 (1.90)	65.62 (2.39)	31.60 (2.19)	45.48 (3.02)
	X	81.49 (1.96)	66.50 (2.44)	38.19 (1.52)	56.01 (2.66)
MVGRL		78.81 (1.73)	70.36 (1.76)	29.58 (0.94)	46.56 (2.84)
	X	80.14(2.41)	70.57(1.56)	37.83(1.32)	55.57 (2.28)
SUGRL		65.35 (2.41)	42.84 (2.57)	31.62 (1.47)	43.20 (1.79)
	X	70.06(1.24)	47.03(3.02)	38.50 (2.19)	51.01 (2.26)
GRACE		76.84 (1.09)	58.40 (3.05)	38.20 (1.38)	53.25 (1.58)
	X	79.15 (1.44)	63.66 (2.96)	51.56 (1.39)	67.39 (2.23)
FIGURE		82.88 (1.42)	70.32 (1.98)	39.38 (1.35)	53.27 (2.40)
	X	82.56 (0.87)	71.25 (2.20)	48.89 (1.55)	65.66 (2.52)

In this section, we analyse the performance of unsupervised baselines using 32-dimensional embeddings with and without RFF projections (see Table 3). Despite extensive hyperparameter tuning, we could not replicate the results reported by SUGRL, so we present the best results we obtained. Two noteworthy observations emerge from these tables. Firstly, it is evident that lower dimensional embeddings can yield meaningful and linearly separable representations when combined with simple RFF projections. Utilising RFF projections enhances performance in almost all cases, highlighting the value captured by MI-based methods even with lower-dimensional embeddings. Secondly, FIGURE consistently achieves superior or comparable performance to the baselines, even in lower dimensions. Notably, this includes SUGRL, purported to excel in such settings. However, there is a 2-3% performance gap between GRACE and our method for the SQUIRREL and CHAMELEON datasets. While GRACE handles heterophily well at lower dimensions, its performance deteriorates with homophilic graphs, unlike FIGURE which captures lower frequency information effectively. Additionally, our method exhibits computational efficiency advantages for specific datasets in lower dimensions. Please refer to the supplementary material for more details. Overall, these findings highlight the potential of RFF projections in extracting useful information from lower dimensional embeddings and reaffirm the competitiveness of FIGURE over the baselines.

6.3 RQ3: Sharing Weights Across Filter Specific Encoders

Table 4: A comparison of the performance on the downstream node classification task using independently trained encoders and weight sharing across encoders is shown. The reported metric is accuracy. In both cases, the embeddings are combined using the method described in 5.2

	CORA	CITeseer	SQUIRREL	CHAMELEON
INDEPENDENT	86.92 (1.10)%	75.03 (1.75)%	50.52 (1.51)%	66.86 (1.85)%
SHARED	87.00 (1.24)%	74.77 (2.00)%	52.23 (1.19)%	68.55 (1.87)%

Our method proposes to reduce the computational load by sharing the encoder weights across all filters. It stands to reason whether sharing these weights causes any degradation in performance. We present the results with shared and independent encoders across the filters in Table 4 to verify this. The findings indicate no significant decrease in performance when using shared weights, and in some cases, it even leads to improvements, validating the use of shared encoders.

6.4 RQ4: Computational Efficiency

To assess the computational efficiency of the different methods, we analyzed the computation time and summarized the results in Table 5. The key metric used in this analysis is the mean epoch time: the average time taken to complete one epoch of training. We compared our method with MI-based methods such as DGI and MVGRL. Due to the increase in the number of augmentation views,

Table 5: Mean epoch time (in milliseconds) averaged across 20 trials with different hyperparameters. A lower number means the method is faster. Even though our method is slower at 512 dimensions, using 128 and 32 dimensional embeddings significantly reduces the mean epoch time. Using RFF as described in 6.2 we are able to prevent the performance drops experienced by DGI and MVGRL.

	DGI	MVGRL	FiGUR _e	FiGUR _e ₂₈	FiGUR _e ₂
CORA	38.53 (0.77)	75.29 (0.56)	114.38 (0.51)	20.10 (0.46)	11.54 (0.34)
CITeseer	52.98 (1.15)	102.41 (0.99)	156.24 (0.56)	30.30 (0.60)	17.16 (0.51)
SQUIRREL	87.06 (2.07)	168.24 (2.08)	257.65 (0.76)	47.72 (1.40)	23.52 (1.14)
CHAMELEON	33.08 (0.49)	64.71 (1.05)	98.36 (0.64)	18.56 (0.39)	11.63 (0.48)

317 there is an expected increase in computation time from DGI to MVGRL to FiGUR_e. However,
 318 as demonstrated in 6.2, using RFF projections allows us to achieve competitive performance even
 319 at lower dimensions. Therefore, we also included comparisons with our method at 128 and 32
 320 dimensions in the table. It is evident from the results that our method, both at 128 and 32
 321 dimensions, exhibits faster computation times compared to DGI and MVGRL, which rely on higher-
 322 dimensional representations to achieve good performance. This result indicates that FiGUR_e
 323 is computationally efficient due to its ability to work with lower-dimensional representations. During
 324 training, our method FiGUR_e₃₂ is 3x faster than DGI and 6x times faster than MVGRL.
 325 Despite the faster computation, FiGUR_e₃₂ also exhibits an average performance improvement of
 326 around 2% across the datasets over all methods considered in our experiments. Please refer to the
 327 supplementary material for additional comparisons to other unsupervised models.

328 6.5 RQ5: Experiments on Other Filter Banks

Table 6: Accuracy percentage results using other filter banks. FiGUR_e. F_{BERNNET}³ refers to the
 F_{BERNNET}¹¹ filter bank (Section 4.1) with K set to 3 and F_{BERNNET}¹¹ refers to K set to 11.

	CORA	CITeseer	SQUIRREL	CHAMELEON
F _{BERNNET} ³	85.13 (1.26)	73.38 (1.81)	37.07 (1.29)	53.95 (2.78)
F _{BERNNET} ¹¹	86.62 (1.59)	73.97 (1.43)	43.48 (3.80)	62.13 (3.66)
F _{GPRGNN}	87.00 (1.24)	74.77 (2.00)	52.23 (1.19)	68.55 (1.87)

329 To showcase the versatility of our proposed framework, we conducted an experiment using Bernstein
 330 filters, as detailed in Table 6. The results indicate that using FiGUR_e leads to better performance
 331 than Bernstein filters. We believe that the reason this is happening is due to the latent characteristics
 332 of the dataset. [10, 20] have shown that datasets like CHAMELEON and SQUIRREL need frequency
 333 response functions that give more prominence to the tail-end spectrum. FiGUR_e are more amenable
 334 to these needs, as demonstrated in [20]. However, datasets requiring frequency response similar to
 335 comb filters may be better approximated by F_{BERNNET}¹¹ as their basis gives uniform prominence on
 336 the entire spectrum. Please refer to the supplementary material, which shows the basis frequency
 337 responses of these two filter banks, with more clarification. Therefore, although FiGUR_e gives
 338 better performance for these datasets, there could be datasets where F_{BERNNET}¹¹ could do better. Hence,
 339 we proposed a general framework that can work with any filter bank.

340 7 Conclusion and Future Work

341 Our work demonstrates the benefits of enhancing contrastive learning methods with filter views and
 342 learning filter-specific representations to cater to diverse tasks from homophily to heterophily. We
 343 have effectively alleviated computational and storage burdens by sharing the encoder across these
 344 filters and focusing on low-dimensional embeddings that utilize high-dimensional projections, a
 345 technique inspired by random feature maps developed for kernel approximations. Future directions
 346 include extending the analysis in [20] to graph contrastive learning and explicitly exploring the linear
 347 separability in low dimensions. This analysis could solidify the connection with the proposed random
 348 feature maps approach.

References

- [1] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. *International conference on knowledge discovery & data mining (KDD)*, pages 2623–2631, 2019.
- [2] S. Arora, H. Khandeparkar, M. Khodak, O. Plevrakis, and N. Saunshi. A theoretical analysis of contrastive unsupervised representation learning. *International Conference on Machine Learning (ICML)*, 2019.
- [3] D. Bo, X. Wang, C. Shi, and H.-W. Shen. Beyond low-frequency information in graph convolutional networks. *Association for the Advancement of Artificial Intelligence (AAAI) 2021*.
- [4] E. Chien, J. Peng, P. Li, and O. Milenkovic. Adaptive universal generalized pagerank graph neural network. *International Conference on Learning Representations (ICLR)*, 2021.
- [5] M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain markov process expectations for large time. *Communications on Pure and Applied Mathematics*, 1975.
- [6] C. Dugas, Y. Bengio, F. Bédoune, C. Nadeau, and R. Garcia. Incorporating second-order functional knowledge for better option pricing. *Neural Information Processing Systems (NeurIPS)*, 2000.
- [7] A. Grover and J. Leskovec. node2vec: Scalable feature learning for networks. *International Conference on Knowledge Discovery and Data Mining (KDD)*, 2016.
- [8] K. Hassani and A. H. Khasahmadi. Contrastive multi-view representation learning on graphs. *International Conference on Machine Learning (ICML)*, 2020.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [10] M. He, Z. Wei, Z. Huang, and H. Xu. Bernnet: Learning arbitrary graph spectral filters via bernstein approximation. *Neural Information Processing Systems (NeurIPS)*, 2022.
- [11] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio. Learning deep representations by mutual information estimation and maximization. *International Conference on Learning Representations (ICLR)*, 2019.
- [12] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Neural Information Processing Systems (NeurIPS)*, 2020.
- [13] P. Kar and H. Karnick. Random feature maps for dot product kernels. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2012.
- [14] D. Kim and A. Oh. How to find your friendly neighborhood: Graph attention design with self-supervision. *International Conference on Learning Representations (ICLR)*, 2021.
- [15] G. Kimeldorf and G. Wahba. Some results on tchebycheff and spline functions. *Journal of Mathematical Analysis and Applications*, 1971.
- [16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2015.
- [17] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *International Conference on Learning Representations (ICLR)*, 2017.
- [18] Z. Li, J.-F. Ton, D. Oglic, and D. Sejdinovic. Towards a unified analysis of random fourier features. *Journal of Machine Learning Research (JMLR)*, 2021.
- [19] D. Lim, F. M. Hohne, X. Li, S. L. Huang, V. Gupta, O. P. Bhalerao, and S.-N. Lim. Large scale learning on non-homophilous graphs: New benchmarks and strong simple methods. *Neural Information Processing Systems (NeurIPS)*, 2021.

- 396 [20] V. Lingam, C. Ekbote, M. Sharma, R. Ragesh, A. Iyer, and S. Sellamanickam. A piece-wise
397 polynomial Itering approach for graph neural networks. European Conference on Machine
398 Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)
399 2022.
- 400 [21] R. Linsker. Self-organization in a perceptual network. Computer, 1988.
- 401 [22] M. McPherson, L. Smith-Lovin, and J. M. Cook. Birds of a feather: Homophily in social
402 networks. Annual Review of Sociology, 2001.
- 403 [23] J. Mercer. Functions of positive and negative type, and their connection with the theory
404 of integral equations. Philosophical Transactions of the Royal Society of London. Series A,
405 Containing Papers of a Mathematical or Physical Character, 1909.
- 406 [24] Y. Mo, L. Peng, J. Xu, X. Shi, and X. Zhu. Simple unsupervised graph representation learning.
407 In Association for the Advancement of Arti cial Intelligence (AAAI), 2022.
- 408 [25] S. Nowozin, B. Cseke, and R. Tomioka. f-gan: Training generative neural samplers using
409 variational divergence minimization. Neural Information Processing Systems (NeurIPS)
410 2016.
- 411 [26] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
412 N. Gimsheine, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning
413 library. In Neural Information Processing Systems (NeurIPS), 2019.
- 414 [27] H. Pei, B. Wei, K. C.-C. Chang, Y. Lei, and B. Yang. Geom-gcn: Geometric graph convolutional
415 networks. In International Conference on Learning Representations (ICLR), 2020.
- 416 [28] J. Pennington, F. X. X. Yu, and S. Kumar. Spherical random features for polynomial kernels. In
417 Neural Information Processing Systems (NeurIPS), 2015.
- 418 [29] B. Perozzi, R. Al-Rfou, and S. Skiena. DeepWalk. In International Conference on Knowledge
419 Discovery and Data Mining (KDD), 2014.
- 420 [30] O. Platonov, D. Kuznedelev, M. Diskin, A. Babenko, and L. Prokhorenkova. A critical look at
421 the evaluation of GNNs under heterophily: Are we really making progress? International
422 Conference on Learning Representations (ICLR), 2023.
- 423 [31] A. Rahimi and B. Recht. Random features for large-scale kernel machines. Neural Informa-
424 tion Processing Systems (NeurIPS), 2007.
- 425 [32] A. Rahimi and B. Recht. Weighted sums of random kitchen sinks: Replacing minimization with
426 randomization in learning. In Neural Information Processing Systems (NeurIPS), 2008.
- 427 [33] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst. The emerging field
428 of signal processing on graphs: Extending high-dimensional data analysis to networks and other
429 irregular domains. IEEE Signal Processing Magazine, 2013.
- 430 [34] N. Tremblay, P. Gonçalves, and P. Borgnat. Design of graph Iters and Iterback. In Cooperative
431 and Graph Signal Processing, 2017.
- 432 [35] P. Velicković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm. Deep graph
433 infomax. In International Conference on Learning Representations (ICLR), 2019.
- 434 [36] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang. Deep Graph Contrastive Representation
435 Learning. In ICML Workshop on Graph Representation Learning and Beyond, 2020.

436 **8 Supplementary Material**

437 **Contents**

438	1 Introduction	1
439	2 Related Work	2
440	3 Problem Setting	2
441	4 Preliminaries	3
442	4.1 Filter Banks	3
443	4.2 Random Feature Maps for Kernel Approximations	3
444	5 Proposed Approach	3
445	5.1 Unsupervised Representation Learning	4
446	5.2 Supervised Representation Learning	5
447	6 Experimental Results	6
448	6.1 RQ1: FiGURe versus SOTA Methods	7
449	6.2 RQ2: RFF Projections on Lower Dimensional Representations	8
450	6.3 RQ3: Sharing Weights Across Filter Specific Encoders	8
451	6.4 RQ4: Computational Efficiency	8
452	6.5 RQ5: Experiments on Other Filter Banks	9
453	7 Conclusion and Future Work	9
454	8 Supplementary Material	12
455	8.1 Reproducibility	12
456	8.2 Datasets	13
457	8.3 Training Details	13
458	8.3.1 Unsupervised Training	14
459	8.3.2 Supervised Training	14
460	8.3.3 Negative Sampling for the Identity Filter	15
461	8.4 Evaluation on large graphs	15
462	8.5 Comparison with other Supervised Methods	15
463	8.6 RFF Projections	16
464	8.7 Computational Comparisons with other Other Unsupervised Methods	16
465	8.8 Choice of Filter Banks	17
466	8.9 Visualising RFF Behavior and Community Structure	18

467 **8.1 Reproducibility**

468 We strive to ensure the reproducibility of our research findings. To facilitate this, we provide the
469 details of our experimental setup, including dataset sources, preprocessing steps, hyperparameters,

470 and model configurations. We also make our code and the datasets used, publicly available at this
 471 [LINK](#), enabling researchers to reproduce our results and build upon our work. We would like to
 472 emphasize that our code is built on top of the existing GRL codebase. For the datasets used
 473 in our evaluation, we provide references to their original sources and any specific data splits that
 474 we employed. This allows others to obtain the same datasets and perform their own analyses using
 475 consistent data. Additionally, we specify the versions of libraries and frameworks used in our
 476 experiments, in Section 8.3, and in the REQUIREMENTS file and the README file, in the codebase,
 477 enabling others to set up a compatible environment. We document any specific seed values or
 478 randomization procedures that may affect the results. By providing these details and resources,
 479 we aim to promote transparency and reproducibility in scientific research. We encourage fellow
 480 researchers to reach out to us if they have any questions or need further clarification on our methods
 481 or results.

482 8.2 Datasets

483 **Homophilic Datasets:** We evaluated our model (as well as baselines) on three homophilic datasets:
 484 CORA, CITESEER, and PUBMED as borrowed from [14]. All three are citation networks, where
 485 each node represents a research paper and the links represent citations. Pubmed consists of medical
 486 research papers. The task is to predict the category of the research paper. We follow the same dataset
 487 setup mentioned in [14] to create 10 random splits for each of these datasets.

488 **Heterophilic Datasets:** In our evaluation, we included four heterophilic datasets: CHAMELEON,
 489 SQUIRREL, ROMAN-EMPIRE, and MINESWEEPER. For CHAMELEON and SQUIRREL, nodes represent
 490 Wikipedia web pages and edges capture mutual links between pages. We utilized the ten random
 491 splits provided in [27], where 48%, 32%, and 20% of the nodes were allocated for the train, validation,
 492 and test sets, respectively. For ROMAN-EMPIRE, each node corresponds to a word in the Roman Empire
 493 Wikipedia article. Two words are connected with an edge if either these words follow each other
 494 in the text, or they are connected in the dependency tree of the sentence. The syntactic role of the
 495 word/node defines its class label. The MINESWEEPER graph is a regular 100x100 grid where each
 496 node is connected to eight neighboring nodes, and the features are one-hot encoded representations
 497 of the number of neighboring mines. The task is to predict which nodes are mines. For both
 498 ROMAN-EMPIRE and MINESWEEPER, we used the ten random splits provided in [30].

499 **Large Datasets:** We also evaluate our method on two large datasets: OGBN-ARXIV (from [12])
 500 and ARXIV-YEAR (from [19]). Both these datasets are from the arxiv citation network. OGBN-
 501 ARXIV, the task is to predict the category of the research paper, and ARXIV-YEAR the task is
 502 to predict the year of publishing. We use the publicly available splits of OGBN-ARXIV [14] and
 503 follow the same dataset setup mentioned in [19] to generate 5 random splits for ARXIV-YEAR. Note
 504 that OGBN-ARXIV is a homophilic dataset while ARXIV-YEAR is a heterophilic datasets.

505 The detailed dataset statistics can be found in Table 7.

Table 7: Dataset Statistics. The table provides information on the following dataset characteristics: number of nodes, number of edges, feature dimension, number of classes, as well as the count of nodes used for training, validation, and testing.

PROPERTIES	HETEROPHILIC DATASETS					HOMOPHILIC DATASETS				
	SQUIRREL	CHAMELEON	ROMAN-EMPIRE	MINESWEEPER	ARXIV-YEAR	OGBN-ARXIV	CITeseer	PUBMED	CORA	
#NODES	5201	2277	22662	10000	169343	169343	3327	19717	2708	
#EDGES	222134	38328	32927	39402	1166243	1335586	12431	108365	13264	
#FEATURES	2089	500	300	7	128	128	3703	500	1433	
#CLASSES	5	5	18	2	5	40	6	3	7	
#TRAIN	2496	1092	11331	5000	84671	90941	1596	9463	1192	
#VAL	1664	729	5665	2500	42335	29799	1065	6310	796	
#TEST	1041	456	5666	2500	42337	48603	666	3944	497	

506 8.3 Training Details

507 We conducted all experiments on a machine equipped with an Intel(R) Xeon(R) CPU E5-2690 v4 @
 508 2.60GHz processor, 440GB RAM, and a Tesla-P100 GPU with 16GB of memory. The experiments
 509 were executed using Python 3.9.12 and PyTorch 1.12.0. To optimize the hyperparameter search,
 510 we employed Optuna [1]. We utilized the Adam optimizer [16] for the optimization process.

511 8.3.1 Unsupervised Training

512 We conducted hyperparameter tuning for all unsupervised methods using 20 Optuna trials. The
513 hyperparameter ranges and settings for each method are as follows:

514 DEEPWALK: We set the learning rate to 0.001, number of epochs to 20 and the varied the random
515 walk length over {8; 9; 10; 11; 12}. Additionally, we varied the context window size over {3; 4; 5}
516 and the negative size (number of negative samples per positive sample) over {0.5; 1; 2}.

517 NODE2VEC: For Node2Vec, we set the learning rate to 0.01 and number of epochs to 100. We varied
518 the number of walks over {5; 10; 15} and the walk length over {40; 50; 60}. The α (return parameter)
519 value was chosen from {0.1; 0.25; 0.5; 1} and q (in-out parameter) value was chosen from {0.5; 1; 2}.

520 DGI: DGI [35] proposes a self-supervised learning framework for graph representation learning by
521 maximizing the mutual information between local and global structural context of nodes, enabling
522 unsupervised feature extraction in graph neural networks. We relied on the authors' code and the
523 prescribed hyperparameter ranges specific to the DGI model, for our experiments.

524 MVGRL: MVGRL [8] proposes a method for learning unsupervised node representations by
525 leveraging two views of the graph data, the graph diffusion view and adjacency graph view. We relied
526 on the authors' code and the prescribed hyperparameter ranges specific to the MVGRL model, for
527 our experiments.

528 GRACE: GRACE [36] proposes a technique where two different perspectives of the graph are
529 created through corruption, and the learning process involves maximizing the consistency between
530 the node representations obtained from these two views. We relied on the authors' code and the
531 prescribed hyperparameter ranges specific to the GRACE model, for our experiments.

532 SUGRL: SUGRL [24] proposes a technique for learning unsupervised representations which capture
533 node proximity, while also utilising node feature information. We relied on the authors' code and
534 the prescribed hyperparameter ranges specific to the SUGRL model, for our experiments.

535 FiGURE: We followed the setting of the MVGRL model, setting the batch size to 2 and number of
536 GCN layers to 1. We further tuned the learning rate over {0.0001; 0.0001; 0.001; 0.01; 0.1} and the
537 sample size (number of nodes selected per batch) over {50; 100; 175; 200; 225; 250}.

538 In each case, we selected the hyperparameters that resulted in the lowest unsupervised training loss.

539 8.3.2 Supervised Training

540 For all unsupervised methods, including the baselines and our method, we perform post-training
541 supervised evaluation using logistic regression with 20 Optuna trials. We set the maximum number
542 of epochs to 1000 and select the epoch and hyperparameters that yield the best validation accuracy.
543 The learning rate is swept over the range {0.0001; 0.0001; 0.001; 0.0015; 0.01; 0.015; 0.1; 0.5; 1; 2},
544 and the weight decay is varied over { 10^{-5} ; 10^{-4} ; 10^{-3} ; 10^{-2} ; 10^{-1} ; 0; 0.5; 1; 3}.

545 FiGURE: Along with the hyperparameters described above, following the approach described in [10],
546 we also tune the combination coefficients (β) with a separate learning rate. This separate learning
547 rate is swept over the range {0.0001; 0.0001; 0.001; 0.0015; 0.01; 0.015; 0.1; 0.5; 1; 2}. In addition,
548 we have a coefficient for masking the incoming embeddings from each layer, which is varied between
549 0 and 1. Furthermore, these coefficients are passed through an activation layer, and we have two
550 options: 'none' and 'exp'. When 'none' is selected, the coefficients are used directly, while 'exp'
551 indicates that they are passed through an exponential function before being used.

552 FiGURE with RFF: For the experiments involving Random Fourier Features (RFF), we use
553 the same hyperparameter ranges as mentioned above. However, we also tune the gamma pa-
554 rameter which is specific to RFF projections. The gamma parameter is tuned within the range
555 of {0.1; 0.2; 0.3; 0.4; 0.5; 0.6; 0.7; 0.8; 0.9; 1.0; 1.1; 1.2}.

¹<https://github.com/PetarV-/DGI.git>

²<https://github.com/kavehassani/mvgrl.git>

³<https://github.com/CRIPAC-DIG/GRACE.git>

⁴<https://github.com/YujieMo/SUGRL.git>

8.3.3 Negative Sampling for the Identity Filter

In our implementation of F_{GPRGNN} or F_{BERNNET} , we follow a specific procedure for handling the filters during training and evaluation. For all filters except the identity filter, we employ the negative sampling approach described in Section 6. However, the identity filter is treated differently. During training, we exclude the identity filter and only include it during evaluation.

During negative sampling, the generation of the negative anchor involves shuffling the node features, followed by premultiplying the shuffled node feature matrix with the filter matrix and computing the mean. On the other hand, for the positive anchor, the same procedure is applied without shuffling the node features. This approach encourages the model to learn meaningful patterns and relationships in the data when the filter matrix is not the identity matrix.

The decision to exclude the identity filter during training is based on the observation that it presents a special case where the positive and negative anchors become the same. As a result, the model would optimize and minimize the same quantity, potentially leading to trivial solutions. To prevent this, we exclude the identity filter during training.

By excluding the identity filter during training, we ensure that the model focuses on the other filters in F_{GPRGNN} or F_{BERNNET} to capture and leverage the diverse information present in the graph. Including the identity filter only during evaluation allows us to evaluate its contribution to the final performance of the model. This approach helps prevent the model from learning trivial solutions and ensures that it learns meaningful representations by leveraging the other filters.

8.4 Evaluation on large graphs

Table 8: Contains node classification accuracy percentages on two large-scale datasets. ARXIV and ARXIV-YEAR have been added. FIGURE₃₂ and FIGURE₁₂₈ refer to FIGURE trained with 32 and 128 dimensional representations, respectively, and then projected using RFF. The remaining models are trained at 512 dimensions. Higher numbers indicate better performance. Blue, Red and Green represent the 1st, 2nd and 3rd best performing models, for a particular dataset.

	HETEROPHILIC DATASETS					HOMOPHILIC DATASETS				Av. gain
	SQUIRREL	CHAMELEON	ROMANEMPIRE	MINESWEEPER	ARXIV-YEAR	CORA	CITSEER	PUBMED	OGBN-ARXIV	
DGI	39.61 (1.81)	59.28 (1.23)	47.54 (0.76)	82.51 (0.47)	40.59 (0.09)	84.57 (1.22)	73.96 (1.61)	86.57 (0.52)	65.58 (0.06)	61
MVGRL	39.90 (1.39)	54.61 (2.29)	68.50 (0.38)	85.60 (0.35)	OOM	86.22 (1.30)	75.02 (1.72)	87.12 (0.35)	OOM	4.39
GRACE	53.15 (1.10)	68.25 (1.77)	47.83 (0.53)	80.22 (0.45)	OOM	84.79 (1.51)	67.60 (2.01)	87.04 (0.43)	OOM	5.54
SUGRL	43.13 (1.36)	58.60 (2.04)	39.40 (0.49)	82.40 (0.58)	36.96 (0.19)	81.21 (2.07)	67.50 (1.62)	86.90 (0.54)	65.80 (0.06)	64
FIGURE ₃₂	48.89 (1.55)	65.66 (2.52)	67.67 (0.77)	85.28 (0.71)	41.30 (0.21)	82.56 (0.87)	71.25 (2.20)	84.18 (0.53)	66.58 (0.08)	18
FIGURE ₁₂₈	48.78 (2.48)	66.03 (2.19)	68.10 (1.09)	85.16 (0.58)	41.94 (0.15)	86.14 (1.13)	73.34 (1.91)	85.41 (0.52)	69.11 (0.02)	11
FIGURE	52.23 (1.19)	68.55 (1.87)	70.99 (0.52)	85.58 (0.49)	42.26 (0.20)	87.00 (1.24)	74.77 (2.00)	88.60 (0.44)	69.68 (0.00)	0

Similar to Table 1, Table 8, provides a comparison with GOTA methods such as DGI, MVGRL, GRACE and SUGRL. However, in this table we also incorporate the large-scale datasets ARXIV-YEAR and OGBN-ARXIV. FIGURE shows good performance on these datasets as well, demonstrating the scalability of our method. The last column, average percentage gain, is updated accordingly. Two baseline methods MVGRL and GRACE run into memory issues on the larger datasets and are accordingly reported OOM in the table. Even the lower-dimensional representations (with RFF projections) are able to beat the baselines on these large scale datasets. FIGURE is consistently able to provide gains over the baselines methods regardless of the kind of graph, homophilic, heterophilic or large-scale. This once again demonstrates the generalizability of FIGURE. It is noteworthy that computational efficiency gained by reducing the dimension size becomes significant with the scale of the dataset. ARXIV-YEAR for example, 128 dimensional embeddings give 1.6x speedup and 32 dimensional embeddings give 1.7x speedup.

8.5 Comparison with other Supervised Methods

Table 9 presents a comparison with common supervised baselines. Specifically, we choose 3 models for comparison, representing three different kinds of supervised methods, standard aggregation models (GCN), spectral filter-based models (SPRGNN) and smart-aggregation models (SAGCN). There are two key observations from this table. Firstly, FIGURE is competitive with the supervised baselines, lagging behind only by a few percentage points in some cases. This suggests that much of the information that is required by the downstream tasks, captured by the supervised models, can be made available through unsupervised methods. FIGURE which uses filter banks. It is important

596 to note that in FiGURe we only utilize logistic regression while evaluating on the downstream task.
 597 This is much more efficient than training a graph neural network end to end. Additionally it is possible
 598 that further gains may be obtained by utilizing a non-linear model like an MLP.

599 Furthermore, as indicated by 9, we can gain further computational efficiency by utilizing lower
 600 dimensional representations like 32 and 128 (with RFF), and still not compromise significantly on
 601 the performance.

602 Overall FiGURe manages to remain competitive despite not having access to task-specific labels and
 603 is computationally efficient as well.

Table 9: Contains node classification accuracy percentages on heterophilic and homophilic datasets. GCN, GPRGNN and H₂GCN are supervised methods. FiGURe₃₂ and FiGURe₁₂₈ refer to FiGURe trained with 32 and 128 dimensional representations, respectively, and then projected using RFF. The remaining models are trained at 512 dimensions. Higher numbers indicate better performance.

	HETEROPHILIC DATASETS					HOMOPHILIC DATASETS				
	SQUIRREL	CHAMELEON	ROMAN-EMPIRE	MINESWEEPER	ARXIV-YEAR	OGBN-ARXIV	CORA	CITSEER	PUBMED	
GCN	47.78 (2.13)	62.83 (1.52)	73.69 (0.74)	89.75 (0.52)	46.02 (0.26)	69.37 (0.00)	87.36 (0.91)	76.47 (1.34)	88.41 (0.46)	
GPRGNN	46.31 (2.46)	62.59 (2.04)	64.85 (0.27)	86.24 (0.61)	45.07 (0.21)	68.44 (0.00)	87.77 (1.31)	76.84 (1.69)	89.08 (0.39)	
H ₂ GCN	37.90 (2.02)	58.40 (2.77)	60.11 (0.52)	89.71 (0.31)	49.09 (0.10)	OOM	87.81 (1.35)	77.07 (1.64)	89.59 (0.33)	
FiGURe ₃₂	48.89 (1.55)	65.66 (2.52)	67.67 (0.77)	85.28 (0.71)	41.30 (0.21)	66.58 (0.00)	82.56 (0.87)	71.25 (2.20)	84.18 (0.53)	
FiGURe ₁₂₈	48.78 (2.48)	66.03 (2.19)	68.10 (1.09)	85.16 (0.58)	41.94 (0.15)	69.11 (0.00)	86.14 (1.13)	73.34 (1.91)	85.41 (0.52)	
FiGURe	52.23 (1.19)	68.55 (1.87)	70.99 (0.52)	85.58 (0.49)	42.26 (0.20)	69.69 (0.00)	87.00 (1.24)	74.77 (2.00)	88.60 (0.44)	

604 8.6 RFF Projections

605 As shown in Section 6.2 and in Section 6.1, RFF projections are a computationally efficient way to
 606 achieve training by preserving the latent class behavior present in lower dimensional embeddings,
 607 by projecting them into a higher dimensional linearly separable space. The natural question that
 608 comes up is how do we compute the RFF projections? We provide an algorithm to compute the
 609 RFF projections in this section, in algorithm 1. Note that this follows [31].

Algorithm 1 Random Fourier Feature Computation

Require: Input data $X \in \mathbb{R}^{N \times d}$, target dimension D , kernel bandwidth σ
 Ensure: Random Fourier Features $Z \in \mathbb{R}^{N \times D}$
 1: Initialize random weight matrix $W \in \mathbb{R}^{d \times D}$ with Gaussian distribution
 2: Initialize random bias vector $b \in \mathbb{R}^D$ uniformly from $[0, 2\sigma]$
 3: Compute scaled input $X^0 = XW + b$
 4: Compute random Fourier features $Z = \frac{2}{D} \cos(X^0)$
 5: return Z

610 8.7 Computational Comparisons with other Other Unsupervised Methods

Table 10: Mean epoch time (in milliseconds) averaged across 20 trials with different hyperparameters. A lower number means the method is faster. Even though our method is slower at 512 dimensions, using 128 and 32 dimensional embeddings significantly reduces the mean epoch time. Using RFF as described in 6.2 we are able to prevent the performance drops experienced by SUGRL and GRACE.

	SUGRL	GRACE	FiGURe	FiGURe ₃₂	FiGURe ₁₂₈
CORA	15.92 (4.10)	51.19 (6.8)	114.38 (0.51)	20.10 (0.46)	11.54 (0.34)
CITSEER	24.37 (4.92)	77.16 (7.2)	156.24 (0.56)	30.30 (0.60)	17.16 (0.51)
SQUIRREL	33.63 (6.94)	355.2 (67.34)	257.65 (0.76)	47.72 (1.40)	23.52 (1.14)
CHAMELEON	16.91 (5.90)	85.05 (14.1)	98.36 (0.64)	18.56 (0.39)	11.63 (0.48)

611 In Section 6.4, we compared the computational time of FiGURe with MVGRL and DGI, as all
 612 three methods fall under the category of unsupervised methods that perform contrastive learning
 613 with representations of the entire graph. However, there is another class of methods, such as

614 SUGRL and GRACE, that contrast against other nodes without the need for graph representation
 615 computation. Consequently, these methods exhibit higher computational efficiency. Hence, as shown
 616 in Table 10 upon initial inspection, it appears that SUGRL (at 512 dimensions) exhibits the highest
 617 computational efficiency, even outperforming FIGURE₁₂₈. However, despite its computational
 618 efficiency, the significant drop in performance across datasets (as discussed in Section 6.1) renders
 619 it less favorable for consideration. In fact, FIGURE₃₂ offers computational cost savings compared
 620 to SUGRL, while also achieving significantly better downstream classification accuracy. Turning
 621 to GRACE, it demonstrates greater computational efficiency than FIGURE (at 512 dimensions) for
 622 low to medium-sized graphs. However, as the graph size increases, due to random node feature
 623 level masking and edge level masking, the computational requirements of GRACE substantially
 624 increase (as evidenced by the results on SQUIRREL). Therefore, for larger graphs with more than
 625 approximately 5000 nodes, FIGURE proves to be more computationally efficient than GRACE (even
 626 at 512 dimensions). Furthermore, considering the performance improvements exhibited by FIGURE,
 627 it is evident that FIGURE (combined with RFF projections) emerges as the preferred method for
 628 unsupervised contrastive learning in graph data.

629 8.8 Choice of Filter Banks

630 In Section 4.1, we explore the flexibility of FIGURE to accommodate various filter banks. When
 631 making a choice, it is crucial to examine the intrinsic properties of the filters contained within different
 632 filter banks. We pick two filter banks F_{BERNNET} and F_{GPRGNN} and provide an overview of the filters
 633 contained in the filter banks. We use these two filter banks as examples to illustrate what should one
 634 be looking for, while choosing a filter bank.

635 **Bernstein Polynomials** Figure 3 illustrates that as the number of Bernstein Basis increases, the
 636 focus on different parts of the eigenspectrum also undergoes changes. With an increase in polynomial
 637 order, two notable effects can be observed. Firstly, the number of filters increases, enabling each
 638 filter to focus on more fine-grained eigenvalues. This expanded set of polynomial filters allows for a
 639 more detailed examination of the eigenspectrum. Secondly, if we examine the first and last Bernstein
 640 polynomials, we observe an outward shift in their shape. This shift results in the enhancement of a
 641 specific fine-grained part at the ends of the spectrum. These observations demonstrate that Bernstein
 642 polynomials offer the capability to selectively target and enhance specific regions of interest within
 643 the eigenspectrum.

644 **Standard Basis** Figure 3 reveals two key observations. Firstly, at a polynomial order of 2, the
 645 standard basis exhibit focus at the ends of the spectrum, in contrast to the behavior of Bernstein
 646 polynomials, which tend to concentrate more on the middle of the eigenspectrum. This discrepancy
 647 highlights the distinct characteristics and emphasis of different polynomial bases in capturing different
 648 parts of the eigenspectrum. Secondly, as the number of polynomials increases (in contrast to
 649 Bernstein polynomials), the lower order polynomials remain relatively unchanged. Instead, additional
 650 polynomials are introduced, offering a more fine-grained focus at the ends of the spectrum. This
 651 expansion of polynomials allows for a more detailed exploration of specific regions of interest within
 652 the the ends of eigenspectrum.

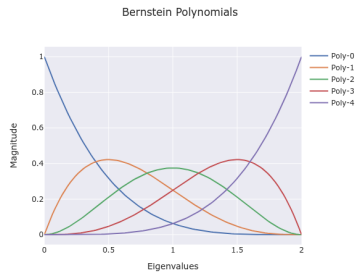
653 In the context of filter banks, previous studies [20, 4] have demonstrated that certain datasets, such as
 654 SQUIRREL and CHAMELEON, benefit from frequency response functions that enhance the tail ends
 655 of the eigenspectrum. This observation suggests that the standard basis, which naturally focuses
 656 on the ends of the spectrum, may outperform Bernstein basis functions at lower orders. However,
 657 as the order of the Bernstein basis increases, as discussed in 4.1, there is a notable improvement in
 658 performance. This can be attributed to the increased focus of Bernstein basis functions on specific
 659 regions, particularly the ends of the spectrum. As a result, higher-order Bernstein filters exhibit
 660 enhanced capability in capturing important information in those regions. It is worth noting that the
 661 choice between F_{GPRGNN} and F_{BERNNET} depends on the specific requirements of the downstream
 662 task. If the task necessitates a stronger focus on the middle of the spectrum or requires a band-pass or
 663 comb-like frequency response, F_{BERNNET} is likely to outperform F_{GPRGNN} . Thus, the selection of the
 664 appropriate filter bank should be based on the desired emphasis on different parts of the eigenspectrum.
 665 Regarding the performance comparison between F_{BERNNET} and F_{GPRGNN} , it is plausible that as we
 666 increase the order of the Bernstein basis, the performance could potentially match that of F_{GPRGNN} .
 667 However, further investigation and experimentation are required to determine the specific conditions
 668 and orders at which this convergence in performance occurs.

669 8.9 Visualising RFF Behavior and Community Structure

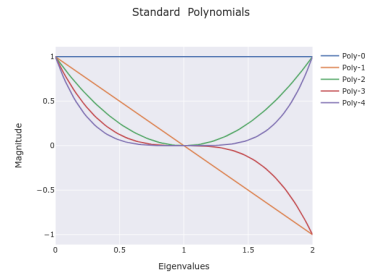
670 As shown in prior sections, FiGURe improves on both computational efficiency as well as perfor-
671 mance by utilising RFF projections. In this section, we aim to gain insights into the behavior of
672 RFF projections and comprehend their underlying operations through a series of simple visualizations.

673 **t-SNE Plots:** Figure 4 offers insights into the structure of the embeddings for the CORA dataset
674 across different dimensions. Remarkably, even at lower dimensions (e.g., 32 dimensions), clear
675 class structures are discernible, indicating that the embeddings capture meaningful information
676 related to the class labels. Furthermore, when employing RFF to project the embeddings into higher
677 dimensions, these distinct class structures are still preserved. This suggests that the role of RFF is not
678 to introduce new information, but rather to enhance the suitability of lower-dimensional embeddings
679 for linear classifiers while maintaining the underlying class-related information. Notably, even at
680 512 dimensions, the class structures remain distinguishable. However, it is worth noting that the
681 class-specific embeddings appear to be more tightly clustered and less dispersed compared to the
682 32-dimensional embeddings or the projected 32-dimensional embeddings. This suggests that learning
683 a 512-dimensional embedding differs inherently from learning a 32-dimensional embedding and
684 subsequently projecting it into higher dimensions.

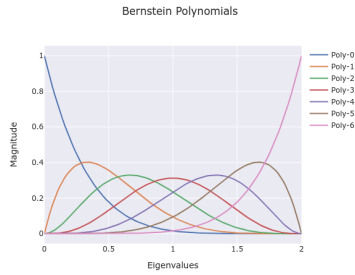
685 **Correlation Plots:** Figure 5 offers insights into the correlation patterns within the embeddings gener-
686 ated from the SQUIRREL dataset across different dimensions. In lower dimensions, the embeddings
687 exhibit high correlation with each other, which can be attributed to the presence of a mixture of
688 topics or latent classes within the dataset. However, when the embeddings are projected to higher
689 dimensions using RFF, the correlation is reduced, and a block diagonal matrix emerges. This block
690 diagonal structure indicates the presence of distinct classes or communities within the dataset. Even at
691 512 dimensions, a more refined block diagonal structure can be observed compared to the correlation
692 matrix of the 32-dimensional embeddings. Furthermore, it is noteworthy that the correlation of
693 the projected embeddings can be regarded as a sparser version of the correlation observed in the
694 512-dimensional embeddings.



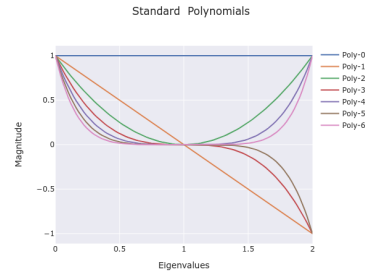
(a) Five Bernstein Basis



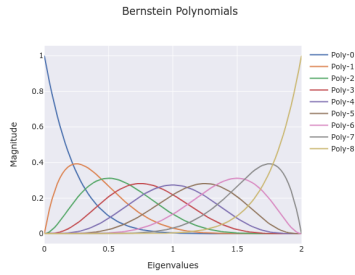
(b) Five Standard Basis



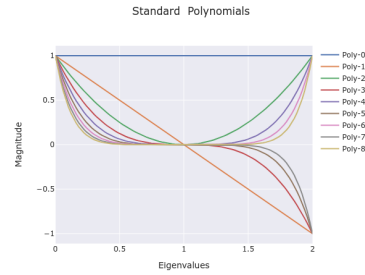
(c) Seven Bernstein Basis



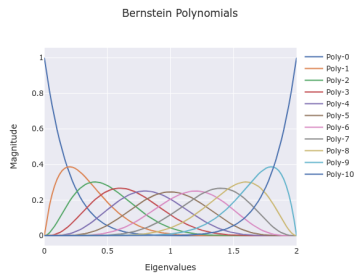
(d) Seven Standard Basis



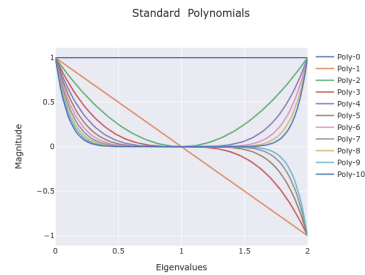
(e) Nine Bernstein Basis



(f) Nine Standard Basis



(g) Eleven Bernstein Basis



(h) Eleven Standard Basis

Figure 3: The figures contain the Bernstein basis as well as standard basis for different degrees. The x-axis of the figures represents the eigenvalues of the Laplacian matrix, while the y-axis represents the magnitude of the polynomials. It is important to note that while plotting the standard polynomials, they are computed with respect to the Laplacian matrix (L_n) rather than the adjacency matrix. As a result, the eigenvalues lie between $[0; 2]$. On the other hand, the Bernstein polynomials are typically defined for the normalised Laplacian matrix, and therefore there is no change in the eigenvalue range (the eigenvalues of the normalised Laplacian matrix typically range from 0 to 2). By using the Laplacian matrix as the basis for plotting the polynomials, we can observe the behavior and magnitude of the polynomials at different eigenvalues, providing insights into their spectral properties and frequency response characteristics.

