

# COMPUTATIONALLY-EFFICIENT GRAPH MODELING WITH REFINED GRAPH RANDOM FEATURES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We propose *refined GRFs* (GRFs++), a new class of *Graph Random Features* (GRFs) for efficient and accurate computations involving kernels defined on the nodes of a graph. GRFs++ resolve some of the long-standing limitations of regular GRFs, including difficulty modeling relationships between more distant nodes. They reduce dependence on sampling long graph random walks via a novel *walk-stitching* technique, concatenating several shorter walks without breaking unbiasedness. By applying these techniques, GRFs++ inherit the approximation quality provided by longer walks but with greater efficiency, trading sequential, inefficient sampling of a long walk for parallel computation of short walks and matrix-matrix multiplication. Furthermore, GRFs++ extend the simplistic GRFs walk termination mechanism (Bernoulli schemes with fixed halting probabilities) to a broader class of strategies, applying general distributions on the walks' lengths. This improves the approximation accuracy of graph kernels, without incurring extra computational cost. We provide empirical evaluations to showcase all our claims and complement our results with theoretical analysis.

## 1 INTRODUCTION & RELATED WORK

Graph modeling plays an important role in several applications of machine learning (ML), such as anomaly, community and fraud detection (Kim et al., 2022; 2026; Beutel et al., 2015; Noble & Cook, 2003; Li et al., 2025; Dong et al., 2025; Chen et al., 2024; Liu et al., 2023), recommender systems (Yang et al., 2025; 2023; Deng et al., 2022; Gao et al., 2023), and computational biology (Mao et al., 2024; Banerjee & Jost, 2009; Zhang et al., 2024). As for Euclidean data, research on graph modeling has spurred the development of many hard-coded/learnable or parameterized (Yanardag & Vishwanathan, 2015) classes of *kernels* (similarity functions), defining relationships between nodes in the graph (Smola & Kondor, 2003; Kondor & Lafferty, 2002) or between the graphs themselves (Vishwanathan et al., 2010; Shervashidze et al., 2009).

In this paper, we focus on graph node kernels  $K : V(G) \times V(G) \rightarrow \mathbb{R}$  defined on the vertices  $V$  of a given graph  $G$ , where the similarity between the nodes is measured via their relationship in the graph – e.g how well-connected the nodes are. Common examples of such kernels include the *d-regularized Laplacian*, *diffusion process*, *p-step random walk*, and *inverse cosine* kernels (Smola & Kondor, 2003; Choromanski, 2023). Computing the corresponding Gram matrices  $\mathbf{K}(G) = [K(v_i, v_j)]_{i,j=1}^N$  for  $v_1, \dots, v_N \in V(G)$  tends to be expensive, since this often requires operations of time complexity cubic in the number of graph nodes  $N$ . For this reason, research has been dedicated to developing efficient approximation strategies. One common approach is to rewrite the graph kernel as a product of two lower rank matrices, linearizing the graph kernel values with some mapping  $\phi : V \rightarrow \mathbb{R}^m$  as follows:

$$\widehat{\mathbf{K}}(v_i, v_j) = \phi(v_i)^\top \phi(v_j). \quad (1)$$

This low-rank factorization unlocks efficient computations with the corresponding kernel matrices. In particular, matrix-multiplication operations no longer require explicit materialization of  $\mathbf{K}$ , exploiting the associativity of matrix multiplication. However, until recently this approach was restricted to ad-hoc learnable graph kernels defined implicitly via learnable  $\phi$  (Wu et al., 2019), rather than approximations of the specific classes listed above.

A series of recent papers proposed a new mechanism called *Graph Random Features* (GRFs) (Choromanski, 2023; Reid et al., 2023; 2024b;a; 2025). GRFs provide unbiased approximation of the classes of graph kernels listed above, with probabilistic mappings  $\phi$  obtained via graph random walks. For every graph node  $v$ , GRFs build a scalar field on the subset of RW-reachable graph nodes  $V(G)$  via incremental (kernel-dependent) updates of the field in the visited nodes. This field is then mapped to a node-embedding  $\phi(v)$ , encoding the relationship of the node to the entire graph  $G$ . Since  $\phi(v)$  is probabilistic, it is referred to as the *graph random feature* corresponding to  $v$ . Though originally introduced to approximate kernels defined between pairs of graph nodes, GRFs were recently *lifted* for unbiased approximation of kernels defined between pairs of graphs (Choromanski et al., 2025).

In this paper, we propose a new class of GRF for efficient and accurate computations involving graph kernels defined on the nodes of the graph, that we refer to as *refined GRFs*, or GRFs++. GRFs++ resolve some of the long-standing challenges that regular GRFs face, such as: difficulty in modeling relationships between more distant nodes. They also reduce the dependence on the longer random walks in the graphs, the workhorse mechanism of the regular GRFs. This is done via the newly-proposed *walk-stitching* technique, where several shorter walks are concatenated to emulate the mechanism of conducting longer random walks. By applying this techniques, GRFs++ inherit the approximation quality provided by longer walks, yet via a much more efficient method, effectively trading sequential and less computationally-efficient mechanism of conducting a long walk for a parallel computation of short walks and matrix-matrix multiplications. Furthermore, GRFs++ extend simplistic GRFs’ walk-termination mechanism leveraging standard Bernoulli schemes with fixed halting probabilities into a class of strategies applying general distributions on the walks’ lengths and maintaining unbiasedness of regular GRFs. This leads to more accurate approximation of the graph kernels under consideration, and with no extra computational cost. We provide empirical evaluations, showcasing all the claims, and complement our results with the theoretical analysis.

This paper is organized as follows:

1. In Sec. 2, we present the refined GRFs++ mechanism, introducing the walk-stitching technique (Sec. 2.2.1), a general termination strategy (Sec. 2.2.2), and their connection to higher-order de-convolutions. In Sec. 2 (continued in Sec. 3), we also provide an intrinsic connection between finding a particular instantiation of the GRFs++ algorithm for a given graph kernel and higher-order (*de-*)convolutions of the discrete series encoding its kernel matrix as a Taylor series involving powers of the graph’s weight matrices.
2. In Sec. 3, we provide theoretical analysis of GRFs++, including its unbiasedness and concentration results. We show that stitching more walks improves approximation.
3. In Sec. 4, we provide thorough experimental evidence comparing GRFs++ to regular GRFs on approximation quality, speed, and several downstream tasks: normal vector field prediction on meshes, clustering and graph classification.
4. We conclude in Sec. 5 and provide all additional results in the Appendix (Sec. A).

## 2 REFINED GRFS (GRFS++)

### 2.1 PRELIMINARIES: REGULAR GRFS

We start by providing an overview of the regular GRF mechanism. We take a weighted undirected graph  $G(V, E, \mathbf{W} = [w(i, j)]_{i, j \in V})$  with  $N$  nodes/vertices, where (1)  $V$  is a set of vertices, (2)  $E \subseteq V \times V$  is a set of undirected edges ( $(i, j) \in E$  indicates that there is an edge between  $i$  and  $j$  in  $G$ ), and (3)  $\mathbf{W} \in \mathbb{R}_{\geq 0}^{N \times N}$  is a weighted adjacency matrix (if no edge exists then the corresponding weight is zero).

We consider the following kernel matrix  $\mathbf{K}_{\alpha}(\mathbf{W}) \in \mathbb{R}^{N \times N}$ , where  $\alpha = (\alpha_k)_{k=0}^{\infty}$  and  $\alpha_k \in \mathbb{R}$ :

$$\mathbf{K}_{\alpha}(\mathbf{W}) = \sum_{k=0}^{\infty} \alpha_k \mathbf{W}^k. \quad (2)$$

For bounded  $(\alpha_k)_{k=0}^{\infty}$  and  $\|\mathbf{W}\|_{\infty}$  small enough, the above sum converges. The matrix  $\mathbf{K}_{\alpha}(\mathbf{W})$  defines a kernel on the nodes of the underlying graph. Interestingly, Eq. 2 covers

all the special cases of graph node kernels we explicitly listed in Sec. 1. It also covers functions that are not positive definite, since  $(\alpha_k)_{k=0}^\infty$  can be chosen arbitrarily. From now on, we will associate graph kernels with sequences  $(\alpha_k)_{k=0}^\infty$ .

GRFs enable one to rewrite  $\mathbf{K}_\alpha(\mathbf{W})$  (in expectation) as  $\mathbf{K}_\alpha(\mathbf{W}) \stackrel{\mathbb{E}}{=} \mathbf{K}_1 \mathbf{K}_2^\top$ , for independently sampled  $\mathbf{K}_1, \mathbf{K}_2 \in \mathbb{R}^{N \times d}$  and some  $d \leq N$ . This factorization enables efficient (sub-quadratic) and unbiased approximation of the matrix-vector products  $\mathbf{K}_\alpha(\mathbf{W})\mathbf{x}$  as  $\mathbf{K}_1(\mathbf{K}_2^\top \mathbf{x})$ , if  $\mathbf{K}_1, \mathbf{K}_2$  are sparse or  $d = o(N)$ . This is often the case in practice. However, if this does not hold, explicitly materializing  $\mathbf{K}_1 \mathbf{K}_2^\top$  enables one to approximate  $\mathbf{K}_\alpha(\mathbf{W})$  in quadratic (c.f. cubic) time. Below, we describe the base GRF method for constructing sparse  $\mathbf{K}_1, \mathbf{K}_2$  for  $d = N$ . Extensions giving  $d = o(N)$ , using the Johnson-Lindenstrauss Transform (Freksen, 2021), can be found in (Choromanski, 2023). Each  $\mathbf{K}_j$  for  $j \in \{1, 2\}$  is obtained by row-wise stacking of the vectors  $\phi_f(i) \in \mathbb{R}^N$  for  $i \in V$ , where  $f$  is the *modulation function*  $f : \mathbb{R} \rightarrow \mathbb{R}$ , specific to the graph kernel being approximated. The procedure to construct random vectors  $\phi_f(i)$  is given in Algorithm 1. Intuitively, one samples an ensemble of RWs from each node  $i \in V$ . Every time a RW visits a node, the scalar value in that node (the so-called *load*) is updated, depending on the modulation function.

---

**Algorithm 1 Regular GRFs:** Construct vectors  $\phi_f(i) \in \mathbb{R}^N$  to approximate  $\mathbf{K}_\alpha(\mathbf{W})$

---

**Input:** weighted adjacency matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$ , vector of unweighted node degrees (number of out-neighbours)  $\text{deg} \in \mathbb{R}^N$ , modulation function  $f : (\mathbb{N} \cup \{0\}) \rightarrow \mathbb{R}$ , termination probability  $p_{\text{halt}} \in (0, 1)$ , node  $i \in \mathcal{N}$ , number of random walks to sample  $m \in \mathbb{N}$ .

**Output:** random feature vector  $\phi_f(i) \in \mathbb{R}^N$

```

1: initialize:  $\phi_f(i) \leftarrow \mathbf{0}$ 
2: for  $w = 1, \dots, m$  do
3:   initialise:  $\text{load} \leftarrow 1$ ,  $\text{current\_node} \leftarrow i$ ,  $\text{terminated} \leftarrow \text{False}$ ,  $\text{walk\_length} \leftarrow 0$ 
4:   while  $\text{terminated} = \text{False}$  do
5:      $\phi_f(i)[\text{current\_node}] \leftarrow \phi_f(i)[\text{current\_node}] + \text{load} \times f(\text{walk\_length})$ 
6:      $\text{walk\_length} \leftarrow \text{walk\_length} + 1$ 
7:      $\text{new\_node} \leftarrow \text{Unif}[\mathcal{N}(\text{current\_node})]$  ▷ assign to one of neighbours
8:      $\text{load} \leftarrow \text{load} \times \frac{\text{deg}[\text{current\_node}]}{1 - p_{\text{halt}}} \times \mathbf{W}[\text{current\_node}, \text{new\_node}]$  ▷ update load
9:      $\text{current\_node} \leftarrow \text{new\_node}$ 
10:     $\text{terminated} \leftarrow (t \sim \text{Unif}(0, 1) < p_{\text{halt}})$  ▷ draw RV  $t$  to decide on termination
11:   end while
12: end for
13: normalize:  $\phi_f(i) \leftarrow \phi_f(i) / m$ 

```

---

After all the walks terminate, the vector  $\phi_f(i)$  is obtained by concatenation of all the scalars/loads from the discrete scalar field, followed by a simple renormalization. It remains to describe how the kernel-dependent modulation function  $f$  is constructed. For unbiased estimation,  $f : \mathbb{N} \rightarrow \mathbb{C}$  needs to satisfy  $\sum_{p=0}^k f(k-p)f(p) = \alpha_k$ , for  $k = 0, 1, \dots$  (see Theorem 2.1 in (Reid et al., 2024b)).

## 2.2 FROM GRFs TO GRFs++

### 2.2.1 WALK-STITCHING MECHANISM

The inherently sequential procedure of constructing random walks is not supported by modern accelerators. This is one of the key weaknesses of regular GRFs. Shortening the walks by increasing  $p_{\text{halt}}$  can in principle mitigate this, at the cost of giving up modeling relationships between more distant nodes in the graph; a graph kernel value between two nodes  $i$  and  $j$  whose corresponding walks do not intersect is approximated by zero.

In GRFs++, we propose a novel *walk-stitching* technique, where several independently-calculated shorter walks are combined to emulate sampling a longer walk. Mathematically,

we unbiasedly approximate graph kernel matrix  $\mathbf{K}_\alpha(\mathbf{W})$  as:

$$\mathbf{K}_\alpha(\mathbf{W}) \stackrel{\mathbb{E}}{=} \prod_{i=1}^l \mathbf{K}_1^{(i)} (\mathbf{K}_2^{(i)})^\top, \quad (3)$$

We refer to  $l \in \mathbb{N}_+$  as the walk-stitching *degree*. Each  $i$  corresponds to one pair of intersecting walks from the regular GRF mechanism. GRFs++ with degree  $l = 1$  are equivalent to regular GRFs. A schematic is given in Fig. 1.

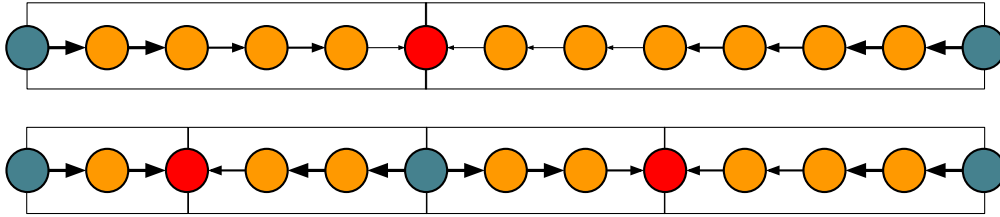


Figure 1: Pictorial description of the *walk-stitching technique*. Each rectangular block corresponds to a random walk and red nodes depict vertices where walks meet. The blue nodes are the communicating ones. The thickness of the arrow, depicting a transition from step  $t$  to step  $t+1$ , indicates the probability that such a transition will occur (a walk can terminate earlier). **Top:** In regular GRFs, two graph vertices communicate via intersecting walks, originating at each vertex. As the nodes become more distant, the probability that such two walks will be constructed decreases. **Bottom:** In GRFs++, two nodes communicate with each other less directly, via proxies (the middle blue node in the picture) and much shorter walks, with lengths that have much higher probability of being realized. The communication is established by stitching several small walks.

Each  $\mathbf{K}_j^{(i)}$ , for  $j \in \{1, 2\}$ , is computed as described in Algorithm 1, but the modulation function changes. The following is true:

**Lemma 2.1** (Unbiased walk-stitching and higher-level convolutions). *Suppose that, for each independent instantiation of Alg. 1, the modulation function  $f$  satisfies:*

$$\alpha_k = \sum_{p_1+p_2+\dots+p_{2l}=k} f(p_1)f(p_2)\dots f(p_{2l}). \quad (4)$$

Then the product  $\prod_{i=1}^l \mathbf{K}_1^{(i)} (\mathbf{K}_2^{(i)})^\top$  provides an unbiased estimation of  $\mathbf{K}_\alpha(\mathbf{W})$ .

We prove Lemma 2.1 (in fact its more general version) in Sec. 3. The condition from Lemma 2.1 is equivalent to saying that coefficients  $\alpha_k$  are obtained via  $2l$ -level discrete convolution

$\overbrace{(f \star f) \dots (f \star f)}^l$  of the modulation function  $f$  with itself. Equivalently, the function  $f$  must be constructed by  $2l$ -*de-convolving* sequence  $\alpha = (\alpha_k)_{k=0}^\infty$  that defines graph kernel.

Interestingly, for several classes of graph kernels this de-convolution can be efficiently calculated. For instance, for graph diffusion kernels with kernel-matrices of the form  $\mathbf{K}_\alpha(\mathbf{W}) = \exp(\lambda \mathbf{W})$ , the correct modulation function for GRFs++ with  $l$ -degree walk-stitching mechanism is given via a simple expression:  $f(p) = \frac{\lambda^p}{(2l)^p p!}$ . In Sec. 3, we provide a general mechanism for finding  $f$  for more arbitrary  $\mathbf{K}_\alpha$ .

### 2.2.2 GOING BEYOND THE BERNOULLI TRIAL SCHEME

Another key building block of GRFs is the *walk termination mechanism*. In regular GRFs, walk lengths are built incrementally, with walkers terminating independently with probability  $p_{\text{halt}}$  at each timestep. This gives the simple update rule in line 10 of Algorithm 1. However, sampling walk lengths from the Bernoulli distribution is not necessarily optimal given fixed computational budget (e.g. fixed average walk length). Here, we propose a very general scheme of RW-length sampling, proposing a simple modification to the update step in Algorithm 1 that improves kernel estimation accuracy.

Take any discrete probabilistic distribution on  $\mathbb{N}$ :  $\mathbf{P} = (P(i))_{i=0}^\infty$ . We will only assume that: (1) sampling  $X \sim \mathbf{P}$  and (2) the computation of  $\mathbb{P}(X \geq k)$  for any given  $k \in \mathbb{N}$  can be conducted efficiently. We modify Algorithm 1 as follows, to obtain Algorithm 2:

1. The  $m$  lengths of walks are sampled:  $s_1, \dots, s_m \stackrel{\text{iid}}{\sim} \mathbf{P}$  before line 2.

2. Line 5 is updated as follows, for  $\tau(k) \stackrel{\text{def}}{=} \mathbb{P}(X \geq k)$ :

$$\phi_f(i)[\text{current\_node}] \leftarrow \phi_f(i)[\text{current\_node}] + \frac{\text{load} \times f(\text{walk\_length})}{\tau(\text{walk\_length})}.$$

3. In line 8, term  $1 - p_{\text{halt}}$  is dropped from the update equation.

4. Line 10 is updated as follows:  $\text{terminated} \leftarrow \mathbb{I}[\text{walk\_length} \geq s_m]$ .

Note that Algorithm 1 is a special instantiation of Algorithm 2, with  $\mathbf{P}$  corresponding to the number of consecutive successes of a Bernoulli scheme with failure probability  $p_{\text{halt}}$ . In Section 3, we show that Lemma 2.1 still holds if Algorithm 1 is replaced by Algorithm 2.

### 2.2.3 PUTTING IT ALL TOGETHER

We are ready to present the complete GRFs++ mechanism, which we will refer to as  $(l, \mathbf{P})$ -GRFs++, where  $l \in \mathbb{N}_+$  and  $\mathbf{P} \in \mathcal{P}(\mathbb{N})$  are the hyperparameters of the mechanism. We first construct  $(\mathbf{K}_1^i, \mathbf{K}_2^i)_{i=1}^l$ , as in Lemma 2.1, but with Algorithm 2 replacing Algorithm 1.

**Option I:** In the most direct approach, the refined random feature vectors are given as rows of the following two matrices  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{N \times N}$ , satisfying  $\mathbf{K}_\alpha(\mathbf{W}) \stackrel{\mathbb{E}}{=} \mathbf{X}\mathbf{Y}^\top$ :

$$\mathbf{X} = \prod_{i=1}^{\frac{l}{2}} \mathbf{K}_1^{(i)} (\mathbf{K}_2^{(i)})^\top, \mathbf{Y} = \prod_{i=l}^{\frac{l}{2}+1} \mathbf{K}_2^{(i)} (\mathbf{K}_1^{(i)})^\top, \text{ if } l \text{ is even} \quad (5)$$

$$\mathbf{X} = \left[ \prod_{i=1}^{\frac{l-1}{2}} \mathbf{K}_1^{(i)} (\mathbf{K}_2^{(i)})^\top \right] \mathbf{K}_1^{(\frac{l+1}{2})}, \mathbf{Y} = \left[ \prod_{i=l}^{\frac{l+3}{2}} \mathbf{K}_2^{(i)} (\mathbf{K}_1^{(i)})^\top \right] \mathbf{K}_2^{(\frac{l+1}{2})}, \text{ if } l \text{ is odd.} \quad (6)$$

We define the product of the empty sequence of matrices as an identity matrix. If all the matrices  $\mathbf{K}_j^{(i)}$  are sparse (i.e. contain only linear in  $N$  number of nonzero entries; note for instance that for the regular  $p_{\text{halt}}$ -termination strategy, the average number of those entries is  $Nm \frac{1-p_{\text{halt}}}{p_{\text{halt}}}$ ), then  $\mathbf{X}, \mathbf{Y}$  can be computed in time  $O(N^2)$  for constant  $l$  and are also sparse. This means the refined random feature vectors are sparse, like their regular counterparts.

**Option II:** Like regular GRFs (Choromanski, 2023), the Johnson-Lindenstrauss Transform (JLT)(Freksen, 2021) can be used to reduce the dimensionality of GRFs++, at the cost of sacrificing their sparsity. The formula for matrices  $\mathbf{X}, \mathbf{Y}$  is analogous to this from Option I, but with matrices  $\mathbf{K}_j^{(i)}$  replaced by their down-projections, obtained with random Gaussian variates. In particular, we take

$$\widehat{\mathbf{K}}_j^{(i)} = \frac{1}{\sqrt{r}} \mathbf{K}_j^{(i)} \mathbf{G}^{(i)}, \quad (7)$$

for independently created Gaussian matrices  $\mathbf{G}^{(i)} \in \mathbb{R}^{N \times r}$ , with entries taken independently at random from  $\mathcal{N}(0, 1)$  and a hyperparameter  $r \in \mathbb{N}$ . For constant  $l, r$ , the computation of all  $\widehat{\mathbf{K}}_j^{(i)}$  can be done in  $O(N^2)$  time (with no sparsity assumption on  $\mathbf{K}_j^{(i)}$ ). Note also, that under this condition, matrices  $\mathbf{X}, \mathbf{Y}$  can be computed in time  $O(N)$ , via matrix associativity property. Since the JLT preserves dot-products in expectation, we conclude that  $\mathbb{E}[\widehat{\mathbf{K}}_1^{(i)} \widehat{\mathbf{K}}_2^{(i)}] = \mathbb{E}[\mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)}]$  for each  $i$ . Thus resulting  $\mathbf{X}, \mathbf{Y}$  still satisfy:  $\mathbf{K}_\alpha(\mathbf{W}) = \mathbb{E}[\mathbf{X}\mathbf{Y}^\top]$ .

**Option III:** In practice, as for regular GRFs, GRFs++ do not always need to be explicitly constructed. In most applications of random feature methods, one only needs access to products between the (approximate) kernel matrix and vectors, rather than the kernel matrix itself. As such, one only needs to support efficient multiplication algorithm for  $\left[ \prod_{i=1}^l \mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)} \right] \mathbf{v}$  for any  $\mathbf{v} \in \mathbb{R}^N$ . This can be done by multiplying with matrices from the

chain  $\prod_{i=1}^l \mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)}$  or the chain  $\prod_{i=1}^l \widehat{\mathbf{K}}_1^{(i)} \widehat{\mathbf{K}}_2^{(i)}$  from right to left, exploiting associativity. If we use the setting from Option I,  $l$  is constant and individual matrices are sparse, so this can be done in time  $O(N)$  (rather than brute-force  $O(N^2)$ ). This is also the case if Option II is applied with constant  $r$ .

**Parallel computations of random walks in GRFs++:** One of the most attractive computational features of GRFs++ is that one can compute short RWs in parallel for different  $i = 1, 2, \dots, l$ . These are in turn put together by walk-stitching, implicitly constructing longer walks. This gives computational gains compared to regular GRFs, since it avoids explicit, sequential sampling of longer walks.

**Re-using the same set of random walks:** Even though for unbiasedness, different matrices  $\mathbf{K}_j^{(i)}$  for  $j \in \{1, 2\}$  ought to use independent sets of random walks, we empirically observe that in practice re-using the same set of random walks also works very well. This is especially the case for larger graphs of higher diameter; see Section 4.

**Walk-stitching with general termination strategies:** To see how walk-stitching helps more distant nodes to connect with each other, consider a termination strategy, where the first transition occurs with probability  $p_0 = 1$  and consequent transitions occur with probability  $p_{\text{next}} < 1$ . In such a setting, the probability of regular GRFs emulating any existing walk of length  $r \geq 2$ , joining two given vertices  $i$  and  $j$  scales with  $p_{\text{next}}$  as  $p_{\text{next}}^{r-2}$ , whereas for walk-stitching of degree  $\lceil \frac{r}{2} \rceil$ , this walk will be emulated via GRFs++ with probability lower-bounded by the expression completely independent of  $p_{\text{next}}$ .

### 3 THEORETICAL ANALYSIS

We are ready to provide a rigorous theoretical analysis of GRFs++. We start by presenting a strengthened version of Lemma 2.1 from Section 2 (proof in App. A.1).

**Lemma 3.1** (Unbiased walk-stitching, higher-order convolutions & general termination). *Lemma 2.1 remains true if Algorithm 1 in its statement is replaced by Algorithm 2.*

Since Algorithm 2 is more general than Algorithm 1 (see Sec. 2.2.2), this also proves Lemma 2.1. The setting with Algorithm 1 and  $l = 1$  is equivalent to regular GRFs.

The formula for the mean squared error (MSE) of the GRFs++-based graph kernel estimator with general degree  $l \geq 1$  in terms of the individual components  $\mathbf{X}_i = \mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)}$  is complicated. However, for degree  $l = 2$ , it has a particularly compact form, provided below.

**Lemma 3.2** (MSE of the approximation via GRFs++ with  $l = 2$ ). *The MSE of the estimator  $\widehat{\mathbf{K}}_\alpha(\mathbf{W})$  of the groundtruth graph kernel matrix  $\mathbf{K}_\alpha(\mathbf{W})$ , leveraging GRFs++ with degree  $l = 2$  satisfies (proof in the Appendix: Sec. A.2,  $\|\cdot\|_F$  stands for the Frobenius norm):*

$$\text{MSE}(\widehat{\mathbf{K}}_\alpha(\mathbf{W})) \stackrel{\text{def}}{=} \mathbb{E}[\|\mathbf{X}_1 \mathbf{X}_2 - \mathbf{K}_\alpha(\mathbf{W})\|_F^2] = \mathbb{E}[\|\mathbf{X}_1^\top \mathbf{X}_1\|_F^2] - \|\mathbf{K}_\alpha(\mathbf{W})\|_F^2 \quad (8)$$

Finally, we show that the approximation of the graph kernel monotonically improves with the GRFs++ degree (for degrees being the powers of two; proof in App. A.3).

**Theorem 3.3.** *If  $\widehat{\mathbf{K}}_\alpha^{(l)}(\mathbf{W})$  stands for the estimator of the groundtruth graph kernel matrix, leveraging GRFs++ of degree  $l$ , then the following holds if standard termination strategy is applied:*

$$\text{MSE}(\widehat{\mathbf{K}}_\alpha^{(1)}(\mathbf{W})) \geq \text{MSE}(\widehat{\mathbf{K}}_\alpha^{(2)}(\mathbf{W})) \geq \text{MSE}(\widehat{\mathbf{K}}_\alpha^{(4)}(\mathbf{W})) \geq \dots \quad (9)$$

#### 3.1 DE-MYSTIFYING $2l$ -LEVEL DE-CONVOLUTIONS

Let us assume that the coefficient  $\alpha = (\alpha_k)_{k=0}^\infty$  defining graph kernel, encode also an analytical function  $g : \mathbb{C} \rightarrow \mathbb{C}$  of the form:  $g(x) = \sum_{i=0}^\infty \alpha_k x^k$ . Assume furthermore that one can compute  $h(x) = g^{\frac{1}{2l}}(x)$  and its Taylor expansion is of the form:  $h(x) = \sum_{i=0}^\infty \beta_i x^i$ . Then it is easy to see that function:  $f(p) = \beta_p$  satisfies Equation 4.

The above observation provides a straightforward algorithm for computing modulation function  $f$  for GRFs++ with hyperparameter  $l$ : (1) map the graph kernel under consideration to function  $g$ , (2) compute its  $(2l)^{\text{th}}$ -root  $h$ , (3) find Taylor series of  $h$  to define  $f$ .

**Remark 3.4.** Now we also see why the formula for  $f$  in the GRFs++ mechanism corresponding to the diffusion graph kernel is particularly simple for any  $l \in \mathbb{N}_+$ . The roots of  $g : x \rightarrow \exp(x)$ , that corresponds to that kernel, are trivial to compute.

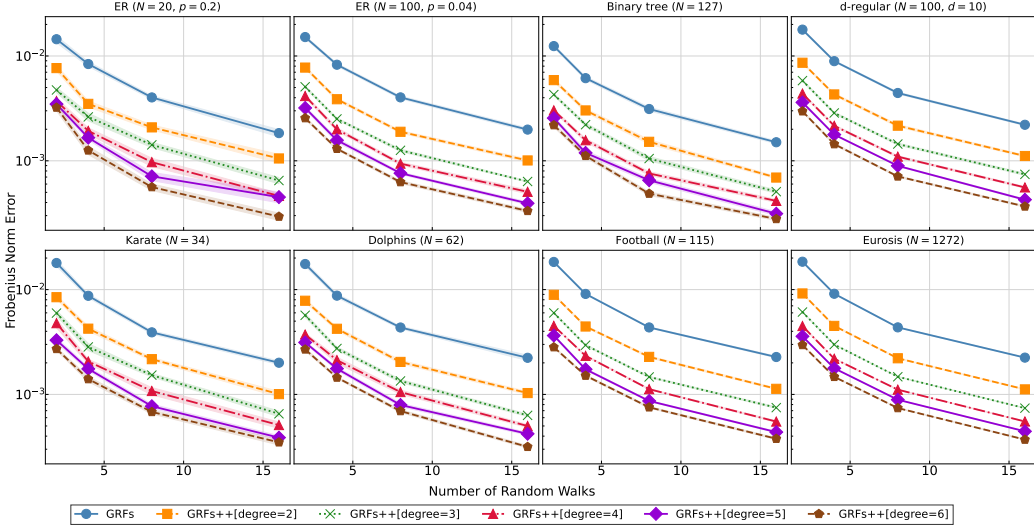


Figure 2: Comparison of different GRF methods for the diffusion kernel estimation. The approximation error (y-axis) improves with the number of walks  $m$  (x-axis) and GRF++ provides a sharper estimate than the previous GRF mechanism. The experiment is repeated  $s = 10$  times.

## 4 EXPERIMENTS

In this section, we showcase the ability of GRFs++ to efficiently approximate graph node kernels (see Sec. 4.1), including with larger diameters. Furthermore, we show downstream applications of GRFs++ in graph classification, node clustering tasks and normal field prediction on meshes (see Sec 4.2). We use the graph diffusion kernel.

### 4.1 ACCURATE ESTIMATION OF GRAPH KERNELS WITH GRFs++

Following (Reid et al., 2024b), we choose eight graphs of varying sizes: (1) Erdős-Rényi graphs of two sizes, (2) a binary tree, (3) a d-regular graph, and (4) four real world examples (karate, dolphins, football and eurosis). Fig. 2 plots the relative Frobenius norm error of the approximation  $\hat{\mathbf{K}}_\alpha(\mathbf{W})$  of the groundtruth kernel matrix  $\mathbf{K}_\alpha(\mathbf{W})$  with GRFs++ (i.e.,  $\|\mathbf{K}_\alpha(\mathbf{W}) - \hat{\mathbf{K}}_\alpha(\mathbf{W})\|_F / \|\mathbf{K}_\alpha(\mathbf{W})\|_F$ ) against the number of random walks  $m$ , showcasing improved estimation accuracy with GRFs++. Next, we show that our method can capture long-distance information more accurately than regular GRFs (see Fig. 8 in Appendix). For this task, we select eight diverse graphs from datasets including Peptides (Dwivedi et al., 2022), CIFAR-10 (Dwivedi et al., 2020), Reddit-Binary (Morris et al., 2020), and Geometric Shapes (Yannick-S, 2025). These graphs exhibit varying degrees of sparsity and heterophily, yet all are characterized by large diameters (up to  $\text{diam} = 159$ ). We again compute kernel estimation error (with  $p_{\text{halt}} = 0.1$ ), but now for node pairs that are a specified distance apart. Again, GRFs++ are more accurate. See App. B.1.1 for details.

**New Termination Strategy:** Next, we investigate the benefits a more general (non-Bernoulli) termination strategy, as described in Sec. 2.2.2. Specifically, we employ a halting probability governed by a **Poisson distribution**  $\mathbf{P}$ . For a fair comparison, the parameters are chosen so that the expected random walk length remains the same as in regular GRFs. Fig. 3 shows that our novel halting strategy improves regular GRF mechanism on a wide range of diverse graphs. We also get a more accurate estimation of kernel values for distant nodes in large diameter graphs (see Fig. 9 in Appendix).

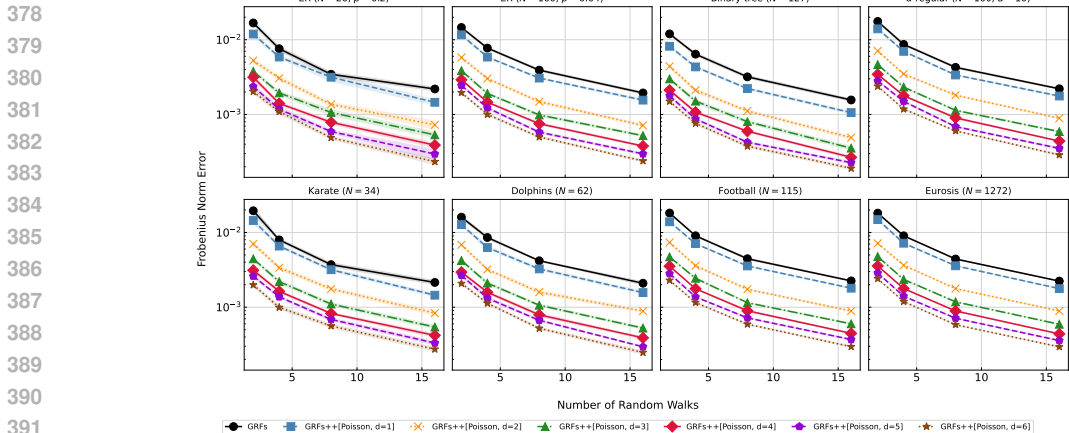


Figure 3: Our novel halting policy based on Poisson distribution provides additional gains over the GRF mechanisms. We run the experiment  $s = 10$  times on different graphs of varying sizes.

**Re-using the same set of random walks:**

Finally, we conduct an ablation study (see Fig. 4) to pinpoint the benefit of the walk-stitching mechanism itself. In this experiment, rather than sampling new independent walks for each component (as before), we take the exact same set of random walks generated for the baseline GRF method and re-use them for each degree of the GRFs++ estimator. This still results in a significant improvement in the Frobenius norm error for all GRFs++ variants, compared to the regular GRFs baseline (see Tab. 8 for a practical application of a downstream experiment).

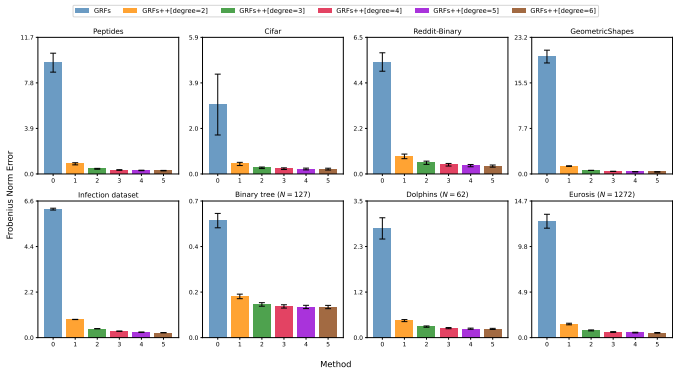


Figure 4: Using the exact same walk as the baseline GRF, repeated multiple times, pinpoints the effectiveness of the walk-stitching algorithm, showing additional computational gains.

**Computational Time:**

We generated random graphs with 500 nodes. We evaluated two configurations using base halting probabilities ( $p_{halt}$ ) of 0.01 and 0.001, which correspond to the regular GRF method (degree  $l = 1$ ). For the GRFs++ methods of degree  $l > 1$ , we used a scaled halting probability of  $p_{halt} \times l$  to ensure a fair comparison. Fig. 5 presents a computational speed analysis of the GRFs++ algorithm, breaking down its performance into two key components: “Walk Time” (the time required for random walk sampling) and “Stitching Time” (the time for the matrix-matrix operations used in the walk-stitching technique). As the degree increases, both the walk time and the stitching time decreases.

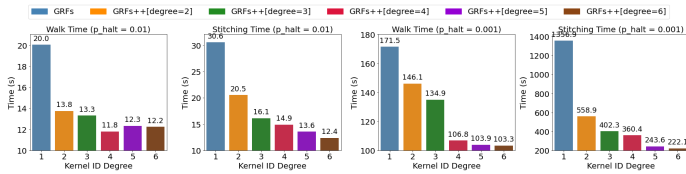


Figure 5: Speed comparison for various GRF-methods: regular GRFs and GRFs++ with different degrees.

4.2 DOWNSTREAM TASKS

In this section, we show the efficacy of our GRFs++ based approximate kernel in various downstream tasks: graph classification, node clustering and vertex normal prediction.

**Graph Classification** : Graph kernels have been widely used for graph classification tasks (Kriege et al., 2020; Nikolentzos et al., 2021). We compare the graph classification results obtained using the approximate kernel from GRF++ with those from the exact diffusion kernel on a wide variety of datasets (Morris et al., 2020). Fig. 6 shows that our method performs on par with the diffusion kernel and outperforms regular GRF (additional details in App. B.2).

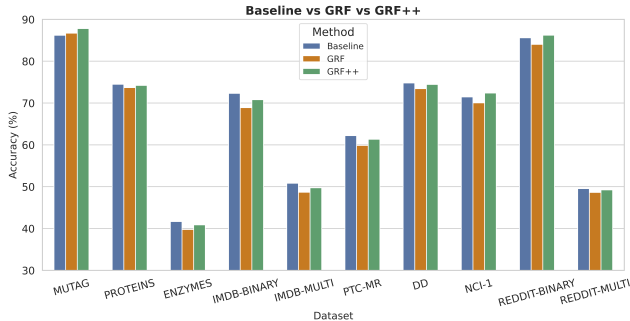


Figure 6: Graph classification using the approximate diffusion kernel from GRF++. Our method performs at par with the baseline diffusion kernel and always beats GRFs.

**Node Clustering**: We also test GRFs++’s utility on the downstream task of node clustering. For this experiment, we perform spectral clustering, implemented using the SciPy library, to group nodes based on the approximated diffusion kernel. We compare the clustering error  $E = (\# \text{ of wrong pairs}) / (N * (N - 1))$  of the baseline GRF with our GRF++[degree=2] method. The results, presented in Tab. 9, show that GRF++ achieves a lower error rate **across all tested datasets** (additional details in App. B.3).

Table 1: Node Clustering: GRFs++ vs regular GRFs.

Name	# Nodes	# clusters	GRF	GRF++[d=2]
Karate	34	2	0.2995	<b>0.2585</b>
Dolphins	62	2	0.0635	<b>0.0323</b>
Polbooks	105	3	0.1060	<b>0.1033</b>
Football	115	12	0.0731	<b>0.0362</b>
Databases	1006	6	0.3528	<b>0.3001</b>
Eurosis	1272	13	0.2248	<b>0.1304</b>

**Normal Prediction** : We test GRF++ on normal vector prediction (mesh interpolation) Every vertex of the graph mesh  $G$  with a vertex-set  $V$ , is associated with spatial coordinates  $x_i \in \mathbb{R}^3$  and a unit normal vector  $F_i \in \mathbb{R}^3$ . Following (Choromanski et al., 2024), we randomly sample a subset  $V' \subset V$  from each mesh with  $|V'| = 0.8|V|$  and mask out their vertex normals. Our goal is to predict the vertex normals of each masked vertex  $i \in V'$  via:  $F_i = \sum_{j \in V \setminus V'} K(i, j)F_j$ , where  $K$  is the diffusion kernel. We report the cosine similarity between predicted and groundtruth vertex normals, averaged over all the nodes. We validate GRF++ over 40 meshes of 3D printed objects of varying sizes from the Thingi10K dataset (Zhou & Jacobson, 2016). Additional details are provided in App. B.4. To save space, we show all **40+** mesh results in Tab. 8 in the Appendix. Here we present a few larger meshes in Table 2. GRFs++ provide consistent gains, compared to regular GRFs.

Table 2: Cosine Similarity results for Meshes. GRF++ matches the baseline kernel (BF) and outperforms GRFs. GRF++r, reusing the same random walk, also outperforms GRF.

MESH SIZE	5985	6577	6911	7386	7953	8011	8261	8449	8800	9603
BF	0.9194	0.9622	0.9769	0.9437	0.9460	0.9382	0.9196	0.9276	0.9836	0.9766
GRF	0.9091	0.9525	0.9682	0.9308	0.9383	0.9233	0.9050	0.9139	0.9778	0.9708
GRF++	0.9154	0.9599	0.9751	0.9374	0.9429	0.9321	0.9145	0.9205	0.9820	0.9748
GRF++r	0.9129	0.9561	0.9701	0.9348	0.9410	0.9269	0.9095	0.9160	0.9805	0.9734
Diff	<b>0.0063</b>	<b>0.0074</b>	<b>0.0069</b>	<b>0.0066</b>	<b>0.0046</b>	<b>0.0088</b>	<b>0.0095</b>	<b>0.0066</b>	<b>0.0042</b>	<b>0.0040</b>

## 5 CONCLUSION

We introduced *refined GRFs* (GRFs++) for improved approximation of graph node kernels. GRFs++ address regular GRFs’ shortcomings, modeling the relationship between distant pairs of nodes more effectively and introducing novel random walk termination strategies. GRFs++ provide more accurate and more efficient kernel approximation, replacing computationally-inefficient and inherently sequential sampling of long random walks with matrix-matrix operations. We complement our algorithm with theoretical analysis, showing that GRFs++ give unbiased approximation. We provide concentration results, as well as detailed empirical evaluation on a wide variety of graph datasets and tasks.

## 6 REPRODUCIBILITY STATEMENT

The paper provides a clear description of the GRFs++ algorithm. In Sec. 2.1, we present detailed description of the regular GRFs algorithm, namely Algorithm 1 box, that GRFs++ build on. This algorithm is also implemented in the github repository mentioned on the first page of (Reid et al., 2024b). In Lemma 2.1, we explain how Algorithm 1 is used in GRFs++ for the general walk-stitching mechanism. Then in Sec. 2.2.2, we provide detailed explanation of the modification of Algorithm 1 that needs to be conducted in order to support arbitrary termination strategies (points 1-3). For all the experiments, we provided the names of all datasets and graphs used (or exact procedures to construct those graphs, e.g. random Erdős-Rényi graph models with explicitly given probabilities  $p$  of edge sampling). All the theoretical statements have all the assumptions clearly stated and the corresponding proofs given (see: Section 2.2.1, Section 3 and Appendix: Section A.1, Section A.2 and Section A.3). Finally, we provide a pointer to the anonymous github repository with code in Appendix: Section B.7.

## REFERENCES

- Anirban Banerjee and Jürgen Jost. Graph spectra as a systematic tool in computational biology. *Discret. Appl. Math.*, 157(10):2425–2431, 2009. doi: 10.1016/J.DAM.2008.06.033. URL <https://doi.org/10.1016/j.dam.2008.06.033>.
- Alex Beutel, Leman Akoglu, and Christos Faloutsos. Graph-based user behavior modeling: From prediction to fraud detection. In Longbing Cao, Chengqi Zhang, Thorsten Joachims, Geoffrey I. Webb, Dragos D. Margineantu, and Graham Williams (eds.), *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-13, 2015*, pp. 2309–2310. ACM, 2015. doi: 10.1145/2783258.2789985. URL <https://doi.org/10.1145/2783258.2789985>.
- Jingyan Chen, Guanghui Zhu, Chunfeng Yuan, and Yihua Huang. Boosting graph anomaly detection with adaptive message passing. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=CanomFZssu>.
- Krzysztof Marcin Choromanski. Taming graph kernels with random features. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 5964–5977. PMLR, 2023. URL <https://proceedings.mlr.press/v202/choromanski23a.html>.
- Krzysztof Marcin Choromanski, Arijit Sehanobish, Somnath Basu Roy Chowdhury, Han Lin, Kumar Avinava Dubey, Tamas Sarlos, and Snigdha Chaturvedi. Fast tree-field integrators: From low displacement rank to topological transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=Eok6HbcSRI>.
- Krzysztof Marcin Choromanski, Isaac Reid, Arijit Sehanobish, and Kumar Avinava Dubey. Optimal time complexity algorithms for computing general random walk graph kernels on sparse graphs. In Yingzhen Li, Stephan Mandt, Shipra Agrawal, and Mohammad Emtiyaz Khan (eds.), *International Conference on Artificial Intelligence and Statistics, AISTATS 2025, Mai Khao, Thailand, 3-5 May 2025*, volume 258 of *Proceedings of Machine Learning Research*, pp. 3457–3465. PMLR, 2025. URL <https://proceedings.mlr.press/v258/choromanski25a.html>.
- Nathan de Lara and Edouard Pineau. A simple baseline algorithm for graph classification, 2018.
- Leyan Deng, Defu Lian, Chenwang Wu, and Enhong Chen. Graph convolution network based recommender systems: Learning guarantee and item mixture powered strategy. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.),

- 540 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural*  
 541 *Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November*  
 542 *28 - December 9, 2022, 2022.* URL [http://papers.nips.cc/paper\\_files/paper/2022/](http://papers.nips.cc/paper_files/paper/2022/hash/18fd48d9cbbf9a20e434c9d3db6973c5-Abstract-Conference.html)  
 543 [hash/18fd48d9cbbf9a20e434c9d3db6973c5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/18fd48d9cbbf9a20e434c9d3db6973c5-Abstract-Conference.html).
- 544 Xiangyu Dong, Xingyi Zhang, Lei Chen, Mingxuan Yuan, and Sib0 Wang. Spacegnn:  
 545 Multi-space graph neural network for node anomaly detection with extremely limited  
 546 labels. In *The Thirteenth International Conference on Learning Representations, ICLR*  
 547 *2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL [https://openreview.](https://openreview.net/forum?id=Syt4fWwVm1)  
 548 [net/forum?id=Syt4fWwVm1](https://openreview.net/forum?id=Syt4fWwVm1).
- 549 Vijay Prakash Dwivedi, Chaitanya K Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua  
 550 Bengio, and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint*  
 551 *arXiv:2003.00982*, 2020.
- 552 Vijay Prakash Dwivedi, Ladislav Rampásek, Mikhail Galkin, Ali Parviz, Guy Wolf,  
 553 Anh Tuan Luu, and Dominique Beaini. Long range graph benchmark. In *Thirty-sixth*  
 554 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track,*  
 555 *2022.* URL <https://openreview.net/forum?id=in7XC5RcjEn>.
- 556 Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of  
 557 graph neural networks for graph classification. In *Proceedings of the 8th International*  
 558 *Conference on Learning Representations (ICLR)*, 2020.
- 559 Casper Benjamin Freksen. An introduction to johnson-lindenstrauss transforms. *CoRR*,  
 560 [abs/2103.00564](https://arxiv.org/abs/2103.00564), 2021. URL <https://arxiv.org/abs/2103.00564>.
- 561 Chen Gao, Yu Zheng, Nian Li, Yinfeng Li, Yingrong Qin, Jinghua Piao, Yuhan Quan,  
 562 Jianxin Chang, Depeng Jin, Xiangnan He, and Yong Li. A survey of graph neural networks  
 563 for recommender systems: Challenges, methods, and directions. *Trans. Recomm. Syst.*, 1  
 564 (1):1–51, 2023. doi: 10.1145/3568022. URL <https://doi.org/10.1145/3568022>.
- 565 William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on  
 566 large graphs. In *Proceedings of the 31st International Conference on Neural Informa-*  
 567 *tion Processing Systems, NIPS’17*, pp. 1025–1035, Red Hook, NY, USA, 2017. Curran  
 568 Associates Inc. ISBN 9781510860964.
- 569 Vladimir Ivashkin and Pavel Chebotarev. Do logarithmic proximity measures outperform  
 570 plain ones in graph clustering? In *International Conference on Network Analysis*, pp.  
 571 87–105. Springer, 2016.
- 572 Hwan Kim, Byung Suk Lee, Won-Yong Shin, and Sungsu Lim. Graph anomaly detection  
 573 with graph neural networks: Current status and challenges. *IEEE Access*, 10:111820–  
 574 111829, 2022. doi: 10.1109/ACCESS.2022.3211306. URL [https://doi.org/10.1109/](https://doi.org/10.1109/ACCESS.2022.3211306)  
 575 [ACCESS.2022.3211306](https://doi.org/10.1109/ACCESS.2022.3211306).
- 576 Hwan Kim, Junghoon Kim, Byung Suk Lee, and Sungsu Lim. Label-based graph aug-  
 577 mentation with metapath for graph anomaly detection. *Expert Syst. Appl.*, 296:129087,  
 578 2026. doi: 10.1016/J.ESWA.2025.129087. URL [https://doi.org/10.1016/j.eswa.](https://doi.org/10.1016/j.eswa.2025.129087)  
 579 [2025.129087](https://doi.org/10.1016/j.eswa.2025.129087).
- 580 Risi Kondor and John D. Lafferty. Diffusion kernels on graphs and other discrete input  
 581 spaces. In Claude Sammut and Achim G. Hoffmann (eds.), *Machine Learning, Proceedings*  
 582 *of the Nineteenth International Conference (ICML 2002), University of New South Wales,*  
 583 *Sydney, Australia, July 8-12, 2002*, pp. 315–322. Morgan Kaufmann, 2002.
- 584 Nils M. Kriege, Fredrik D. Johansson, and Christopher Morris. A survey on graph ker-  
 585 nels. *Applied Network Science*, 5(1), January 2020. ISSN 2364-8228. doi: 10.1007/  
 586 [s41109-019-0195-3](http://dx.doi.org/10.1007/s41109-019-0195-3). URL <http://dx.doi.org/10.1007/s41109-019-0195-3>.
- 587 Jingham Li, Yuan Gao, Jinda Lu, Junfeng Fang, Congcong Wen, Hui Lin, and Xiang Wang.  
 588 Diffgad: A diffusion-based unsupervised graph anomaly detector. In *The Thirteenth In-*  
 589 *ternational Conference on Learning Representations, ICLR 2025, Singapore, April 24-28,*  
 590 *2025.* OpenReview.net, 2025. URL <https://openreview.net/forum?id=AhcYq4CnfF>.

- 594 Yixin Liu, Kaize Ding, Qinghua Lu, Fuyi Li, Leo Yu Zhang, and Shirui Pan. To-  
595 wards self-interpretable graph-level anomaly detection. In Alice Oh, Tristan Nau-  
596 mann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-  
597 vances in Neural Information Processing Systems 36: Annual Conference on Neural In-  
598 formation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December  
599 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
600 1c6f06863df46de009a7a41b41c95cad-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/1c6f06863df46de009a7a41b41c95cad-Abstract-Conference.html).
- 601 Zetian Mao, Jiawen Li, Chen Liang, Diptesh Das, Masato Sumita, and Koji Tsuda. Molecule  
602 graph networks with many-body equivariant interactions. *CoRR*, abs/2406.13265,  
603 2024. doi: 10.48550/ARXIV.2406.13265. URL [https://doi.org/10.48550/arXiv.  
604 2406.13265](https://doi.org/10.48550/arXiv.2406.13265).
- 605  
606 Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and  
607 Marion Neumann. Tudataset: A collection of benchmark datasets for learning with  
608 graphs. In *ICML 2020 Workshop on Graph Representation Learning and Beyond (GRL+  
609 2020)*, 2020. URL [www.graphlearning.io](http://www.graphlearning.io).
- 610  
611 Giannis Nikolentzos, Giannis Siglidis, and Michalis Vazirgiannis. Graph kernels: A survey.  
612 *Journal of Artificial Intelligence Research*, 72:943–1027, November 2021. ISSN 1076-9757.  
613 doi: 10.1613/jair.1.13225. URL <http://dx.doi.org/10.1613/jair.1.13225>.
- 614  
615 Caleb C. Noble and Diane J. Cook. Graph-based anomaly detection. In Lise Getoor,  
616 Ted E. Senator, Pedro M. Domingos, and Christos Faloutsos (eds.), *Proceedings of the  
617 Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Min-  
618 ing, Washington, DC, USA, August 24 - 27, 2003*, pp. 631–636. ACM, 2003. doi:  
10.1145/956750.956831. URL <https://doi.org/10.1145/956750.956831>.
- 619  
620 Isaac Reid, Adrian Weller, and Krzysztof Marcin Choromanski. Quasi-monte  
621 carlo graph random features. In Alice Oh, Tristan Naumann, Amir Globerson,  
622 Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neu-  
623 ral Information Processing Systems 36: Annual Conference on Neural Information  
624 Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10  
625 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/paper/2023/hash/  
626 2f9b3ee2bcea04b327c09d7e3145bd1e-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/2f9b3ee2bcea04b327c09d7e3145bd1e-Abstract-Conference.html).
- 627  
628 Isaac Reid, Eli Berger, Krzysztof Marcin Choromanski, and Adrian Weller. Repelling  
629 random walks. In *The Twelfth International Conference on Learning Representations,  
630 ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. URL <https://openreview.net/forum?id=31I0mrnP4>.
- 631  
632 Isaac Reid, Krzysztof Marcin Choromanski, Eli Berger, and Adrian Weller. General graph  
633 random features. In *The Twelfth International Conference on Learning Representations,  
634 ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024b. URL <https://openreview.net/forum?id=viftsX50Rt>.
- 635  
636 Isaac Reid, Kumar Avinava Dubey, Deepali Jain, William F. Whitney, Amr Ahmed, Joshua  
637 Ainslie, Alex Bewley, Mithun George Jacob, Aranyak Mehta, David Rendleman, Connor  
638 Schenck, Richard E. Turner, René Wagner, Adrian Weller, and Krzysztof Marcin Choro-  
639 manski. Linear transformer topological masking with graph random features. In *The  
640 Thirteenth International Conference on Learning Representations, ICLR 2025, Singa-  
641 pore, April 24-28, 2025*. OpenReview.net, 2025. URL [https://openreview.net/forum?  
642 id=6MBqQLp17E](https://openreview.net/forum?id=6MBqQLp17E).
- 643  
644 Nino Shervashidze, S. V. N. Vishwanathan, Tobias Petri, Kurt Mehlhorn, and Karsten M.  
645 Borgwardt. Efficient graphlet kernels for large graph comparison. In David A. Van  
646 Dyk and Max Welling (eds.), *Proceedings of the Twelfth International Conference on  
647 Artificial Intelligence and Statistics, AISTATS 2009, Clearwater Beach, Florida, USA,  
April 16-18, 2009*, volume 5 of *JMLR Proceedings*, pp. 488–495. JMLR.org, 2009. URL  
<http://proceedings.mlr.press/v5/shervashidze09a.html>.

- 648 Alexander J. Smola and Risi Kondor. Kernels and regularization on graphs. In Bernhard  
649 Schölkopf and Manfred K. Warmuth (eds.), *Computational Learning Theory and Kernel  
650 Machines, 16th Annual Conference on Computational Learning Theory and 7th Kernel  
651 Workshop, COLT/Kernel 2003, Washington, DC, USA, August 24-27, 2003, Proceedings*,  
652 volume 2777 of *Lecture Notes in Computer Science*, pp. 144–158. Springer, 2003. doi: 10.  
653 1007/978-3-540-45167-9\\_12. URL [https://doi.org/10.1007/978-3-540-45167-9\\_12](https://doi.org/10.1007/978-3-540-45167-9_12).
- 655 S. V. N. Vishwanathan, Nicol N. Schraudolph, Risi Kondor, and Karsten M. Borgwardt.  
656 Graph kernels. *J. Mach. Learn. Res.*, 11:1201–1242, 2010. doi: 10.5555/1756006.1859891.  
657 URL <https://dl.acm.org/doi/10.5555/1756006.1859891>.
- 658  
659 Lingfei Wu, Ian En-Hsu Yen, Zhen Zhang, Kun Xu, Liang Zhao, Xi Peng, Yinglong Xia,  
660 and Charu C. Aggarwal. Scalable global alignment graph kernel using random features:  
661 From node embedding to graph embedding. In Ankur Teredesai, Vipin Kumar, Ying  
662 Li, Rómer Rosales, Evimaria Terzi, and George Karypis (eds.), *Proceedings of the 25th  
663 ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD  
664 2019, Anchorage, AK, USA, August 4-8, 2019*, pp. 1418–1428. ACM, 2019. doi: 10.1145/  
665 3292500.3330918. URL <https://doi.org/10.1145/3292500.3330918>.
- 666 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph  
667 neural networks? In *International Conference on Learning Representations*, 2019. URL  
668 <https://openreview.net/forum?id=ryGs6iA5Km>.
- 669 Pinar Yanardag and S. V. N. Vishwanathan. Deep graph kernels. In Longbing Cao,  
670 Chengqi Zhang, Thorsten Joachims, Geoffrey I. Webb, Dragos D. Margineantu, and  
671 Graham Williams (eds.), *Proceedings of the 21th ACM SIGKDD International Confer-  
672 ence on Knowledge Discovery and Data Mining, Sydney, NSW, Australia, August 10-  
673 13, 2015*, pp. 1365–1374. ACM, 2015. doi: 10.1145/2783258.2783417. URL <https://doi.org/10.1145/2783258.2783417>.
- 674  
675 Haoran Yang, Xiangyu Zhao, Yicong Li, Hongxu Chen, and Guandong Xu. An empirical  
676 study towards prompt-tuning for graph contrastive pre-training in recommendations. In  
677 Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey  
678 Levine (eds.), *Advances in Neural Information Processing Systems 36: Annual Confer-  
679 ence on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans,  
680 LA, USA, December 10 - 16, 2023*, 2023. URL [http://papers.nips.cc/paper\\_files/  
681 paper/2023/hash/c6af791af7ef0f3e02bccef011211ca5-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/c6af791af7ef0f3e02bccef011211ca5-Abstract-Conference.html).
- 682  
683 Xin Yang, Xingrun Li, Heng Chang, Jinze Yang, Xihong Yang, Shengyu Tao, Ningkan  
684 Chang, Maiko Shigeno, Junfeng Wang, Dawei Yin, and Erxue Min. Hgformer: Hyperbolic  
685 graph transformer for recommendation. *ICML 2025*, abs/2502.15693, 2025. doi: 10.48550/  
686 ARXIV.2502.15693. URL <https://doi.org/10.48550/arXiv.2502.15693>.
- 687 Yannick-S. `geometric_shapes`: Representation of geometric shapes. [https://github.com/  
688 Yannick-S/geometric\\_shapes](https://github.com/Yannick-S/geometric_shapes), 2025. Commit snapshot as of November 21, 2025.
- 689  
690 Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure  
691 Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Pro-  
692 ceedings of the 32nd International Conference on Neural Information Processing Systems*,  
693 NIPS’18, pp. 4805–4815, Red Hook, NY, USA, 2018. Curran Associates Inc.
- 694  
695 Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep  
696 learning architecture for graph classification. *Proceedings of the AAAI Conference on  
697 Artificial Intelligence*, 32(1), Apr. 2018. doi: 10.1609/aaai.v32i1.11782. URL <https://ojs.aaai.org/index.php/AAAI/article/view/11782>.
- 698  
699 Yang Zhang, Zhewei Wei, Ye Yuan, Chongxuan Li, and Wenbing Huang. Equipocket: an  
700 e(3)-equivariant geometric graph neural network for ligand binding site prediction. In  
701 *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?id=  
1vGN3CSxVs](https://openreview.net/forum?id=1vGN3CSxVs).

702 Qingnan Zhou and Alec Jacobson. Thingi10k: A dataset of 10,000 3d-printing models.  
703 *arXiv preprint arXiv:1605.04797*, 2016.  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A APPENDIX

### A.1 PROOF OF LEMMA 3.1

*Proof.* By using similar analysis, as in the proof of Theorem 2.1 from (Reid et al., 2024b), we obtain:

$$\mathbb{E}[\mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)}(a, b)] = \sum_v \sum_{p=0}^{\infty} \sum_{t=0}^{\infty} \mathbf{W}^p(a, v) \mathbf{W}^t(v, b) f(p) f(t) \mathbb{P}[X \geq p] \mathbb{P}[Y \geq t] \frac{1}{\tau(p)} \frac{1}{\tau(t)}, \quad (10)$$

where  $X, Y \stackrel{\text{iid}}{\sim} \mathbf{P}$ . Thus, from the fact that pairs  $(\mathbf{K}_1^{(i)}, \mathbf{K}_2^{(i)})$  are constructed independently for different  $i = 1, \dots, l$ , we obtain:

$$\mathbb{E} \left[ \prod_{i=1}^l \mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)}(a, b) \right] = \sum_{v_1, v_2, \dots, v_{2l-1}} \sum_{p^{(1)}}^{\infty} \dots \sum_{p^{(l)}}^{\infty} \sum_{t^{(1)}}^{\infty} \dots \sum_{t^{(l)}}^{\infty} \mathbf{W}^{p^{(1)}}(a, v_1) \mathbf{W}^{t^{(1)}}(v_1, v_2) \dots \mathbf{W}^{p^{(l)}}(v_{2l-2}, v_{2l-1}) \mathbf{W}^{t^{(l)}}(v_{2l-1}, b) f(p^{(1)}) f(t^{(1)}) \dots f(p^{(l)}) f(t^{(l)}) \quad (11)$$

Therefore, we obtain:

$$\mathbb{E} \left[ \prod_{i=1}^l \mathbf{K}_1^{(i)} \mathbf{K}_2^{(i)} \right] = \sum_{i=0}^{\infty} \left[ \sum_{p^{(1)}+t^{(1)}+\dots+p^{(l)}+t^{(l)}=i} f(p^{(1)}) f(t^{(1)}) \dots f(p^{(l)}) f(t^{(l)}) \right] \mathbf{W}^i \quad (12)$$

That completes the proof, because of Equation 4.  $\square$

### A.2 PROOF OF LEMMA 3.2

We now provide the proof of Lemma 3.2, that we re-state here for reader's convenience:

**Lemma A.1** (MSE of the GRFs++ based graph estimator with GRFs++ degree  $l = 2$ ). *The MSE of the estimator  $\widehat{\mathbf{K}}_{\alpha}(\mathbf{W})$  of the groundtruth graph kernel matrix  $\mathbf{K}_{\alpha}(\mathbf{W})$ , leveraging GRFs++ with degree  $l = 2$  satisfies (proof in the Appendix):*

$$\text{MSE}(\widehat{\mathbf{K}}_{\alpha}(\mathbf{W})) \stackrel{\text{def}}{=} \mathbb{E}[\|\mathbf{X}_1 \mathbf{X}_2 - \mathbf{K}_{\alpha}(\mathbf{W})\|_{\mathbb{F}}^2] = \|\mathbb{E}[\mathbf{X}_1^{\top} \mathbf{X}_1]\|_{\mathbb{F}}^2 - \|\mathbf{K}_{\alpha}(\mathbf{W})\|_{\mathbb{F}}^2 \quad (13)$$

*Proof.* We have the following:

$$\begin{aligned} \mathbb{E}[\|\mathbf{X}_1 \mathbf{X}_2 - \mathbf{K}_{\alpha}(\mathbf{W})\|_{\mathbb{F}}^2] &= \mathbb{E}[\|\mathbf{X}_1 \mathbf{X}_2 - \mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2] = \mathbb{E}[\|\mathbf{X}_1 \mathbf{X}_2\|_{\mathbb{F}}^2] - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \\ &= \mathbb{E}[\text{tr}((\mathbf{X}_1 \mathbf{X}_2)(\mathbf{X}_1 \mathbf{X}_2)^{\top})] - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \mathbb{E}[\text{tr}(\mathbf{X}_1 \mathbf{X}_2 \mathbf{X}_2^{\top} \mathbf{X}_1^{\top})] - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \\ &= \mathbb{E}[\text{tr}(\mathbf{X}_1^{\top} \mathbf{X}_1 \mathbf{X}_2 \mathbf{X}_2^{\top})] - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \text{tr}(\mathbb{E}[\mathbf{X}_1^{\top} \mathbf{X}_1 \mathbf{X}_2 \mathbf{X}_2^{\top}]) - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \\ &= \text{tr}(\mathbb{E}[\mathbf{X}_1^{\top} \mathbf{X}_1] \mathbb{E}[\mathbf{X}_2 \mathbf{X}_2^{\top}]) - \|\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2]\|_{\mathbb{F}}^2 = \|\mathbb{E}[\mathbf{X}_1^{\top} \mathbf{X}_1]\|_{\mathbb{F}}^2 - \|\mathbf{K}_{\alpha}(\mathbf{W})\|_{\mathbb{F}}^2 \end{aligned} \quad (14)$$

In the series of equalities above, we applied several facts:

1. unbiasedness of the estimator:  $\mathbb{E}[\mathbf{X}_1 \mathbf{X}_2] = \mathbf{K}_{\alpha}(\mathbf{W})$ ,
2. standard formula for the scalar variance:  $\mathbb{E}[(Z - \mathbb{E}[Z])^2] = \mathbb{E}[Z^2] - (\mathbb{E}[Z])^2$ , lifted to the matrix space via Frobenius norm,
3. the following formula:  $\|\mathbf{Z}\|_{\mathbb{F}}^2 = \text{tr}(\mathbf{Z}\mathbf{Z}^{\top})$ , where  $\text{tr}$  denotes *trace* of the input matrix,
4. cyclic property of the trace:  $\text{tr}(\mathbf{ABCD}) = \text{tr}(\mathbf{BCDA})$ ,
5. the symmetry of  $\mathbf{X}_1^{\top} \mathbf{X}_1$  and  $\mathbf{X}_2 \mathbf{X}_2^{\top}$ ,
6. the equality:  $\mathbb{E}[\mathbf{X}_1^{\top} \mathbf{X}_1] = \mathbb{E}[\mathbf{X}_2 \mathbf{X}_2^{\top}]$  that comes from the definition of  $\mathbf{X}_1$  and  $\mathbf{X}_2$ ,
7. independence of  $\mathbf{X}_1$  and  $\mathbf{X}_2$ .

$\square$

## A.3 PROOF OF THEOREM 3.3

We will now provide a proof of Theorem 3.3.

*Proof.* Without loss of generality, we will assume that  $m = 1$ . Take two vertices:  $i$  and  $j$  of a fixed graph  $G$ . We will consider two GRFs++ estimators of the value  $\mathbf{K}_\alpha(\mathbf{W})[i, j]$  of the graph kernel between them:  $\widehat{\mathbf{K}}_\alpha^{(2^t)}(\mathbf{W})[i, j]$  and  $\widehat{\mathbf{K}}_\alpha^{(2^{t+1})}(\mathbf{W})[i, j]$ , applying GRFs++ mechanism with degree  $2^t$  and  $2^{t+1}$  respectively (for  $t \geq 0$ ). Note that since both estimators are unbiased, it only suffices to prove the following:

$$\mathbb{E}[(\widehat{\mathbf{K}}_\alpha^{(2^{t+1})}(\mathbf{W})[i, j])^2] \leq \mathbb{E}[(\widehat{\mathbf{K}}_\alpha^{(2^t)}(\mathbf{W})[i, j])^2] \quad (15)$$

Note that estimator  $\widehat{\mathbf{K}}_\alpha^{(2^t)}(\mathbf{W})[i, j]$  can be re-written as:

$$\sum_{\substack{\omega \in \Omega(i, j) \\ i=p_0, v_1, p_1, \dots, v_{2^t}, p_{2^t}=j}} X_{f^{(2^t)}}^{(\omega)}(p_0, v_1) \dots X_{f^{(2^t)}}^{(\omega)}(p_{2^t-1}, v_{2^t}) X_{f^{(2^t)}}^{(\omega)}(p_1, v_1) \dots X_{f^{(2^t)}}^{(\omega)}(p_{2^t}, v_{2^t}), \quad (16)$$

where:

1.  $\Omega(i, j)$  is the set of all the walks between  $i$  and  $j$
2.  $p_0, \dots, p_{2^t}$  are some vertices (potentially with repetitions) from  $\omega$ , visited in that order along the walk  $\omega$ , as going from  $i$  to  $j$
3.  $X_f^{(\omega)}(a, b)$  is a random variable that is equal to

$$f(l(\omega(a, b))) W(a, b) \left( \prod_{v \in \omega(a, b)} \deg(v) \right) (1 - p_{\text{halt}})^{-l(\omega(a, b))}$$

(for  $\omega(a, b)$  denoting vertices on  $\omega$  from  $a$ , but not including  $b$ ,  $W(a, b)$  denoting the corresponding product of edge weights and  $l(\omega(a, b))$  being the number of edges of the part of  $\omega$  from  $a$  to  $b$ ) if a random walk from  $a$  reaches  $b$ , as a prefix of  $\omega$  starting from  $a$  and going to  $b$  and is zero otherwise.

4.  $f^{(l)}$  is a modulation function for the GRFs++ mechanism of degree  $l$ .

Similarly, one can write  $\widehat{\mathbf{K}}_\alpha^{(2^{t+1})}(\mathbf{W})[i, j]$  as:

$$\sum_{\substack{\omega \in \Omega(i, j) \\ i=p_0, v_1, p_1, \dots, v_{2^t}, p_{2^t}=j \\ u_1, u_2, \dots, u_{2^{t+1}-1}, u_{2^{t+1}}}} X_{f^{(2^{t+1})}}^{(\omega)}(p_0, v_1) \dots X_{f^{(2^{t+1})}}^{(\omega)}(p_{2^t-1}, v_{2^t}) X_{f^{(2^{t+1})}}^{(\omega)}(p_1, v_1) \dots X_{f^{(2^{t+1})}}^{(\omega)}(p_{2^t}, v_{2^t}) \\ X_{f^{(2^{t+1})}}^{(\omega)}(p_0, u_1) X_{f^{(2^{t+1})}}^{(\omega)}(p_1, u_3) \dots X_{f^{(2^{t+1})}}^{(\omega)}(p_{2^t-1}, u_{2^{t+1}-1}) \\ X_{f^{(2^{t+1})}}^{(\omega)}(v_1, u_1) X_{f^{(2^{t+1})}}^{(\omega)}(v_2, u_3) \dots X_{f^{(2^{t+1})}}^{(\omega)}(v_{2^t}, u_{2^{t+1}-1}) \\ X_{f^{(2^{t+1})}}^{(\omega)}(v_1, u_2) X_{f^{(2^{t+1})}}^{(\omega)}(v_2, u_4) \dots X_{f^{(2^{t+1})}}^{(\omega)}(v_{2^t}, u_{2^{t+1}}) \\ X_{f^{(2^{t+1})}}^{(\omega)}(p_1, u_2) X_{f^{(2^{t+1})}}^{(\omega)}(p_2, u_4) \dots X_{f^{(2^{t+1})}}^{(\omega)}(p_{2^t}, u_{2^{t+1}}) \quad (17)$$

Denote the sub-sum of the above sum, corresponding to the particular choice of:  $p_1, \dots, p_{2^t}, v_1, \dots, v_{2^t}$  as:  $\Psi(p_1, \dots, p_{2^t}, v_1, \dots, v_{2^t})$ .

It suffices to prove that for any two sequences  $p_1, p_2, \dots, p_{2^t}, v_1, \dots, v_{2^t}$  and  $p'_1, p'_2, \dots, p'_{2^t}, v'_1, \dots, v'_{2^t}$ , the following holds:

864  
865  
866  
867  
868  
869  
870

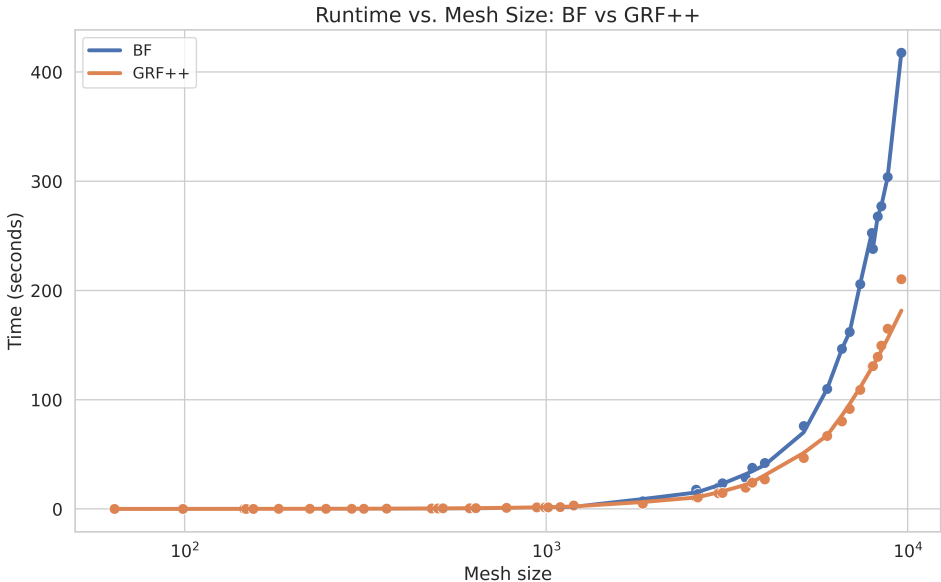
$$\mathbb{E}\left[\left(X_{f^{(2^t)}}^{(\omega)}(p_0, v_1) \dots X_{f^{(2^t)}}^{(\omega)}(p_{2^t-1}, v_{2^t}) X_{f^{(2^t)}}^{(\omega)}(p_1, v_1) \dots X_{f^{(2^t)}}^{(\omega)}(p_{2^t}, v_{2^t})\right) \left(X_{f^{(2^t)}}^{(\omega)}(p'_0, v'_1) \dots X_{f^{(2^t)}}^{(\omega)}(p'_{2^t-1}, v'_{2^t}) X_{f^{(2^t)}}^{(\omega)}(p'_1, v'_1) \dots X_{f^{(2^t)}}^{(\omega)}(p'_{2^t}, v'_{2^t})\right)\right] \geq \mathbb{E}\left[\Psi(p_1, \dots, p_{2^t}, v_1, \dots, v_{2^t}) \Psi(p'_1, \dots, p'_{2^t}, v'_1, \dots, v'_{2^t})\right] \tag{18}$$

871 This however follows from the convolutional properties of the modulation function  $f$  (Lemma  
872 2.1) and the fact that the product of two  $X$ -variables corresponding to some walk starting  
873 at some fixed vertex of a graph  $G$  is not identically zero if and only if one of the walks is a  
874 prefix of another one.  $\square$

875  
876 **B ADDITIONAL EXPERIMENTAL DETAILS**  
877

878 In this section, we provide additional details regarding the experimental setup and present  
879 additional results. Moreover we show a plot (Fig. 7) highlighting the speed gains by GRF++  
880 over the baseline kernel. At moderate mesh sizes (2K–4K), GRF++ gives  $\sim 1.7 \times -2.2 \times$   
881 speedup while for larger meshes (6K–10K), GRF++ consistently delivers  $\sim 3 \times$  faster run-  
882 time. Thus GRF++ scales significantly better as mesh complexity increases.

883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903



904 Figure 7: Time comparison between GRF++ and baseline diffusion kernel over various mesh sizes.  
905 GRF++ is significantly faster than the baseline as the graphs get larger.  
906

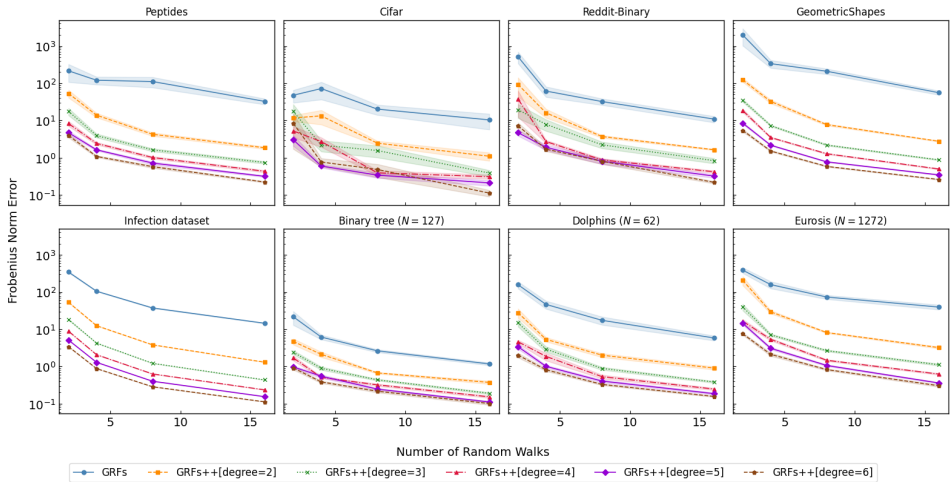
907  
908  
909 **B.1 ACCURATE ESTIMATION OF GRAPH KERNELS**

910 We follow the exact setup as (Reid et al., 2024b). For computational comparison we used a  
911 randomly generated connected graph with 500 nodes. To have fair comparison we derived  
912 the relevant  $p_{halt} = p_{base} * degree - of - kernel$ . We fixed the number of random walks to  
913 256.  
914

915 **B.1.1 EXPERIMENTS ON GRAPHS WITH LARGE DIAMETERS**  
916

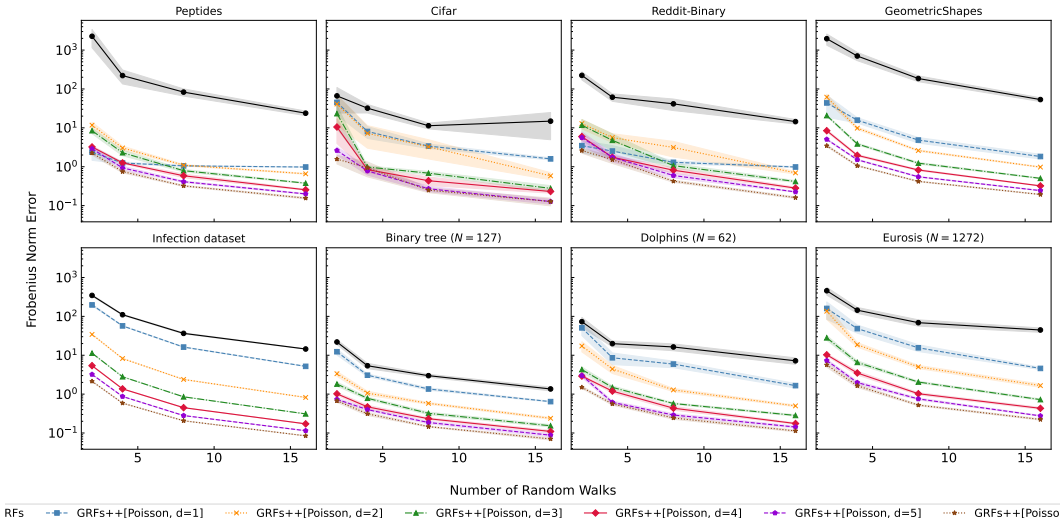
917 In this subsection, we provide details on the graphs used for estimating longer walks. For  
this task, we pick graphs from Peptides (Dwivedi et al., 2022), CIFAR-10 (Dwivedi et al.,

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933



934 Figure 8: Estimation of the kernel values for distant nodes for the diffusion kernel. GRF++  
935 provides a more accurate estimation in various graphs of large diameters.  
936

937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953



954 Figure 9: As in Fig. 8, but with Poisson termination strategy activated. This novel halting  
955 strategy, proposed in this paper, further improves approximation quality.  
956

957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

2020), Reddit-Binary (Morris et al., 2020), Geometric Shapes (Yannick-S, 2025) as well from the dataset considered by Reid et al. (2024b).

For each of these datasets, we remove isolated nodes and select the graphs with the largest diameters from the subset of the connected graphs. Finally to threshold the graphs to select the longer walks, we compute the shortest path distance via the Floyd-Warshall algorithm. We then select the pair of nodes where the walks are longer than the specified distance away. We then use this information to mask (i.e. zero out) all entries in the diffusion kernel.

We provide the walk threshold for these graphs in Table 3.

## B.2 GRAPH CLASSIFICATION EXPERIMENTS

In this subsection we provide additional details about our graph classification experiments. The statistics of our datasets is provided in Table 4 with additional details provided in (Morris et al., 2020). We follow the framework proposed by (Errica et al., 2020) to evaluate the

Table 3: Statistics of datasets used in experiments for estimate the accuracy to capture long walks. The walk length column refers to the fact that we are estimating all walks  $\geq k$ .

Dataset	Diameter	# Nodes	Walk Length (k)
Peptides	159	434	3
CIFAR	11	128	4
Reddit-Binary	19	436	5
GeometricShapes	28	864	5
Infection	4	500	4
Binary Tree	12	127	3
Dolphins	8	62	4
Eurosis	10	1272	4

Table 4: Statistics of the graph classification datasets used in this paper.

DATASETS	# Graphs	# Labels	Avg. # Nodes	Avg. # Edges	# Node Labels	# Node Attributes
MUTAG	188	2	17.93	19.79	7	-
PTC-MR	344	2	14.29	14.69	19	-
ENZYMES	600	6	32.63	62.14	3	18
PROTEINS	1113	2	39.06	72.82	3	1
D&D	1178	2	284.32	715.66	82	-
IMDB BINARY	1000	2	19.77	96.53	-	-
IMDB MULTI	1500	3	13.0	65.94	-	-
NCII	4110	2	29.87	32.30	37	-
REDDIT BINARY	2000	2	429.63	497.75	-	-
REDDIT MULTI-5K	4999	5	508.52	594.87	-	-

performance of the diffusion kernel as well as the approximate kernels obtained by GRF and GRF++. In particular, we use 10-fold cross-validation to obtain an estimate of the generalization performance of the methods.

Finally, we follow the approach by (de Lara & Pineau, 2018) to create graph features by using the smallest  $k$ -eigenvalues of the corresponding kernels. These features are then passed to a random forest classifier for classification.  $k$  is independently for the baseline as well as for GRF and GRF++.

We did a small hyperparameter sweep over  $\{.5, .6, .8, .9\}$  to find the width of the diffusion kernel. For GRF and GRF++, we fix the halting probability to be .1 and do a hyperparameter sweep over the number of walks. The degree of GRF++ is chosen to be 2.

Finally we posit our results in the context of various kernel methods as well as GNNs. In particular, we compare our methods against 20 popular kernels like (1) Vertex Histogram kernel (VH), (2) Random Walk kernel (RW), (3) Shortest Path kernel (SP), (4) Graphlet kernel (GR), (5) Weisfeiler-Lehman sub-tree kernel (WL-VH), (6) Weisfeiler-Lehman shortest path kernel (WL-SP), (7) Weisfeiler-Lehman pyramid match kernel (WL-PM), (8) Weisfeiler-Lehman optimal assignment kernel (WLOA), (9) Neighborhood Hash kernel (NH), (10) Neighborhood subgraph pairwise distance kernel (NSPDK), (11) Lovász  $\vartheta$  kernel (Lo- $\vartheta$ ), (12) SVM- $\vartheta$  kernel (SVM- $\vartheta$ ), (13) Ordered Decompositional DAGs with subtree kernel (ODD-STh), (14) Pyramid Match kernel (PM), (15) GraphHopper kernel (GH), (16) Subgraph Matching kernel (SM), (17) Propagation kernel (PK), (18) Multiscale Laplacian kernel (ML), (19) Core Weisfeiler-Lehman subtree kernel (CORE-WL-VH), and (20) Core Shortest Path kernel (CORE-SP). Note that this kernel computes relationship between graphs whereas our GRF++ computes relationship between nodes in a given graph. Finally we also compare our methods to popular GNNs like (1) DGCNN (Zhang et al., 2018), (2) GraphSAGE (Hamilton et al., 2017), (3) DiffPool (Ying et al., 2018), and (4) GIN (Xu et al., 2019). All these baseline numbers are taken from (Nikolentzos et al., 2021). Our results show that diffusion kernel is competitive across all the datasets. Our GRF++ mechanism maintains similar performance to that of the baseline diffusion kernel while being more computationally efficient.

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

Methods	DATASETS						
	MUTAG	ENZYMES	NCII	PTC-MR	D&D	PROTEINS	
Kernels	VH	69.1 ( $\pm$ 4.1)	20.0 ( $\pm$ 4.8)	55.7 ( $\pm$ 2.0)	57.1 ( $\pm$ 9.6)	74.8 ( $\pm$ 3.7)	71.1 ( $\pm$ 4.4)
	RW	81.4 ( $\pm$ 8.9)	16.7 ( $\pm$ 1.8)	TIMEOUT	54.4 ( $\pm$ 9.8)	OUT-OF-MEM	69.5 ( $\pm$ 5.1)
	SP	82.4 ( $\pm$ 5.5)	37.3 ( $\pm$ 8.7)	72.5 ( $\pm$ 2.0)	60.2 ( $\pm$ 9.4)	77.9 ( $\pm$ 4.5)	74.9 ( $\pm$ 3.6)
	WL-VH	86.7 ( $\pm$ 7.3)	50.7 ( $\pm$ 7.3)	85.2 ( $\pm$ 2.2)	64.9 ( $\pm$ 6.4)	78.7 ( $\pm$ 2.3)	76.2 ( $\pm$ 3.5)
	WL-SP	81.4 ( $\pm$ 8.7)	27.3 ( $\pm$ 7.4)	60.8 ( $\pm$ 2.4)	54.5 ( $\pm$ 9.8)	76.0 ( $\pm$ 3.5)	72.1 ( $\pm$ 3.1)
	WL-PM	88.3 ( $\pm$ 7.1)	57.5 ( $\pm$ 6.8)	85.6 ( $\pm$ 1.7)	65.1 ( $\pm$ 7.5)	OUT-OF-MEM	75.9 ( $\pm$ 3.8)
	WL-OA	87.2 ( $\pm$ 5.4)	58.0 ( $\pm$ 5.0)	86.3 ( $\pm$ 1.6)	65.7 ( $\pm$ 9.6)	77.6 ( $\pm$ 3.0)	76.2 ( $\pm$ 3.9)
	NH	88.3 ( $\pm$ 6.3)	54.5 ( $\pm$ 3.6)	84.7 ( $\pm$ 1.9)	63.4 ( $\pm$ 9.2)	74.6 ( $\pm$ 3.5)	75.0 ( $\pm$ 4.2)
	NSPDK	85.6 ( $\pm$ 8.9)	42.2 ( $\pm$ 8.0)	74.3 ( $\pm$ 2.1)	59.1 ( $\pm$ 7.3)	78.9 ( $\pm$ 4.7)	72.5 ( $\pm$ 2.9)
	ODD-STh	80.4 ( $\pm$ 8.8)	32.3 ( $\pm$ 4.8)	75.2 ( $\pm$ 2.0)	59.4 ( $\pm$ 9.8)	76.4 ( $\pm$ 4.5)	70.9 ( $\pm$ 4.1)
	PM	85.1 ( $\pm$ 5.8)	43.2 ( $\pm$ 5.3)	73.5 ( $\pm$ 1.9)	60.2 ( $\pm$ 8.2)	77.9 ( $\pm$ 3.7)	70.9 ( $\pm$ 4.4)
	GH	82.5 ( $\pm$ 5.8)	37.2 ( $\pm$ 6.6)	71.0 ( $\pm$ 2.3)	60.2 ( $\pm$ 9.4)	TIMEOUT	74.8 ( $\pm$ 2.4)
	SM	85.7 ( $\pm$ 5.8)	35.7 ( $\pm$ 5.5)	TIMEOUT	60.2 ( $\pm$ 6.8)	OUT-OF-MEM	OUT-OF-MEM
	PK	76.6 ( $\pm$ 5.2)	44.0 ( $\pm$ 6.3)	82.1 ( $\pm$ 2.1)	65.1 ( $\pm$ 5.6)	77.7 ( $\pm$ 4.2)	73.1 ( $\pm$ 4.7)
	ML	87.2 ( $\pm$ 7.5)	48.5 ( $\pm$ 7.8)	79.7 ( $\pm$ 1.8)	64.5 ( $\pm$ 5.8)	78.6 ( $\pm$ 4.0)	74.2 ( $\pm$ 4.4)
	CORE-WL-VH	85.6 ( $\pm$ 6.5)	51.7 ( $\pm$ 7.0)	85.2 ( $\pm$ 2.2)	65.5 ( $\pm$ 5.6)	79.5 ( $\pm$ 3.2)	76.5 ( $\pm$ 4.4)
CORE-SP	85.1 ( $\pm$ 6.8)	39.5 ( $\pm$ 9.3)	73.8 ( $\pm$ 1.4)	57.3 ( $\pm$ 9.7)	79.3 ( $\pm$ 3.8)	76.5 ( $\pm$ 3.9)	
GNNs	Diffusion	86.2 ( $\pm$ 4.2)	41.7 ( $\pm$ 3.6)	71.5 ( $\pm$ 1.3)	62.2 ( $\pm$ 5.4)	74.8 ( $\pm$ 2.2)	74.5 ( $\pm$ 2.4)
	GRF	86.7 ( $\pm$ 4.4)	39.8 ( $\pm$ 3.2)	70.0 ( $\pm$ 1.6)	59.7 ( $\pm$ 5.8)	73.4 ( $\pm$ 2.4)	73.7 ( $\pm$ 2.3)
	GRF++ (ours)	87.8 ( $\pm$ 5.2)	40.9 ( $\pm$ 3.9)	72.4 ( $\pm$ 1.2)	61.3 ( $\pm$ 5.7)	74.5 ( $\pm$ 2.6)	74.3 ( $\pm$ 2.8)
	DGCNN	84.0 ( $\pm$ 7.1)	46.3 ( $\pm$ 6.3)	76.4 ( $\pm$ 1.7)	59.5 ( $\pm$ 6.9)	76.6 ( $\pm$ 4.3)	73.2 ( $\pm$ 3.2)
	GraphSAGE	83.6 ( $\pm$ 9.6)	46.1 ( $\pm$ 5.4)	76.0 ( $\pm$ 1.8)	61.7 ( $\pm$ 4.9)	72.9 ( $\pm$ 2.0)	74.3 ( $\pm$ 3.8)
	DiffPool	79.8 ( $\pm$ 6.7)	50.7 ( $\pm$ 8.7)	76.9 ( $\pm$ 1.9)	61.1 ( $\pm$ 5.6)	75.0 ( $\pm$ 3.5)	72.5 ( $\pm$ 3.5)
	GIN	84.7 ( $\pm$ 6.7)	44.5 ( $\pm$ 4.1)	80.0 ( $\pm$ 1.4)	59.1 ( $\pm$ 7.0)	75.3 ( $\pm$ 2.9)	72.8 ( $\pm$ 3.6)

Table 5: Average classification accuracy ( $\pm$  standard deviation) on the 6 classification datasets containing node-labeled graphs.

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

Methods	DATASETS						
	IMDB BINARY	IMDB MULTI	REDDIT BINARY	REDDIT MULTI-5K	REDDIT MULTI-12K	COLLAB	
Kernels	VH	50.0 ( $\pm$ 0.0)	33.3 ( $\pm$ 0.0)	50.0 ( $\pm$ 0.0)	20.0 ( $\pm$ 0.0)	21.7 ( $\pm$ 1.5)	52.0 ( $\pm$ 0.1)
	RW	64.1 ( $\pm$ 4.5)	44.6 ( $\pm$ 4.1)	TIMEOUT	TIMEOUT	TIMEOUT	68.0 ( $\pm$ 1.7)
	SP	58.2 ( $\pm$ 4.7)	39.2 ( $\pm$ 2.3)	81.7 ( $\pm$ 2.5)	47.9 ( $\pm$ 1.9)	TIMEOUT	58.8 ( $\pm$ 1.2)
	GR	66.1 ( $\pm$ 2.7)	39.5 ( $\pm$ 2.7)	76.1 ( $\pm$ 2.6)	34.7 ( $\pm$ 2.0)	23.0 ( $\pm$ 1.4)	73.0 ( $\pm$ 2.0)
	WL-VH	70.7 ( $\pm$ 6.8)	51.3 ( $\pm$ 4.4)	67.8 ( $\pm$ 3.5)	50.5 ( $\pm$ 1.6)	38.7 ( $\pm$ 1.7)	78.3 ( $\pm$ 2.1)
	WL-SP	58.2 ( $\pm$ 4.7)	39.2 ( $\pm$ 2.3)	TIMEOUT	TIMEOUT	TIMEOUT	58.8 ( $\pm$ 1.2)
	WL-PM	73.6 ( $\pm$ 3.4)	49.1 ( $\pm$ 5.5)	OUT-OF-MEM	OUT-OF-MEM	OUT-OF-MEM	OUT-OF-MEM
	WL-OA	72.6 ( $\pm$ 5.5)	51.1 ( $\pm$ 4.3)	89.0 ( $\pm$ 1.3)	54.0 ( $\pm$ 1.2)	TIMEOUT	80.5 ( $\pm$ 2.0)
	NH	71.6 ( $\pm$ 4.5)	50.5 ( $\pm$ 5.0)	81.2 ( $\pm$ 2.0)	49.9 ( $\pm$ 2.4)	39.6 ( $\pm$ 1.4)	81.1 ( $\pm$ 2.4)
	NSPDK	67.4 ( $\pm$ 3.3)	44.6 ( $\pm$ 3.8)	TIMEOUT	TIMEOUT	TIMEOUT	TIMEOUT
	Lo- $\theta$	51.0 ( $\pm$ 4.2)	39.8 ( $\pm$ 2.6)	TIMEOUT	TIMEOUT	TIMEOUT	TIMEOUT
	SVM- $\theta$	52.3 ( $\pm$ 4.0)	39.5 ( $\pm$ 2.7)	74.8 ( $\pm$ 2.6)	31.4 ( $\pm$ 1.1)	22.9 ( $\pm$ 0.9)	52.0 ( $\pm$ 0.1)
	ODD-STh	65.0 ( $\pm$ 4.0)	46.7 ( $\pm$ 3.4)	52.1 ( $\pm$ 3.2)	43.1 ( $\pm$ 1.8)	30.0 ( $\pm$ 1.6)	52.0 ( $\pm$ 0.1)
	PM	66.3 ( $\pm$ 4.2)	46.1 ( $\pm$ 3.8)	86.5 ( $\pm$ 2.1)	48.3 ( $\pm$ 2.5)	41.1 ( $\pm$ 0.6)	74.0 ( $\pm$ 2.4)
	GH	59.4 ( $\pm$ 3.4)	39.5 ( $\pm$ 2.6)	TIMEOUT	TIMEOUT	TIMEOUT	60.0 ( $\pm$ 1.4)
	SM	TIMEOUT	TIMEOUT	OUT-OF-MEM	OUT-OF-MEM	OUT-OF-MEM	TIMEOUT
PK	51.7 ( $\pm$ 3.7)	34.5 ( $\pm$ 3.0)	63.9 ( $\pm$ 3.0)	34.9 ( $\pm$ 1.7)	23.9 ( $\pm$ 1.2)	57.0 ( $\pm$ 1.2)	
ML	69.9 ( $\pm$ 4.8)	47.7 ( $\pm$ 3.2)	89.4 ( $\pm$ 2.1)	35.4 ( $\pm$ 2.0)	OUT-OF-MEM	75.6 ( $\pm$ 1.6)	
CORE-WL-VH	73.5 ( $\pm$ 6.1)	51.7 ( $\pm$ 4.1)	73.0 ( $\pm$ 4.5)	51.1 ( $\pm$ 1.6)	40.2 ( $\pm$ 1.8)	84.5 ( $\pm$ 2.0)	
CORE-SP	68.5 ( $\pm$ 3.9)	51.0 ( $\pm$ 3.5)	91.0 ( $\pm$ 1.8)	TIMEOUT	OUT-OF-MEM	TIMEOUT	
GNNs	Diffusion	72.3 ( $\pm$ 1.8)	50.8 ( $\pm$ 2.0)	85.6 ( $\pm$ 1.2)	49.5 ( $\pm$ 0.4)	35.9 ( $\pm$ 0.7)	75.1 ( $\pm$ 0.2)
	GRF	68.9 ( $\pm$ 1.9)	48.7 ( $\pm$ 2.4)	84.0 ( $\pm$ 1.1)	48.6 ( $\pm$ 0.7)	35.1 ( $\pm$ 0.9)	70.8 ( $\pm$ 0.6)
	GRF++	71.2 ( $\pm$ 2.1)	49.1 ( $\pm$ 2.9)	84.4 ( $\pm$ 1.5)	49.2 ( $\pm$ 0.9)	36.8 ( $\pm$ 0.9)	73.6 ( $\pm$ 0.9)
	DGCNN	69.2 ( $\pm$ 3.0)	45.6 ( $\pm$ 3.4)	87.8 ( $\pm$ 2.5)	49.2 ( $\pm$ 1.2)	43.9 ( $\pm$ 1.0)	71.2 ( $\pm$ 1.9)
	GraphSAGE	68.8 ( $\pm$ 4.5)	47.6 ( $\pm$ 3.5)	84.3 ( $\pm$ 1.9)	50.0 ( $\pm$ 1.3)	43.5 ( $\pm$ 1.0)	73.9 ( $\pm$ 1.7)
	DiffPool	68.4 ( $\pm$ 3.3)	45.6 ( $\pm$ 3.4)	89.1 ( $\pm$ 1.6)	53.8 ( $\pm$ 1.4)	44.4 ( $\pm$ 1.4)	68.9 ( $\pm$ 2.0)
	GIN	71.2 ( $\pm$ 3.9)	48.5 ( $\pm$ 3.3)	89.9 ( $\pm$ 1.9)	56.1 ( $\pm$ 1.7)	48.3 ( $\pm$ 1.6)	75.6 ( $\pm$ 2.3)

Table 6: Average classification accuracy ( $\pm$  standard deviation) on the 6 classification datasets containing unlabeled graphs.

1075

1076

1077

1078

1079

Table 7: Final Ranking of Methods on the Unlabeled Datasets (Lower Average Rank is Better).

Rank	Method	Average Rank
1	GIN	4.50
2	CORE-WL-VH	5.00
3	NH	6.67
4	Diffusion	7.00
5	WL-OA	7.17
6	WL-VH	7.33
7	GraphSAGE	8.83
7	GRF++ (ours)	8.83
7	DiffPool	8.83
10	DGCNN	9.50
11	PM	10.33
12	GRF	11.00
13	ML	11.83
14	GR	15.83
15	ODD-STh	16.33
16	CORE-SP	16.50
17	SVM- $\vartheta$	19.17
17	PK	19.17
19	SP	19.33
20	WL-PM	19.50
21	VH	21.50
22	RW	22.17
23	GH	23.17
24	NSPDK	23.50
25	WL-SP	24.17
26	Lo- $\vartheta$	25.33
27	SM	27.00

### B.3 NODE CLUSTERING

We follow the same setup as (Reid et al., 2024b) except that the number of clusters is based upon the actual number of different classes. Thus we use  $p_{halt} = 0.1$ ,  $m = 16$ . To get details of the dataset, please see Ivashkin & Chebotarev (2016). In table 9, we compare GRF++ with 5 other methods : Louvain, Spectral, GRF, KCenters and Propagation. Our method is competitive across all datasets and outperforms all methods on 3 out of 5 datasets.

### B.4 VERTEX NORMAL PREDICTION EXPERIMENTS

In this sub-section, we present implementation details for vertex normal prediction experiments. All the experiments are run on free Google Colab with 12Gb of RAM.

For this task, we choose 40 meshes for 3D-printed objects of varying sizes from the Thingi10K dataset. Following (Choromanski et al., 2024), we choose the following meshes corresponding to the ids given by :

[60246, 85580, 40179, 964933, 1624039, 91657, 79183, 82407, 40172, 65414, 90431, 74449, 73464, 230349, 40171, 61193, 77938, 375276, 39463, 110793, 368622, 37326, 42435, 1514901, 65282, 116878, 550964, 409624, 101902, 73410, 87602, 255172, 98480, 57140, 285606, 96123, 203289, 87601, 409629, 37384, 57084]

We do a small search for the width of the kernel  $\sigma \in \{.5, .6, .8\}$  for the baseline runs. For both GRF and GRF++, the number of walks are chosen from the subset  $\{4, 8, 16\}$  and the halting probability of the walk is .1. The degree of GRF++ is chosen to be 2.

Table 8: Cosine Similarity for Meshes. GRF++ matches the performance of the baseline kernel (BF). GRF++r reuses the same walk and still outperform GRF.

MESH SIZE	64	99	146	148	155	182	222	246	290	313
BF	0.4255	0.7675	0.9424	0.4325	0.7095	0.9654	0.8715	0.7464	0.8895	0.5514
GRF	0.3083	0.6786	0.9348	0.3813	0.6831	0.9466	0.8569	0.6722	0.8651	0.5309
GRF++	0.3889	0.7163	0.9367	0.4434	0.6892	0.9611	0.8684	0.7377	0.8679	0.5432
GRF++r	0.3905	0.7163	0.9327	0.4435	0.6867	0.9564	0.8510	0.6878	0.8464	0.5178
MESH SIZE	362	482	502	518	614	639	777	942	992	1012
BF	0.5884	0.9830	0.8881	0.4956	0.9172	0.8958	0.8022	0.8559	0.7206	0.9366
GRF	0.5751	0.9737	0.8673	0.4486	0.8866	0.8739	0.7812	0.8369	0.6967	0.9136
GRF++	0.5821	0.9807	0.8843	0.4830	0.9084	0.8914	0.8039	0.8496	0.7144	0.9234
GRF++r	0.5688	0.9775	0.8772	0.4734	0.8982	0.8866	0.8019	0.8406	0.7085	0.9200
MESH SIZE	1094	1192	1849	2599	2626	2996	3072	3559	3715	4025
BF	0.9236	0.8297	0.9265	0.4987	0.8927	0.9326	0.4796	0.9356	0.9619	0.9669
GRF	0.9001	0.7987	0.9101	0.4065	0.8778	0.9171	0.4637	0.9208	0.9508	0.9588
GRF++	0.9170	0.8167	0.9202	0.4615	0.8820	0.9277	0.4731	0.9293	0.9546	0.9646
GRF++r	0.9098	0.8134	0.9146	0.4209	0.8730	0.9210	0.4606	0.9281	0.9561	0.9620
MESH SIZE	5155	5985	6577	6911	7386	7953	8011	8261	8449	8800
BF	0.9011	0.9194	0.9622	0.9769	0.9437	0.9460	0.9382	0.9196	0.9276	0.9836
GRF	0.8833	0.9091	0.9525	0.9682	0.9308	0.9383	0.9233	0.9050	0.9139	0.9778
GRF++	0.8931	0.9154	0.9599	0.9751	0.9374	0.9429	0.9321	0.9145	0.9205	0.9820
GRF++r	0.8896	0.9129	0.9561	0.9701	0.9348	0.9410	0.9269	0.9095	0.9160	0.9805

Table 9: Node Clustering: Comparing GRFs++ against various clustering algorithms. Our method outperforms all other baselines on 3 out of 5 datasets.

Name	# Nodes	# clusters	Propagation	KCenters	Spectral	Louvain	GRF	GRF++[d=2]
Karate	34	2	0.4011	0.2585	0.2585	0.4367	0.2995	<b>0.2585</b>
Dolphins	62	2	0.0936	0.0323	0.0408	0.3432	0.0635	<b>0.0323</b>
Polbooks	105	3	<b>0.0621</b>	0.0621	0.0692	0.1154	0.1060	0.1033
Football	115	12	0.1986	0.0477	0.0551	<b>0.0174</b>	0.0731	0.0362
Databases	1006	6	0.4762	0.4918	0.3214	0.5298	0.3528	<b>0.3001</b>
Eurosis	1272	13	0.7418	0.1891	0.1295	0.1599	0.2248	<b>0.1304</b>

## B.5 VISION TRANSFORMER ATTENTION MASKING

In this experiment, we investigate the utility of GRFs++ for introducing inductive biases into the self-attention mechanism of Vision Transformers (ViTs). We treat the input image patches as nodes in a regular 2D grid graph (lattice), where edges connect spatially adjacent patches. We employ GRFs++ (with degree  $l = 2$ ) to efficiently approximate the diffusion kernel matrix  $K_{\text{diff}}$  on this grid graph. This approximate kernel is subsequently normalized and utilized as a soft mask  $M$ , which is added to the standard attention logic:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^{\top}}{\sqrt{d}} + \lambda M \right) V \quad (19)$$

We choose  $\lambda = 0.5$  in our setup. Our empirical results on ImageNet show that the ViT-B-16 model with GRF++ masking achieves a top-1 accuracy of **80.71%**, surpassing the baseline ViT-B-16 performance of 80.16%.

## B.6 GRF++ BEYOND DIFFUSION KERNELS

In this section, we showcase the ability of GRFs++ to efficiently approximate graph node kernels beyond the diffusion kernel. We consider the following kernels: (1)  $\mathbf{K}(\mathbf{W}) = (\mathbf{I} - \mathbf{W})^{-4}$ , (2)  $\mathbf{K}(\mathbf{W}) = (\mathbf{I} - \mathbf{W}^2)^{-4}$ , (3)  $\mathbf{K}(\mathbf{W}) = \exp(\mathbf{W}^2)$ , with the corresponding modulation functions being given by Taylor-series expansions of the following matrix functions: (1)  $\mathbf{W} \rightarrow (\mathbf{I} - \mathbf{W})^{-1}$ , (2)  $\mathbf{W} \rightarrow (\mathbf{I} - \mathbf{W}^2)^{-1}$ ,  $\mathbf{W} \rightarrow \exp(\frac{\mathbf{W}^2}{4})$ . The results are presented in Fig. 10,11,12. GRFs++ outperform regular GRFs for all three kernels.

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

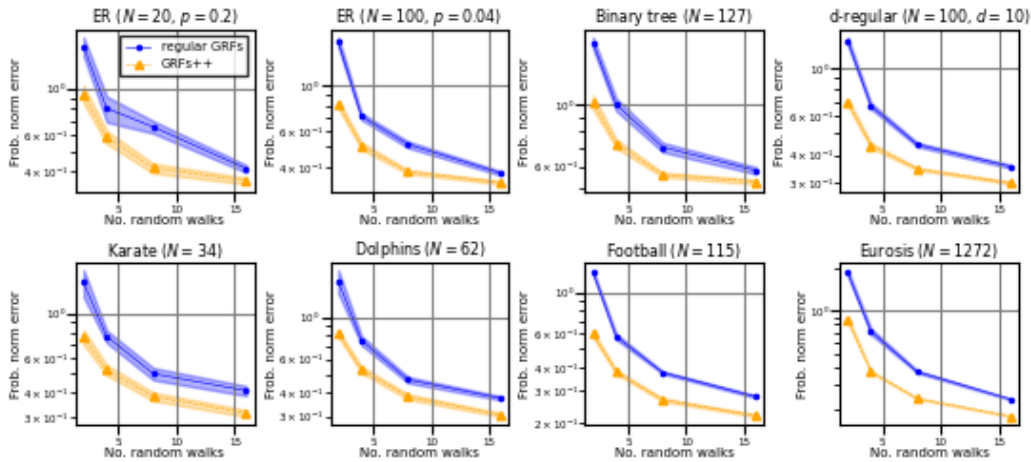


Figure 10: Experiment analogous to this from Fig. 2, but with a fixed degree  $l = 2$  used in GRFs++ and graph kernel of the form  $\mathbf{K}(\mathbf{W}) = (\mathbf{I} - \mathbf{W})^{-4}$ .

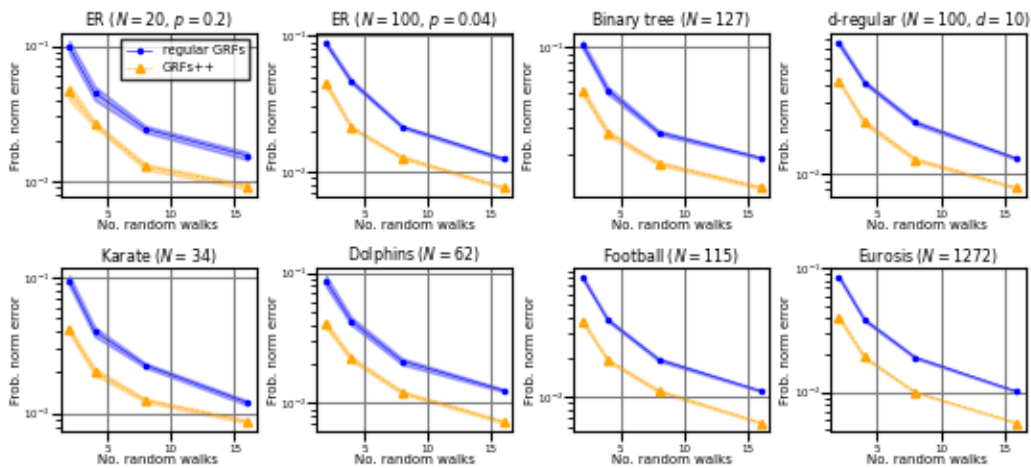


Figure 11: Experiment analogous to this from Fig. 2, but with a fixed degree  $l = 2$  used in GRFs++ and graph kernel of the form  $\mathbf{K}(\mathbf{W}) = (\mathbf{I} - \mathbf{W}^2)^{-4}$ .

## B.7 CODE

The code is available at [https://anonymous.4open.science/r/super\\_random\\_features\\_graphs-2782/README.md](https://anonymous.4open.science/r/super_random_features_graphs-2782/README.md)

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

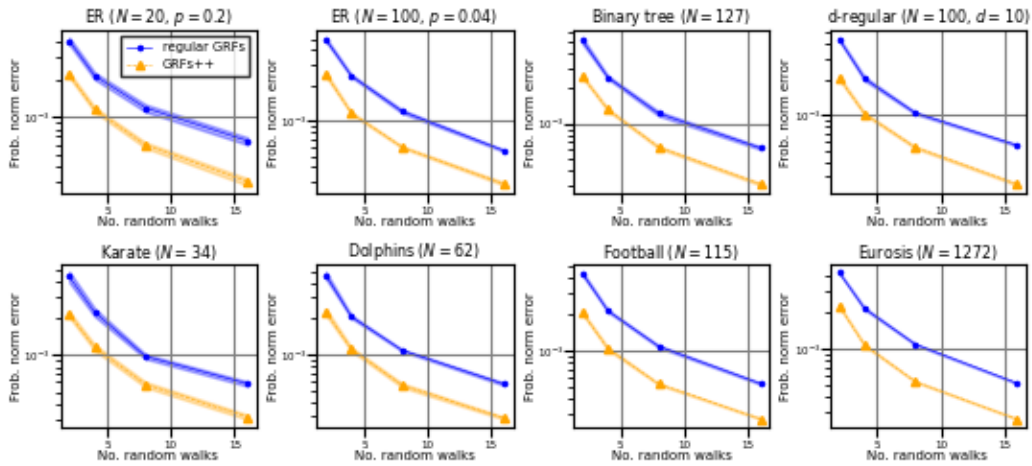


Figure 12: Experiment analogous to this from Fig. 2, but with a fixed degree  $l = 2$  used in GRFs++ and graph kernel of the form  $\mathbf{K}(\mathbf{W}) = \exp(\mathbf{W}^2)$ .