# Graph Kernel Attention Transformers

**Krzysztof Choromanski** *
Google Brain Robotics & Columbia University
kchoro@google.com

**Han Lin** *
Columbia University
hl3199@columbia.edu

**Haoxian Chen** *
Columbia University
hc3136@columbia.edu

**Jack Parker-Holder**
University of Oxford
jack.parker-holder@spc.ox.ac.uk

## Abstract

We introduce a new class of graph neural networks (GNNs), by combining several concepts that were so far studied independently - graph kernels, attention-based networks with structural priors and more recently, efficient Transformers architectures applying small memory footprint implicit attention methods via low rank decomposition techniques. The goal of the paper is twofold. Proposed by us Graph Kernel Attention Transformers (or GKATs) are much more expressive than SOTA GNNs as capable of modeling longer-range dependencies within a single layer. Consequently, they can use more shallow architecture design. Furthermore, GKAT attention layers scale **linearly** rather than quadratically in the number of nodes of the input graphs, even when those graphs are dense, requiring less compute than their regular graph attention counterparts. They achieve it by applying new classes of graph kernels admitting random feature map decomposition via random walks on graphs. As a byproduct of the introduced techniques, we obtain a new class of learnable graph sketches, called *graphots*, compactly encoding topological graph properties as well as nodes' features. We conducted exhaustive empirical comparison of our method with **nine** different GNN classes on tasks ranging from motif detection through social network classification to bioinformatics challenges, showing consistent gains coming from GKATs.

## 1 Introduction

Graphs are prevalent in a variety of real-world datasets. From social networks to bioinformatics (individual proteins, protein-protein interactions and even genomics), they play prominent roles in several high impact settings, while also making it possible to model relative inductive bias [37] and define non-Euclidean metrics that are relevant in navigation and planning problems in Robotics (e.g. probabilistic roadmaps) [24]. As such, there has recently been significant interest in Graph Neural Networks (GNNs) [3, 35, 20, 4, 8, 30], specifically designed for processing graph-structured data.

Despite the tremendous success of modern GNNs in understanding graph data, there are still several challenges with processing it efficiently. Classic *spectral methods* suffer from computational complexity that is cubic in the number of graph vertices (spectral analysis of graphs' Laplacians) which becomes a problem when the graphs under consideration are larger. One of the possible strategies to mitigate this is to apply a rich theory of spectral sketches [15], yet in practice the computed sketches need to be very large. Conversely, algorithms working in a *spatial domain* avoid expensive spectral computations, yet in order to model longer-range dependencies, they rely on deep GNN architectures for the signal to propagate from distant nodes since individual layers model only local interactions.

---

*equal contribution

Furthermore, for dense graphs they require quadratic memory just to store adjacency matrices of the input graphs for all graph-related computations (i.e.to compute attention layers as in GATs [35]).

Given these contrasting issues, the key question we seek to answer in this paper is as follows:

*Is it possible to design GNNs with densified individual layers modeling explicitly longer-range node-to-node relationships in a graph, that enable more shallow architectures and at the same time scale up to larger (not necessarily sparse) graphs?*

We answer it affirmatively by introducing Graph Kernel Attention Transformers (GKATs). GKATs take inspiration from recent research on dense linear attention Transformers [6, 18, 22, 17, 27], which have demonstrated the advantages of kernel-based approaches over sparse attention layers. GKATs model graph attention within each layer as a Hadamard product of the kernel matrix of the nodes' feature vectors and a graph kernel matrix, where the latter kernel is defined on graph **nodes**. GKATs' attention matrices combine a continuous signal coming from the vertices of the graph with a discrete graph topology signal, modulating standard graph-agnostic attention on nodes' representations via the more refined (vs. neighborhood relationship) topological signal encoded by a graph kernel.

GKATs have several advantages over previous methods: **(1)** they are capable of modeling dependencies between distant nodes within a single attention layer, **(2)** they can do this in **linear** time & memory, not requiring explicitly storing graph representations (such as adjacency matrices) to compute attention layers if a graph kernel has finite (at least on expectation) (random) feature representation, **(3)** they are highly flexible since they can be used with various graph kernels. In the paper, we will show how to design expressive and efficiently computable graph kernels with finite (random) feature representations and how they relate to well-established graph kernels, such as the graph diffusion kernels, that are notoriously difficult to apply due to their cubic computational cost.

**Graphots:** Such a re-designing of GNNs has profound consequences. In GKATs, the embedding of every node in any layer can be obtained by a simple linear transformation from a latent graph state in that layer, which we call a *graphot*. Even more importantly, graphots' sizes **do not depend** on the number of the nodes of the graph and play the role of learnable general use-case graph sketches. They can be interpreted as compact associative memories [21] storing compressed representations of all the nodes in the graph as well as their mutual relationships determined by graph's topology.

We confirm our theoretical results empirically, where GKATs consistently beat their competitors, including **nine** other GNN classes. As we show in Section 5, this performance is consistent on a wide range of tasks from motif detection, through social neetwork classification to bioinformatics.

## 2   Related Work

The number of different GNN architectures is vast [3, 35, 20, 4, 8, 30], yet most of them admit a modular structure of distinct blocks responsible for: (1) signal propagation between nodes (defining an information aggregation mechanism), (2) sampling nodes for massive graphs that cannot be entirely stored in memory and (3) pooling for hierarchical representations. The propagation modules most often can be categorized as working in spectral or spatial domain. In the former, a graph signal is translated into spectral domain via graph Fourier transform (FT), convolution is applied via the matrix of eigenvectors of the normalized graph Laplacian and the resulting signal is translated back to the spatial domain via inverse FT [31, 5, 4]. In the latter, then convolution is applied directly on the graph via weighted adjacency matrices. Examples include: learnable graph convolutional networks (LGCNs) [12], diffusion convolutional neural networks [1], Neural FPs [10] and more.

Another line of work on GNNs incorporates attention mechanisms that were widely popularized by Transformers [34, 26] which became SOTA architectures for processing sequential data. Papers on that subject build on an intrinsic connection between Transformers and GNNs. The former can be thought of as special instantiations of GNNs with the corresponding fully-connected graph topologies. Graph attention neural networks (GATs) [35, 38, 42] leverage this connection by replacing a regular attention matrix modeling relationships between all the tokens/nodes by the one with sparsity priors determined by the topology of the input graph. Thus in a fixed layer, every node attends only to its neighbors. As Transformers, GATs apply multi-head attention mechanism to robustify training.

# 3 Preliminaries

We consider weighted undirected graphs $G = (V, E, W, H)$, where $V$ is the set of nodes/vertices, $E$ is the edge set, $W = \{w_{i,j} : \{i,j\} \in E\}$ is the set of edge weights and $H = \{h_i \in \mathbb{R}^d : i \in V\}$ is the set of nodes' feature vectors. We will denote by $\mathrm{Adj}(G)$ the (weighted) adjacency matrix of $G$.

A regular graph attention mechanism (GAT) [35] processes such a data with standard attention layers [34], where the attention matrix $\mathbf{A}_i$ of the $i$th layer is given as $\mathbf{A}_i = \mathrm{Adj}(G) \odot \mathbf{A}_H^i$ (for the Hadamard/element-wise product $\odot$), $\mathbf{A}_H^i$ is the attention matrix for nodes' features, i.e.

$$\mathbf{A}_H^i(k, l) = \frac{\exp(\mathbf{q}_k^i (\mathbf{k}_l^i)^\top) w_{k,l}}{\sum_{\{k,r\} \in E} \exp(\mathbf{q}_k^i (\mathbf{k}_r^i)^\top) w_{k,r}}, \tag{1}$$

and $\mathbf{q}_t^i, \mathbf{k}_t^i \in \mathbb{R}^d$ for $t \in V$ (row-vector representations of queries/keys respectively) are learnable linear projections of nodes' feature vectors $\mathbf{h}_t^i$ in the $i$th layer. Thus the topology of the graph provides an inductive bias defining explicitly modeled relationships within a layer (those between neighboring tokens) and thus consequently introducing sparse attention patterns. The output of the attention module is computed as: $\mathbf{A}_i \mathbf{V}^i$, where the rows $\mathbf{v}_t^i \in \mathbb{R}^d$ for $t \in V$ of $\mathbf{V} \in \mathbb{R}^{N \times d}$ are other learnable linear projections of the feature vectors $\mathbf{h}_t^i$, and $N = |V|$ is the number of the vertices of the graph.

**Cost of the regular attention:** That cost is $O(Md)$, where $M = |E|$ is the number of edges and the required memory is $O(M + Nd)$ (note that depending on the sparsity of the graph, a matrix can be stored explicitly as an adjacency matrix or as a list of $N$ vertex adjacency lists). Thus time and space complexity is **quadratic** in the number of nodes for dense graphs with $M = \Omega(N^2)$.

# 4 GKATS: Towards Expressive and Scalable Graph Processing

## 4.1 Replacing Local with Longer-Range Attention for Graphs

The fundamental idea underpinning the GKAT architecture is the observation that the attention matrix $\mathbf{A}_i$ from GAT can be "densified" as follows:

$$\mathbf{A}_i(k, l) = \frac{K(\mathbf{q}_k, \mathbf{k}_l) T(k, l)}{\sum_{r \in V} K(\mathbf{q}_k, \mathbf{k}_r) T(k, r)}, \tag{2}$$

where $K : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ is a kernel defined on feature vectors in nodes of the graph and $T : V \times V \to \mathbb{R}$ is another kernel defined on the nodes of the graph $G$ (that for simplicity will depend only on the topology of the graph, but not feature vectors in those nodes). In particular, the former kernel can be a softmax kernel given as $K_{\mathrm{smax}}(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x} \mathbf{y}^\top)$ that is already used in the regular attention.

This method replaces the discrete mechanism, modulating attention on features in nodes based on filtering out nonadjacent pairs of nodes, by a smooth one, where filtering is substituted with weighting via a graph kernel $T$. The smoothing enables the attention to model longer-range dependencies within a single layer since it allows distant nodes in the graph to directly attend to each other.

Note that a GAT method is a special instantiation of that paradigm with $T_{\mathrm{GAT}} : V \times V \to \mathbb{R}$ defined as $T_{\mathrm{GAT}}(k, l) = \phi_{\mathrm{GAT}}(k) \widehat{\phi}_{\mathrm{GAT}}^\top(l)$, where $\phi_{\mathrm{GAT}}(i), \widehat{\phi}_{\mathrm{GAT}}(i) \in \mathbb{R}^M$ are row-vectors indexed by the edges of $G$, given as: $\phi_{\mathrm{GAT}}(i)(\{u,v\}) = 1[u = i]\sqrt{w_{u,v}}$, $\widehat{\phi}_{\mathrm{GAT}}(i)(\{u,v\}) = 1[v = i]\sqrt{w_{u,v}}$.

**Graph Diffusion Kernels:** A prominent example of kernels defined on graph nodes is the class of the so-called *graph diffusion kernels* $T_{\mathrm{diff}}^\lambda$ (GDKs) [33] with kernel matrices $\mathcal{T}_{\mathrm{diff}}^\lambda \overset{\mathrm{def}}{=} [T_{\mathrm{diff}}^\lambda(k, l)]_{k,l \in V}$ given as: $\mathcal{T}_{\mathrm{diff}}^\lambda = \exp(\lambda \mathbf{M})$ for a kernel hyperparameter $\lambda > 0$ and, where $\mathbf{M}$ is either a (weighted) adjacency matrix of the input graph or (for the spectral domain variant) its Laplacian (here $\exp$ stands for matrix exponentiation).

## 4.2 Decomposable Longer-Range Attention for Graphs

Even though densified graph attention from Sec. 4.1 does not need several layers to communicate signal between nonadjacent nodes, at first glance it seems that it is at the price of substantial extra cost per layer. It requires the computation of the graph kernel matrix $\mathcal{T} = [T(k, l)]_{k,l \in V}$ which is of cubic in $N$ time complexity for most interesting graph kernels (in particular graph diffusion kernels).
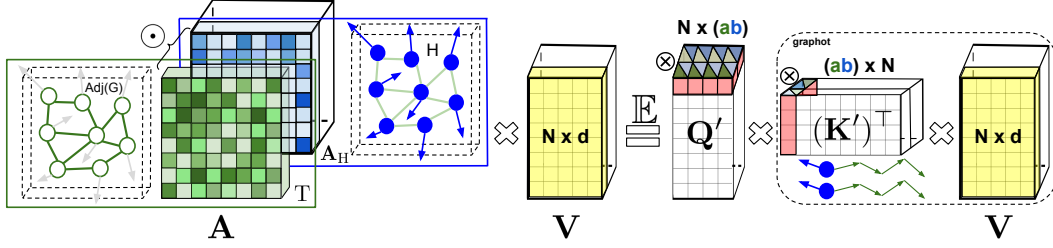
Figure 1: Decomposable longer-range attention in GKATs. For the clarifty of the exposition, we skip row-normalization of the attention matrix $\mathbf{A}$. Matrix $\mathbf{A}$ is a Hadamard product of two matrices: $\mathbf{A}_H$ using kernel K defined on query-key pairs of the nodes' emebeddings from H and kernel matrix corresponding to the graph kernel T. If both K and T admit finite kernel feature decomposition (at least on expectation), $a$ features corresponding to T and $b$ corresponding to K can be used to construct $a \times b$ cartesian features (the highlighted slices of $\mathbf{Q}'$ and $\mathbf{K}'$). These, via matrix associativity property, lead to attention computation scheme avoiding explicit calculation of $\mathbf{A}$ and instead relying on the graph sketch called by us a *graphot*.

It was observed in [6] that if a regular (non-graph) attention kernel $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ (generalizing softmax-kernel $K_{\text{soft}}(\mathbf{x}, \mathbf{y}) = \exp(\mathbf{x}\mathbf{y}^\top)$ from Eq. 1) satisfies: $K(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\phi(\mathbf{x})\phi(\mathbf{y})^\top]$ for arbitrary row-vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and some (deterministic or randomized) mapping $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ then an attention matrix $\mathbf{A}$ leveraging that kernel can be approximated as: $\widehat{\mathbf{A}} = \text{Diag}^{-1}(\mathbf{Q}'(\mathbf{K}')^\top \mathbf{1}_N)\mathbf{Q}'(\mathbf{K}')^T$, where $N$ stands for the length of the attention input sequence, $\mathbf{1}_N \in \mathbb{R}^N$ is the all-one vector, $\text{Diag}(\mathbf{r})$ stands for the diagonal matrix with diagonal $\mathbf{r}$ and furthermore the rows of $\mathbf{Q}', \mathbf{K}' \in \mathbb{R}^{N \times m}$, are $\phi$-transformations of queries and keys respectively. Furthermore, $\mathbf{Q}'(\mathbf{K}')^\top$ is an unbiased approximation of the unnormalized attention matrix $\mathbf{A}_{\text{unnorm}} = [K(\mathbf{q}_i, \mathbf{k}_j)]_{i,j=1,...,N}$. Then the output of the attention block can be approximated as: $\text{Diag}^{-1}(\mathbf{Q}'((\mathbf{K}')^\top \mathbf{1}_N))\mathbf{Q}'((\mathbf{K}')^T \mathbf{V})$, where brackets indicate the order of computations. Thus, the attention matrix is never explicitly materialized and space/time complexity is only linear in $N$. It is also linear in $m$, but as long as $m \ll N$, this mechanism brings computational gains. We call this method an *implicit attention*.

We observe that a function $K_{\text{prod}}(k, l) : V \times V \rightarrow \mathbb{R}$ defined as: $K_{\text{prod}}(k, l) \stackrel{\text{def}}{=} K(\mathbf{q}_k, \mathbf{k}_l)T(k, l)$ is a kernel that can be expressed as:

$$K_{\text{prod}}(k, l) = \mathbb{E}[(\phi_K(\mathbf{q}_k) \otimes \phi_T(k))(\phi_K(\mathbf{k}_l) \otimes \phi_T(l))^\top] \tag{3}$$

if $K(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\phi_K(\mathbf{x})\phi_K(\mathbf{y})]$ for $\phi_K : \mathbb{R}^d \rightarrow \mathbb{R}^{m_1}$ and $T(k, l) = \mathbb{E}[\phi_T(k)\phi_T(l)]$ for $\phi_T : V \rightarrow \mathbb{R}^{m_2}$, $\phi_K$ and $\phi_T$ are constructed independently and where $\mathbf{u} \otimes \mathbf{v}$ defines row-vectorized cartesian product of $\mathbf{u}$ and $\mathbf{v}$. This is a standard finite cartesian features mechanism for the compositional kernels [7]. We conclude that $K_{\text{prod}}$ admits (on expectation) a finite kernel feature decomposition required for the implicit attention.

Note that the GKAT mechanism presented in Eq. 2 is an attention method using kernel $K_{\text{prod}}$ and thus it can leverage computationally efficient implicit attention mechanism (see also: Fig. 1), as long as one can find finite mappings $\phi_1$ and $\phi_2$. Efficient mappings $\phi_1$ exist for several kernels used in attention in practice (often defining them), in particular for the most popular softmax kernel [6]. Next we show how to define expressive kernels on graph nodes that lead to efficient mappings $\phi_2$.

### 4.3 Random Walks Graph-Nodes Kernels

Instead of working directly with graph diffusion kernels defined in Sec. 4.1, we will use introduced by us class of the so-called *Random Walks Graph-Nodes Kernels* or RWGNKs, defined on graph **nodes** (as opposed to exhaustively explored random walk graph kernels defined on graphs [16]). Intuitively, the value of the RWGNK for two nodes is given as a dot-product of two *frequency vectors* that record visits in graph nodes of random walks beginning in the two nodes of interest. More formally, for the hyperparameters $\lambda, \alpha \geq 0$, and two random walks $\omega(k), \omega(l)$ with stopping probability $0 \leq p \leq 1$ (or of a fixed length) starting at $k$ and $l$ respectively, the RWGNK is given as:

$$K_{\text{RWGNK}}^{\lambda, \alpha, p}(k, l) = \frac{\mathbb{E}_{\omega(k)}[f_k^{\omega(k), \lambda}]}{\|\mathbb{E}_{\omega(k)}[f_k^{\omega(k), \lambda}]\|_2^\alpha} \left( \frac{\mathbb{E}_{\omega(l)}[f_l^{\omega(l), \lambda}]}{\|\mathbb{E}_{\omega(l)}[f_l^{\omega(l), \lambda}]\|_2^\alpha} \right)^\top. \tag{4}$$

4

The (row) frequency vector $f_h^{\omega(h),\lambda}$ for $h \in V$ can be defined in various ways. In its base variant $f_h^{\omega(h),\lambda} \in \mathbb{R}^N$ is given as $f_h^{\omega(h),\lambda}(i) \stackrel{\text{def}}{=} \sum_{l \in L^{\omega(h)}(i)} \lambda^l$, where $L^{\omega(h)}(i)$ is the set of lengths of those *prefix sub-walks* of a given random walk $\omega(h)$ that end at $i$ (where the prefix sub-walk of the walk $(j_1, j_2, ..., j_t)$ is any walk of the form $(j_1, ..., j_r)$ for some $r \leq t$ or an empty walk). Note that Eq. 4 leads to the desired representation of $K_{\text{RWGNK}}^\lambda$ as $K_{\text{RWGNK}}^\lambda(k,l) = \Psi(k)\Psi(l)^\top$, where $\Psi(h)$ is the renormalized expected frequency vector. In practice expectations are replaced by Monte Carlo samplings over few random walks and vectors $\Psi(h)$ are not stored explicitly, but in the form of weighted lookup tables. In the *anchor-points* variant, frequency vectors $f_h^{\omega(h),\lambda} \in \mathbb{R}^s$ are indexed by a selected (e.g. randomly) set of $s$ anchor-points in the graph. The advantage of the anchor-point version is that no lookup tables are needed for efficient implementation and decomposable attention can be realized simply via (heavily optimized on TPUs) dense matrix-matrix multiplication operations.
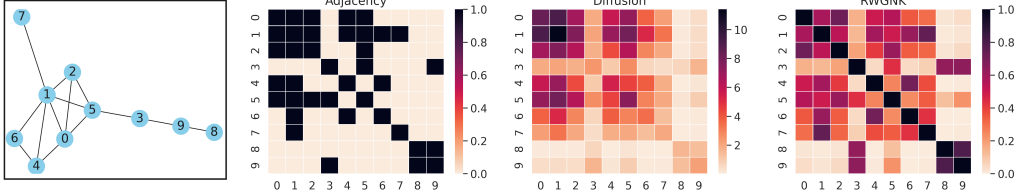


Figure 2: From left to right: unweighted graph G, its adjacency matrix $\text{Adj}(G)$, its GDK matrix $\exp(\text{Adj}(G))$ and RWGNK-matrix with walk length of 3 and $\alpha = 1$. Colored cells measure the relationship among pairs of nodes (darker is stronger). Last two matrices can be thought of as continuous smoothings of $\text{Adj}(G)$.

Figure 2 shows an intuition behind RWGNKs. Attending only to neighbors within a fixed layer via adjacency matrices limits the network's capability of aggregating fast distinct signal. Unlike adjacency matrix, diffusion kernel reveals long-range dense relationships among all the nodes. Equipped with continuous rather than binary values, it can also quantify their strengths. However, it requires cubic computation time $O(N^3)$. Our RWGNKs model dense attention but scale only linearly in $N$.

Next we explore properties of RWGNKs and their connections with GDKs (proof in the Appendix). We denote by $d_{\max}(G), d_{\min}(G)$ the maximum and minimum degree of a vertex in G respectively.

**Theorem 1** (RWGNKs count discounted numbers of walks). *The following is true for the kernel matrix* $\mathcal{K}_{\text{RWGNK}}^{\lambda,\alpha,p}(G) = [K_{\text{RWGNK}}^{\lambda,\alpha,p}(k,l)]_{k,l \in V(G)}$ *of the* RWGNK *kernel with* $0 \leq \lambda \leq 1$, $\alpha = 0$ *and* $0 < p < 1$ *for a graph* G *with vertex set* V(G) *of size* N *(element-wise matrix inequality):*

$$\Gamma\left(\frac{(1-p)\lambda}{d_{\max}(G)}\text{Adj}(G)\right) \leq \mathcal{K}_{\text{RWGNK}}^{\lambda,\alpha,p}(G) \leq \Gamma\left(\frac{(1-p)\lambda}{d_{\min}(G)}\text{Adj}(G)\right), \tag{5}$$

where $\Gamma(\mathbf{A}) = \sum_{i=0}^{\infty}(i+1)\mathbf{A}^i$. Using the fact that $\text{Adj}^i(G)$ encodes the number of walks of length $i$ between pairs of vertices in G, we conclude that $K_{\text{RWGNK}}^{\lambda,0,p}(k,l) = \sum_{i=0}^{\infty} c_{k,l}^i r_{k,l}(i)$, where: $r_{k,l}(i)$ is the number of walks of length $i$ between nodes: $k$ and $l$ and $\frac{\sqrt[i]{i+1}(1-p)\lambda}{d_{\max}(G)} \leq c_{k,l} \leq \frac{\sqrt[i]{i+1}(1-p)\lambda}{d_{\min}(G)}$. Note that GDK with parameter $\lambda$ satisfies: $\text{GDK}^\lambda(k,l) = \sum_{i=0}^{\infty} \tilde{c}^i(k,l) r_{k,l}(i)$, where: $\tilde{c}(k,l) = \frac{\lambda}{\sqrt[i]{i!}}$. In practice, it suffices to have random walks of fixed length (instead of taking $p > 0$) (see: Sec 5). Furthermore, by taking $\alpha > 0$ (e.g. $\alpha = 1$) we can guarantee that kernel values are bounded.

### 4.4 GKAT Algorithm

We are ready to summarize GKAT's graph attention layers (for the complete algorithmic box, see: Appendix, Sec. 7.2). The algorithm is parameterized by two kernels: $K : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$, $T : V \times V \to \mathbb{R}$. Graph kernel $T$ is chosen to be a RWGNK kernel (as in Sec. 4.3), parameterized by discount factor $\lambda > 0$, $\alpha \geq 0$ and stopping probability $0 < p < 1$ (or alternatively fixed walk length $\tau$). For each random walk, in each node with probability $1 - p$ its neighbor is chosen uniformly at random for the next step. Furthermore, $K$ is chosen to admit (potentially on expectation) finite (potentially random) kernel feature decomposition via mapping $\phi_K$ (e.g. regular softmax attention [6]). For smaller graphs, the attention matrix $\mathbf{A}$ can be computed explicitly as in Eq. 2 and the layer outputs $\mathbf{AV}$ (in that setting, in principle any graph kernel defined on nodes can be applied, e.g. GDK). For larger graphs, $\mathbf{A}$ is not explicitly materialized, but instead the output of the attention module is calculated, as explained in Sec. 4.2, with the use of introduced there cartesian features (for $\phi_T(l) = \Psi(l)$ as in Sec. 4.3). Full GKAT's architecture consists of several blocks, where each block

(as in the regular Transformer [34]), is built from the attention layer and standard MLP layer. In practice, because of longer-range dependencies modeling within individual attention layers by GKAT, the total number of blocks can be small (see: Sec. 5), resulting in more shallow end-to-end pipelines.

**Space & Time Complexity:** Denote by $\tau$ the average length of the random walk and by $t$ the number of random walks used per node. Furthermore, let $m$ be: the average number of nodes visited from a given node in a graph by conducting $t$ random walks for the no-anchor-points version and the number of anchor-points for the anchor-points variant. For the former, we clearly have: $m \leq \tau t$. Computing all random walks takes time $O(N\tau t)$ (sampling uniformly next node of a random walk can be done in constant time regardless of a degree of a current node). This is a **one-time cost** per input graph. For the anchor-points version, there is an additional **one-time cost** of computing anchor points ($O(m)$ if those are chosen uniformly at random). Vectors $\Psi$ can be efficiently stored as weighted lookup tables of average length $m$, so the corresponding memory usage is $O(Nm)$. Given vectors $\Psi$, efficient attention computation described in Sec. 4.2 can be conducted in time $O(Nmrd)$ and space $O(Nmr + Nd + mrd)$, where $d$ is the dimensionality of nodes' embeddings and $r$ is the number of features used to linearize kernel K. Thus if $\tau, t, m, r$ are small constants, space and time complexity is effectively **linear** instead of quadratic in $N$ (see: Section 5).

### 4.5 Graphots - New Neural Network Graph Sketches

Note that our decomposable longer-range attention from Sec. 4.2 can be interpreted as a graph compression technique. The new embeddings of nodes can be directly computed be left-multiplying a matrix $(\mathbf{K}')^\top \mathbf{V} \in \mathbb{R}^{m \times d}$ and a vector $(\mathbf{K}')^\top \mathbf{1}_N \in \mathbb{R}^m$ (storing approximate attention renormalization terms) by the rows of $\mathbf{Q}'$ (corresponding to different tokens), where $m$ stands for the total number of random features used per embedding. Crucially, the sizes of those matrices/vectors do not explicitly depend on the number of graph nodes $N$. Thus one can think about pair $\mathrm{S} = ((\mathbf{K}')^\top \mathbf{V}, (\mathbf{K}')^\top \mathbf{1}_N)$ as a graph sketch. We call it a *graphot* (see: Fig. 1). The graphot stores in a compact way all learned nodes-relationships (involving both the topological signal and features in nodes) so that new embeddings of nodes can be obtained simply by interacting with that (graph-size independent) sketch rather than explicitly with other nodes.

## 5 Experiments

We tested GKAT with RWGNK kernels leading to the decomposable attention (see: Sec. 4.2). We conducted an exhaustive set of experiments ranging from purely combinatorial to bioinformatics tasks and benchmarked **10** different methods. For GKAT, we used few random walks per node. All experiments were run on a single Tesla P100 GPU with 16GB memory.

### 5.1 Combinatorial Classification

In this section we focus on the problem of detecting local patterns in graphs. A model takes a graph G as an input and decides whether it contain some graph from the given family of graphs $\mathcal{H}$ as a subgraph (not necessarily induced) or is $\mathcal{H}$-free. This benchmark tests the abilities of different methods to solve purely combinatorial tasks.
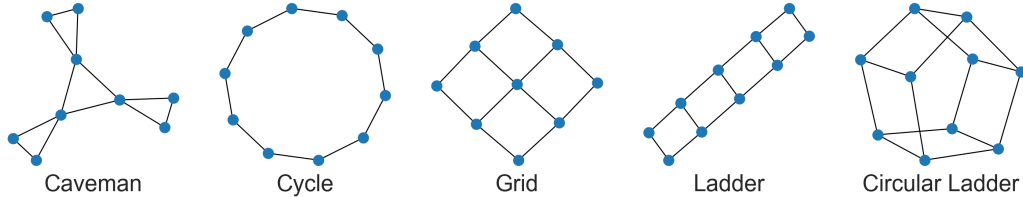


Caveman     Cycle     Grid     Ladder     Circular Ladder

Figure 3: Five motifs (patterns) used in the first class of the combinatorial classification experiments. For each pattern H, an algorithm is trained to distinguish between graphs G containing H and H-free. A naive brute-force algorithm for conducting this task would have time complexity $\Omega(N^h)$, where $h$ is the number of nodes of the motif, prohibitively expensive for all these motifs (even for small graphs G) since we have here: $h \geq 9$.

#### 5.1.1 Erdős-Rényi Random Graph with Motifs

**Data Generation:** Following procedure from [25], we used five binary classification datasets consisting of random Erdős-Rényi (ER) graphs connected with motifs (positive example) or other smaller ER graphs with the same average degree as a motif (negative example), see: Fig. 3 (details in the Appendix, Sec. 7.3). For each dataset we constructed: $S = 2048$ positive and $S$ negative examples.

**Tested Algorithms & Parameter Setting:** We tested our GKAT, graph convolution networks (GCNs)([20]), spectral graph convolution networks (SGCs)([8]) and graph attention networks

(GATs)([35]). A feature vector in each vertex was of length $l = 5$ and contained top ordered $l$ degrees of its neighbors (if there are fewer than $l$ neighbors, we padded zeroes). A dataset for each motif was randomly split into $75\%/25\%$ training/validation set. We chose: the number of epochs $E = 500$, batch size $B = 128$, used Adam optimizer with learning rate $\eta = 0.001$ and early-stopped training if neither the validation loss nor validation accuracy improved for $c = 80$ continuous epochs.

We applied 2-layer architectures. For GCNs and SGCs, we used $h = 32$ nodes in the hidden layer. For SGC, we furthermore binded each hidden layer with 2 polynomial localized filters. For GAT and GKAT, we used 2 attention heads, with $h = 9$ nodes in the hidden layer to make all models of comparable sizes. In GKAT we used random walks of length $\tau = 3$. The results are presented in Fig. 4. We see that GKAT outperforms all other methods for **all** the motifs.
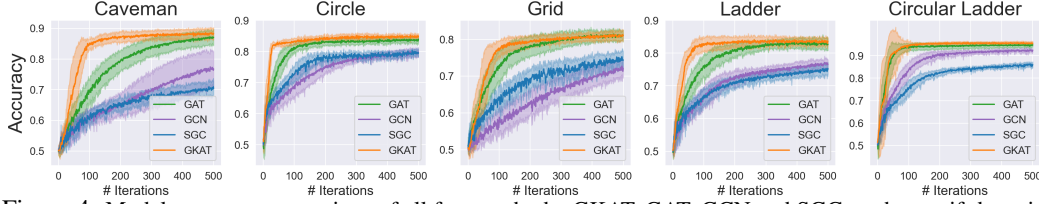


Figure 4: Model accuracy comparison of all four methods: GKAT, GAT, GCN and SGC on the motif-detection task. All tested architectures are 2-layer. GKAT outperforms other algorithms on all the tasks.

### 5.1.2 Detecting Long Induced Cycles & Deep versus Dense Attention Tests

Next we took as $\mathcal{H}$ an infinite family of motifs rather than just a single motif. The algorithm needed to decide decide whether a graph contains an induced cycle of length $> T$ for a given constant $T$. Thus the motif itself became a global property that cannot be detected by exploring just a close neighborhood of a node. In this experiment we focused also on the "depth versus density" trade-off. Shallow neural networks with dense attention from Eq. 2 are capable of modeling deeper networks relying on sparse layers, yet the price is extra computational cost per layer. The question arises whether architectures that apply RWGNK kernels leveraging efficient decomposable long-range attention from Sec. 4.2 can also replace deeper counterparts or they lose their expressiveness.
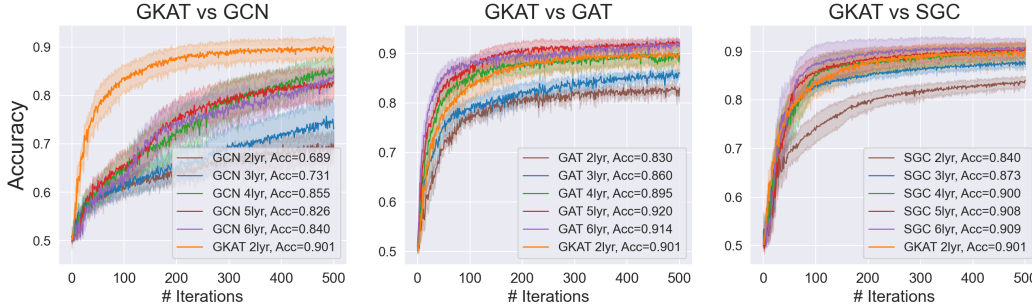


Figure 5: Comparison of the two-layer GKAT with different variants of GCNs, GATs and SGCs, varying by the number of hidden layers. Shallow GKAT architecture has expressiveness of deeper version of its counterparts and in fact outperforms many of them (e.g. graph convolution networks.)

**Dataset Generation:** We created $S = 2048$ random binary trees, each having 50 nodes, with $75\%/25\%$ for training/validation. For each tree, we constructed a positive example, by connecting two nodes with the farthest distance from each other (a negative example was obtained by connecting two random vertices, but not farthest from each other). Note that a positive example constructed in such a way has shortest induced cycle of length $P + 1$, where $P$ is the diameter of the tree.

**Tested Algorithms & Parameter Setting:** We used the same algorithms as before. This time we run detailed ablation studies on the depth of the competitors of our GKAT, by comparing two-layer GKAT with GATs, GCNs and SGCs of up to six hidden layers.

To make fair comparison, we used models with comparable number of parameters. For the two-layer GKAT, we applied 8 heads in the first layer, and 1 head in the second layer. The dimension of each head was $d = 4$. The last layer was fully-connected with output dimensionality $o = 2$ for binary classification. We applied random walk length of $\tau = 6$. For GCN, GAT and SGC, we tested number of layers ranging from 2 to 6. We controlled the number of nodes in the hidden layer(s) for GCN, GAT and SGC, and the number of attention heads in each head for GAT so that their total number of trainable parameters was comparable with this of our two-layer GKAT. All other parameters were chosen as in Sec. 5.1.1. More details on parameter settings and additional ablation tests over random

7

walk length of GKAT are given in Table 5 and Fig. 6 in the Appendix (Sec. 7.3). Our main results are presented in Fig. 5.

We see that a shallow two-layer GKAT beats all GCN-variants (also deeper ones) as well as GATs and SGCs with $< 4$ layers by a wide margin. A two-layer GKAT is asymptotically equivalent to the four-layer GAT and SGC, yet as we show in Sec. 5.3, is faster to train and run inference on.

## 5.2 Graph Neural Networks for Bioinformatics Tasks & Social Networks Data

**Datasets:** We tested GKAT for graph classification tasks on 9 standard and publicly available bioinformatics and social networks datasets [19] using a carefully designed model selection and assessment framework for a fair comparison [11]. The former include: D&D [9], PROTEINS [2], NCI1 [36] and ENZYMES [14], and the latter: IMDB-BINARY, IMDB-MULTI, REDDIT-BINARY, REDDIT-5K and COLLAB [39] (for detailed description, see Table 6 in the Appendix, Sec. 7.4.1).

Table 1: Performance of different algorithms on the bioinformatics datasets. For each dataset, we highlihgted the best performing method(s) and underlined the second best. GKAT is the best on three out of four tasks.

|  | **D&D** | **NCI1** | **Proteins** | **Enzymes** |
|---|---|---|---|---|
| **Baseline** | 78.4 ±4.5% | 69.8±2.2% | **75.8±3.7%** | 65.2±6.4% |
| **DGCNN** | 76.6±4.3% | 76.4±1.7% | 72.9±3.5% | 38.9±5.7% |
| **DiffPool** | 75.0±3.5% | **76.9±1.9%** | 73.7±3.5% | 59.5±5.6% |
| **ECC** | 72.6±4.1% | 76.2±1.4% | 72.3±3.4% | 29.5±8.2% |
| **GraphSAGE** | 72.9±2.0% | 76.0±1.8% | 73.0±4.5% | 58.2±6.0% |
| **RWNN** | 77.6±4.7% | 71.4±1.8% | 74.3±3.3% | 56.7±5.2% |
| **GKAT** | **78.6±3.4%** | 75.2±2.4% | **75.8 ±3.8%** | **69.7 ±6.0%** |

Table 2: Performance of different algorithms on the social network datasets. As for the previous table, for each dataset we marked the best performing methods. GKAT is among two top methods for four out of five tasks.

|  | **IMDB-B** | **IMDB-M** | **REDDIT-B** | **REDDIT-5K** | **COLLAB** |
|---|---|---|---|---|---|
| **Baseline** | 70.8±5.0% | **49.1 ±3.5%** | 82.2±3.0% | 52.2±1.5% | 70.2±1.5% |
| **DGCNN** | 69.2±5.0% | 45.6±3.4% | 87.8±2.5% | 49.2±1.2% | 71.2±1.9% |
| **DiffPool** | 68.4±3.3% | 45.6±3.4% | 89.1±1.6% | 53.8±1.4% | 68.9±2.0% |
| **ECC** | 67.7±2.8% | 43.5±3.1% | out of memory | out of memory | out of memory |
| **GraphSAGE** | 68.8±4.5% | 47.6±3.5% | 84.3±1.9% | 50.0±1.3% | **73.9±1.7%** |
| **RWNN** | 70.8±4.8% | 47.8±3.8% | **90.4±1.9%** | 51.7±1.5% | 71.7±2.1% |
| **GKAT** | **71.4±2.6%** | 47.5±4.5% | 89.3±2.3% | **55.3±1.6%** | 73.1±2.0% |

**Tested Algorithms:** We compared GKAT with some of the top GNN methods used previously for that data: DCGNN [43], DiffPool [40], ECC [32], GraphSAGE [13] and RWNN [25], which are selected based on their popularity and architectural differences. For bioinformatics datasets, but ENZYMES, we used the Molecular Fingerprint (MF) method [28, 23] as a baseline. It first applies global sum pooling and then a single-layer MLP with ReLU activations. For social datasets and ENZYMES, we applied the DeepMultisets (DM) method [41] as a baseline. It is a single-layer MLP followed by global sum pooling and another single-layer MLP for classification. These two baselines use no topological graph signal, so are good to quantify the benefits coming from GNNs.

**GKAT Setting:** We applied a two-layer GKAT followed by the baseline layers: we first applied an attention layer with $k$ heads (a hyperparameter to be tuned), and then another one with one head to aggregate topological information on graphs. Next, we applied either the MF method or the DM method to further process the aggregated information. The random walk length $\tau$ in each GKAT layer satisfied: $\tau \leq 4$ and depended on the evaluated datasets. A long random walk could in principle capture more information, but at the cost of adding contributions also from non-relevant nodes. The average graph diameter (the average value of the longest shortest path for each pair of nodes and over all graphs in a dataset) shown in Table 6 in the Appendix helps to calibrate walk length. We chose it to balance the pros of using a shallow architecture and the cons of information loss from such dense layer compression. Our two-layer GKAT with multi-heads and small number of trainable parameters in each head increased the number of the baseline's parameters only by a negligible fraction.

**Training Details:** We used a 10-fold CV for model assessment, and an inner holdout with $90\%/10\%$ training/validation split for model selection following the same settings [11]. We then trained the

whole training-fold three times, randomly holding out $10\%$ of the data for early stopping after model selection in each fold. The average score of these three runs was reported in Table 1 and Table 2.

The results from Table 1 and Table 2 show that GKAT is the best on three out of four bioinformatics datasets and is among two best methods on four out of five social network datasets. Furthermore, it's the only GNN method that consistently outperforms baseline on all but one bioinformatics dataset (bio-data benefits more than others from efficient longer-range attention modeling as showed in [6]). In the Appendix (Sec. 7.5) we included additional comparisons of GKAT with GAT on citation networks, where GKAT outperforms GAT on two out of three tested datasets.

### 5.3    Space & Time Complexity Gains of GKAT

Finally, we measured speed and memory improvements coming from GKAT with decomposable attention mechanism turned on (called here shortly: GKAT+) as compared to GAT as well as accuracy loss in comparison to GKAT with decomposability not leveraged. The results are presented in Table 3. We decided to report relative rather than absolute numbers since the former are transferable across different computational setups. We see that the accuracy gaps of the corresponding GKAT and GKAT+ models (obtained after the same number of epochs) are marginal, yet GKAT+ yields consistent speed and memory gains as compared to GAT per attention layer, particularly substantial for very large graphs as those from Citeseer and Pubmed.

Table 3: Speed & Space Complexity gains coming from GKAT with decomposable attention mechanism turned on (GKAT+). First row: memory compression of the graph via graphot (lower better). Second & third row: speedup in training and inference respectively per one attention layer as compared to GAT. Last row: accuracy loss as compared to GKAT not applying decomposable attention mechanism. We used four datasets from Sec. 5.1.1, a dataset from Sec. 5.1.2 (Tree) and two citation network datasets (see: Sec. 7.5): Citeseer and Pubmed with graphs of much larger sizes and on which GKAT also outperforms GAT. We applied $r$ random features to linearize softmax kernel for features in nodes with $r = 256$ for citation network datasets, $r = 16$ for datasets from Sec. 5.1.1 and $r = 8$ for a dataset from Sec. 5.1.2.

|                          | Cavem. | Circle | Grid  | Ladder | Tree  | Citeseer | Pubmed |
|--------------------------|--------|--------|-------|--------|-------|----------|--------|
| **graphot size/graph size** | 0.54   | 0.53   | 0.55  | 0.52   | 0.95  | 0.18     | 0.07   |
| **train speedup vs GAT** | 1.40x  | 1.41x  | 1.42x | 1.40x  | 1.10x | 5.10x    | 9.50x  |
| **inf speedup vs GAT**   | 1.46x  | 1.49x  | 1.49x | 1.47x  | 1.12x | 5.21x    | 9.54x  |
| **GKAT - GKAT+ (accur.)** | 0.07%  | 0.09%  | 0.08% | 0.07%  | 0.06% | 0.05%    | 0.06%  |

Below we show that regular GKAT is also faster that its counterparts (GCN, GAT and SGC) in terms of wall clock time needed to reach particular accuracy levels. We illustrate it by comparing accuracy levels reached by different models in a given wall clock time budget (the clock time that GKAT needs to complete first 100 epochs). The results are presented in Table 4.

Table 4: Running time of training different networks on datasets from Sec. 5.1.1 and Sec. 5.1.2. For GCN, GAT and SGC, we reported the accuracy with 2 layers. For GKAT, we used a 2-layer architecture and reported the accuracy with a fixed walk length of 6 for Induced Cycle Detection, and of 3 for five motifs from Sec. 5.1.1.

|      | Induced Cycle | Caveman | Circle | Grid  | Ladder | Circle Ladder |
|------|---------------|---------|--------|-------|--------|---------------|
| GCN  | 63.2%         | 62.1%   | 71.4%  | 59.3% | 66.7%  | 87.4%         |
| GAT  | 77.0%         | 69, 1%  | 80.6%  | 73.8% | 75.9%  | 93.7%         |
| SGC  | 56.6%         | 55.4%   | 64.7%  | 58.2% | 59.1%  | 66.5%         |
| **GKAT** | **83.6%**     | **85.1%** | **83.3%** | **77.1%** | **82.4%** | **94.6%**     |

## 6    Conclusion

We presented Graph Kernel Attention Transformers (GKATs), new attention-based graph neural networks leveraging graph kernel methods and recent advances in scalable attention to provide more expressive models for graph data that are also characterized by low time complexity and memory footprint. Furthermore, we demonstrated that they outperform other techhniques on a wide range of tasks from purely combinatorial problems, through social network data to bioinformatics challenges.

# References

[1] J. Atwood and D. Towsley. Diffusion-convolutional neural networks. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1993–2001, 2016.

[2] K. M. Borgwardt, C. S. Ong, S. Schönauer, S. V. N. Vishwanathan, A. J. Smola, and H.-P. Kriegel. Protein function prediction via graph kernels. *Bioinformatics*, 21(1):47–56, Jan. 2005.

[3] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst. Geometric deep learning: Going beyond euclidean data. *IEEE Signal Process. Mag.*, 34(4):18–42, 2017.

[4] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. In Y. Bengio and Y. LeCun, editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.

[5] X. Chen. A note on spectral graph neural network. *arXiv: Spectral Theory*, 2020.

[6] K. Choromanski, V. Likhosherstov, D. Dohan, X. Song, A. Gane, T. Sarlós, P. Hawkins, J. Davis, A. Mohiuddin, L. Kaiser, D. Belanger, L. Colwell, and A. Weller. Rethinking attention with performers. *ICLR 2021*, 2020.

[7] A. Daniely, R. Frostig, V. Gupta, and Y. Singer. Random features for compositional kernels. *CoRR*, abs/1703.07872, 2017.

[8] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *CoRR*, abs/1606.09375, 2016.

[9] P. D. Dobson and A. J. Doig. Distinguishing enzyme structures from non-enzymes without alignments. *Journal of molecular biology*, 330(4):771—783, July 2003.

[10] D. Duvenaud, D. Maclaurin, J. Aguilera-Iparraguirre, R. Gómez-Bombarelli, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 2224–2232, 2015.

[11] F. Errica, M. Podda, , D. Bacciu, and A. Micheli. A fair comparison of graph neural networks for graph classification. *ICLR 2020*, 2020.

[12] H. Gao, Z. Wang, and S. Ji. Large-scale learnable graph convolutional networks. In Y. Guo and F. Farooq, editors, *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, pages 1416–1424. ACM, 2018.

[13] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[14] S. Ida, C. Antje, E. Christian, G. Marion, H. Christian, H. Gregor, and S. Dietmar. Brenda, the enzyme database: updates and major new developments. *Nucleic Acids Research*, 32(suppl 1):D431–D433, 01 2004.

[15] A. Jambulapati and A. Sidford. Efficient $\tilde{O}(n/\epsilon)$ spectral sketches for the laplacian and its pseudoinverse. SODA '18, page 2487–2503, USA, 2018. Society for Industrial and Applied Mathematics.

[16] U. Kang, H. Tong, and J. Sun. Fast random walk graph kernel. In *Proceedings of the Twelfth SIAM International Conference on Data Mining, Anaheim, California, USA, April 26-28, 2012*, pages 828–838. SIAM / Omnipress, 2012.

[17] J. Kasai, H. Peng, Y. Zhang, D. Yogatama, G. Ilharco, N. Pappas, Y. Mao, W. Chen, and N. A. Smith. Finetuning pretrained transformers into rnns. *CoRR*, abs/2103.13076, 2021.

[18] A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on*

*Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5156–5165. PMLR, 2020.

[19] K. Kersting, N. M. Kriege, C. Morris, P. Mutzel, and M. Neumann. Benchmark data sets for graph kernels, 2016. *URL https://ls11-www.cs.tu-dortmund.de/staff/morris/graphkerneldatasets*, 795, 2016.

[20] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.

[21] D. Krotov and J. J. Hopfield. Dense associative memory for pattern recognition. In D. D. Lee, M. Sugiyama, U. von Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1172–1180, 2016.

[22] V. Likhosherstov, K. Choromanski, J. Davis, X. Song, and A. Weller. Sub-linear memory: How to make performers slim. *CoRR*, abs/2012.11346, 2020.

[23] E. Luzhnica, B. Day, and P. Liò. On graph classification networks, datasets and baselines. *arXiv preprint arXiv:1905.04682*, 2019.

[24] J. D. Marble and K. E. Bekris. Computing spanners of asymptotically optimal probabilistic roadmaps. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2011, San Francisco, CA, USA, September 25-30, 2011*, pages 4292–4298. IEEE, 2011.

[25] G. Nikolentzos and M. Vazirgiannis. Random walk graph neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16211–16222. Curran Associates, Inc., 2020.

[26] N. Parmar, A. Vaswani, J. Uszkoreit, L. Kaiser, N. Shazeer, A. Ku, and D. Tran. Image transformer. In J. G. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 4052–4061. PMLR, 2018.

[27] H. Peng, N. Pappas, D. Yogatama, R. Schwartz, N. A. Smith, and L. Kong. Random feature attention. *CoRR*, abs/2103.02143, 2021.

[28] L. Ralaivola, S. J. Swamidass, H. Saigo, and P. Baldi. Graph kernels for chemical informatics. *Neural networks*, 18(8):1093–1110, 2005.

[29] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Gallagher, and T. Eliassi-Rad. Collective classification in network data.

[30] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson. Structured sequence modeling with graph convolutional recurrent networks. In L. Cheng, A. C. Leung, and S. Ozawa, editors, *Neural Information Processing - 25th International Conference, ICONIP 2018, Siem Reap, Cambodia, December 13-16, 2018, Proceedings, Part I*, volume 11301 of *Lecture Notes in Computer Science*, pages 362–373. Springer, 2018.

[31] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Process. Mag.*, 30(3):83–98, 2013.

[32] M. Simonovsky and N. Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3693–3702, 2017.

[33] D. Tran-Van, A. Sperduti, and F. Costa. Link enrichment for diffusion-based graph node kernels. In A. Lintas, S. Rovetta, P. F. M. J. Verschure, and A. E. P. Villa, editors, *Artificial Neural Networks and Machine Learning - ICANN 2017 - 26th International Conference on Artificial Neural Networks, Alghero, Italy, September 11-14, 2017, Proceedings, Part II*, volume 10614 of *Lecture Notes in Computer Science*, pages 155–162. Springer, 2017.

[34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. von Luxburg, S. Bengio, H. M. Wallach, R. Fergus, S. V. N. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.

[35] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio. Graph attention networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018.

[36] N. Wale, I. A. Watson, and G. Karypis. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowl. Inf. Syst.*, 14(3):347–375, 2008.

[37] X. Wang and A. Gupta. Videos as space-time region graphs. In V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part V*, volume 11209 of *Lecture Notes in Computer Science*, pages 413–431. Springer, 2018.

[38] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu. Heterogeneous graph attention network. In L. Liu, R. W. White, A. Mantrach, F. Silvestri, J. J. McAuley, R. Baeza-Yates, and L. Zia, editors, *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, pages 2022–2032. ACM, 2019.

[39] P. Yanardag and S. Vishwanathan. Deep graph kernels. KDD '15, page 1365–1374, New York, NY, USA, 2015. Association for Computing Machinery.

[40] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec. Hierarchical graph representation learning with differentiable pooling. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

[41] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola. Deep sets. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[42] J. Zhang, X. Shi, J. Xie, H. Ma, I. King, and D. Yeung. Gaan: Gated attention networks for learning on large and spatiotemporal graphs. In A. Globerson and R. Silva, editors, *Proceedings of the Thirty-Fourth Conference on Uncertainty in Artificial Intelligence, UAI 2018, Monterey, California, USA, August 6-10, 2018*, pages 339–349. AUAI Press, 2018.

[43] M. Zhang, Z. Cui, M. Neumann, and Y. Chen. An end-to-end deep learning architecture for graph classification. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

# 7 Appendix: Graph Kernel Attention Transformers

## 7.1 Code Pointers

We include pointers to this part of the code that does not include sensitive/proprietary information. The core GKAT framework is here: https://anonymous.4open.science/r/GKAT-Experiments-BD32. We used (deterministic and random) feature map mechanisms corresponding to the features defined in graph nodes from this repository: https://github.com/google-research/google-research/tree/master/performer.

## 7.2 Graph Kernel Attention Transformers: Algorithmic box for the Attention Layer

Below we present main algorithmic box for a fixed attention layer computation in GKAT (and for a single head). Version II corresponds to the base GKAT version not leveraging decomposability of the attention and Version I to the one that uses it.

---

**Algorithm 1** Attention Layer in GKAT [computations per head]

---

**Hyperparams:** $0 < \lambda \le 1$, $\alpha \ge 0$, $0 < p < 1$ (or fixed walk length $\tau$), number of walks $t$ per node, kernel K defined on features in nodes or its corresponding kernel feature map $\phi_K$.
**Input:** graph G, $\mathbf{H} \in \mathbb{R}^{N \times d}$: matrix of rows $\mathbf{h}_i \in \mathbb{R}^d$ encoding nodes' embeddings, trainable: $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_{QK}}, \mathbf{W}_V \in \mathbb{R}^{d \times d}$.
**Output:** new matrix of embeddings $\widehat{\mathbf{H}} \in \mathbb{R}^{N \times d}$.
**Preprocessing [one-time computation per sample and all attention layers]:** Version I: For every node $h$, compute graph kernel random feature map $\Psi(h)$ corresponding to graph kernel $\mathrm{T} = \mathrm{K}_{\mathrm{RWGNK}}^{\lambda, \alpha, p}$ and defined as: $\Psi(h) = \frac{\mathbf{v}}{\|\mathbf{v}\|^\alpha}$ for $\mathbf{v} = \frac{1}{t} \sum_{i=1}^{t} f_h^{\omega(i), \lambda}$, where $\omega(1), ..., \omega(t)$ are $t$ independent random walks from $h$ constructed either with stopping probability $p$ or of fixed length $\tau$ (see: Sec. 4.3). Version II: Compute kernel matrix $\mathcal{T}(\mathrm{G}) = [\mathrm{T}(k, l)]_{k, l \in \mathrm{V}(\mathrm{G})}$.
**Main computation:**
1. Compute $\mathbf{Q} = \mathbf{H}\mathbf{W}_Q, \mathbf{K} = \mathbf{H}\mathbf{W}_K, \mathbf{V} = \mathbf{H}\mathbf{W}_V$.
2. Version I: Compute $\widehat{\mathbf{Q}}'$ and $\mathbf{K}'$ with rows defined as: $\phi_K(\mathbf{q}_l) \otimes \Psi(l)$ and $\phi_K(\mathbf{k}_l) \otimes \Psi(l)$ respectively for $l \in \mathrm{V}(\mathrm{G})$ and output: $\widehat{\mathbf{H}} = \mathrm{Diag}^{-1}(\mathbf{Q}'((\mathbf{K}')^\top \mathbf{1}_N))\mathbf{Q}'((\mathbf{K}')^T \mathbf{V})$ (see: Sec. 4.2).
2. Version II: Compute attention matrix $\mathbf{A}$ as in Eq. 2 and output: $\widehat{\mathbf{H}} = \mathbf{A}\mathbf{V}$.

---

## 7.3 Combinatorial Classification Experiments: Additional Details

The data for the motif detection task from Section 5.1.1 was generated as follows:

- Firstly we created five simple motifs as shown in Fig. 3. Note that each motif has $\ge 9$ vertices so a brute-force combinatorial algorithm for motif-detection would take time $\Omega(N^9)$, prohibitively expensive even for small graphs G.

- We then generated for each motif $S$ small Erdős-Rényi graphs with the same number of nodes as that motif and the same average degree.

- For each motif, we also generated $S$ larger Erdős-Rényi random graphs, each of 100 vertices, again of the same average degree.

- We obtained positive/negative samples by connecting each larger Erdős-Rényi random graph with the motif/previously generated smaller Erdős-Rényi random graph (with certain edge probability).

In Table 5 we present additional details regarding architectures used in the experiments from Section 5.1.2, in particular the number of parameters and heads / polynomial filters used in different layers. Ablation tests over GKAT random walk length for Section 5.1.2 are presented in Fig. 6.

## 7.4 GNNs for Bioinformatics Tasks & Social Networks Data: Additional Details

### 7.4.1 Datasets Descriptions

Detailed profiles of the datasets used in the experiments from Sec. 5.2 are given in Table 6.
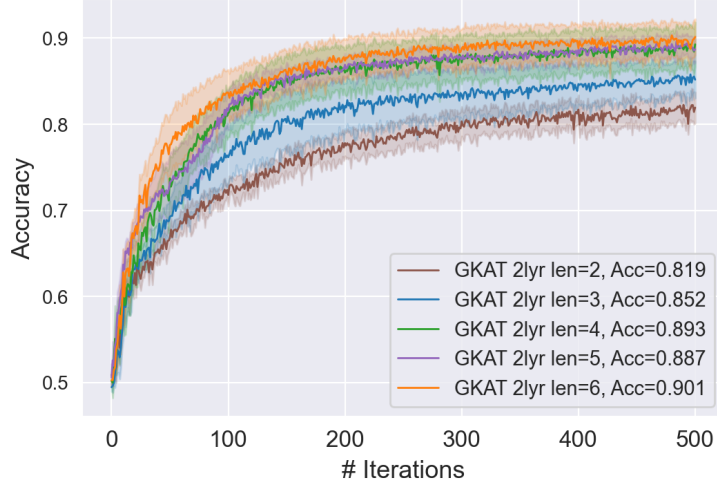
Figure 6: Ablation tests over random walk length of GKAT in Sec. 5.1.2.

Table 5: Additional details regarding architectures used in Section 5.1.2. For GKAT, we applied 8 heads in the first layer, and 1 head in the second layer, with 4 hidden units in each attention head. The total number of trainable parameters was 242. For GAT, we tested the number of layers from 2 to 6, changing the number of attention heads in each layer, but with the same number of hidden units in each attention head. For GCN, we modified the number of hidden units in each layer. For SGC, we modified the number of polynomial filters and the number of hidden units in each layer. The number of attention heads in GAT, as well as the number of hidden units in each layer in GCN and SGC were chosen to make their total number of trainable parameters comparable with the corresponding number of GKAT.

|  |  | #Heads | Dim. Head | #Parameters |
|---|---|---|---|---|
| **GKAT** | **2 layers** | $[8, 1]$ | 4 | 242 |
| **GAT** | **2 layers** | $[8, 1]$ | 4 | 242 |
|  | **3 layers** | $[4, 2, 1]$ | 4 | 242 |
|  | **4 layers** | $[4, 2, 1, 1]$ | 4 | 266 |
|  | **5 layers** | $[2, 2, 2, 1, 1]$ | 4 | 258 |
|  | **6 layers** | $[3, 2, 1, 1, 1, 1]$ | 4 | 270 |
|  |  |  | **Dim. Layer** | **#Parameters** |
| **GCN** | **2 layers** |  | $[14, 14]$ | 268 |
|  | **3 layers** |  | $[10, 10, 10]$ | 262 |
|  | **4 layers** |  | $[8, 8, 8, 8]$ | 250 |
|  | **5 layers** |  | $[10, 8, 6, 6, 6]$ | 260 |
|  | **6 layers** |  | $[10, 6, 6, 6, 6, 6]$ | 268 |
|  |  | **#Polynomial Filters** | **Dim. Layer** | **#Parameters** |
| **SGC** | **2 layers** | $[4, 2]$ | $[10, 8]$ | 236 |
|  | **3 layers** | $[4, 2, 2]$ | $[8, 6, 6]$ | 234 |
|  | **4 layers** | $[8, 2, 2, 2]$ | $[8, 5, 4, 4]$ | 247 |
|  | **5 layers** | $[8, 2, 2, 2, 2]$ | $[6, 5, 4, 4, 4]$ | 245 |
|  | **6 layers** | $[8, 2, 2, 2, 2, 2]$ | $[6, 4, 4, 4, 4, 3]$ | 249 |

For each dataset, we chose graphs with the number of nodes close to the average number of nodes shown in Table 6. Examples of bioinformatics-graphs from these datasets are given in Fig. 7. Examples of social network graphs from these datasets are given in Fig. 8.

### 7.4.2 Hyperparameter Selection

In this section, we present details regarding hyperparameter selection in Section 5.2 (see: Table 7).

The tunable parameters included: general parameters like batch size, learning rate, dropout ratio, global pooling methods, regularization rate, data normalization methods, as well as parameters specific to our GKAT layers, which included: number of GKAT layers, number of attention heads

Table 6: Bioinformatics and Social Dataset descriptions. #NODES, #EDGES and Diameter columns contain values averaged over all graphs in a given dataset.

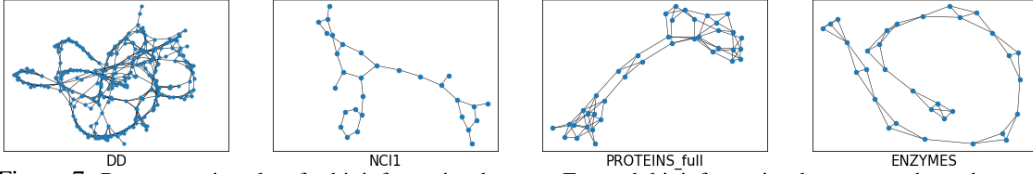| | | #Graphs | #Classes | #Nodes | #Edges | Diameter | #Features |
|---|---|---|---|---|---|---|---|
| BIOINF. | D&D | 1178 | 2 | 284.32 | 715.66 | 19.90 | 89 |
| | ENZYMES | 600 | 6 | 32.63 | 64.14 | 10.86 | 3 |
| | NCI1 | 4110 | 2 | 29.87 | 32.30 | 13.26 | 37 |
| | PROTEINS | 1113 | 2 | 39.06 | 72.82 | 11.48 | 3 |
| SOCIAL | COLLAB | 5000 | 3 | 74.49 | 2457.78 | 1.86 | 1 |
| | IMDB-BINARY | 1000 | 2 | 19.77 | 96.53 | 1.86 | 1 |
| | IMDB-MULTI | 1500 | 3 | 13.00 | 65.94 | 1.47 | 1 |
| | REDDIT-BINARY | 2000 | 2 | 429.63 | 497.75 | 9.72 | 1 |
| | REDDIT-5K | 4999 | 5 | 508.82 | 594.87 | 11.96 | 1 |



Figure 7: Representative plots for bioinformatics datasets. For each bioinformatics dataset, we chose the graph with number of nodes most similar to the average number of nodes shown in Table 6.

and dimension of each head in a GKAT layer. We also tuned other options: whether to add a fully-connected layer after data normalization, but before GKAT layers, and dimension of fully-connected layers (both preceding and coming after GKAT layers). Due to the large amount of tunable parameters, we decided to first conduct a rough search for each parameter using only one random CV fold, select one/several parameter combination(s) with best performance, and then reused on all other folds.

For all other methods, we reported the best scores conducted via an extensive hyperparameters grid search [11]. For GKAT, we fixed the number of epochs to $E = 1000$, early stopping patience as $500$ epochs, the criterion for early stopping as validation accuracy, global pooling method as summation, and used Adam optimizer. Then we performed hyperparameter tuning for: batch size $B \in \{32, 128\}$, learning rate $\eta \in \{0.01, 0.001, 0.0001\}$, dropout ratio $\in \{0.0, 0.1, 0.2, 0.4\}$, $L_2$-regularization rate $\in \{0.001, 0.005\}$, dimension of attention head in the first GKAT layer $h \in \{4, 8, 16, 32\}$, number of attention heads in the first GKAT layer $\in \{1, 4, 8, 12\}$, number of nodes in the MLP layer $\in \{32, 64, 128\}$, GKAT random walk length $\tau \in \{1, 2, 3, 4\}$, whether to use a fully-connected layer before the first GKAT layer, and whether to apply batch normalization to pre-prosess data before feeding the model with it.

For some of the datasets (e.g. D&D), we selected the best hyperparameter set optimized over one random CV fold, and used it across all cross-validation outer folds.

Table 7: Hyperparameter settings for the bioinformatics and social network datasets from Section 5.2.

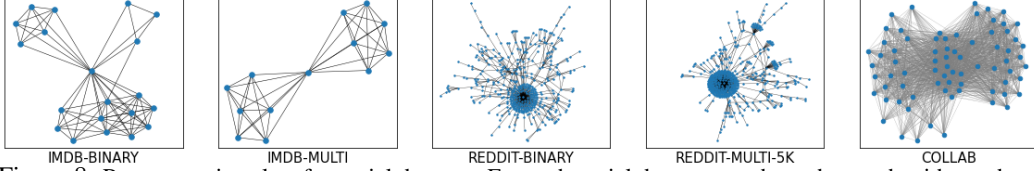| | | BS | #Heads | $d_{Head}$ | $d_{FC}$ | $Len_{rw}$ | Drop | L2 | add FC | Norm |
|---|---|---|---|---|---|---|---|---|---|---|
| CHEM. | D&D | 32 | 8 | 16 | 128 | 2 | − | 0.005 | No | BN |
| | NCI1 | 32 | 8 | 32 | 64 | 4 | 0.1 | 0.001 | Yes | BN |
| | PROTEINS | 32 / 128 | 4 / 8 | 8 / 32 | 32 / 128 | 3 | − | 0.001 | Yes | BN |
| | ENZYMES | 32 | 4 / 8 | 16 / 32 | 32 / 64 | 3 | 0.1 | 0.001 | Yes | BN |
| SOCIAL | COLLAB | 32 | 12 | 4 | 128 | 2 | − | 0.005 | No | No |
| | IMDB-BINARY | 32 | 8 / 12 | 4 / 8 | 64 / 128 | 1 | − | 0.005 | No | No / BN |
| | IMDB-MULTI | 32 | 8 / 12 | 4 / 8 | 64 / 128 | 1 | − | 0.005 | No | No / BN |
| | REDDIT-BINARY | 32 | 4 | 4 | 128 | 2 | − | 0.005 | No | BN |
| | REDDIT-5K | 32 | 8 | 8 | 64 | 2 | − | 0.005 | No | BN |

15

Figure 8: Representative plots for social datasets. For each social dataset, we chose the graph with number of nodes most similar to the average number of nodes shown in Table 6.

## 7.5 Experiments with Citation Networks Datasets

### 7.5.1 Datasets Descriptions

**Datasets:** To directly compare GKAT with GAT, we also tested both algorithms on three publicly available citation networks datasets: Cora, Citeseer and Pubmed ([29]) with the same data splits as in [35]. Datasets descriptions are given in Table 8.

Table 8: Citation Networks Datasets Descriptions.

|  | Cora | Citeseer | Pubmed |
|---|---|---|---|
| **#Nodes** | 2708 | 3327 | 19717 |
| **#Edges** | 5419 | 4732 | 44338 |
| **#Features** | 1433 | 3703 | 500 |
| **#Classes** | 7 | 6 | 3 |
| **#Training Nodes** | 140 | 120 | 60 |
| **#Validation Nodes** | 500 | 500 | 500 |
| **#Test Nodes** | 1000 | 1000 | 1000 |

### 7.5.2 Comparison with GAT

**Experiment Settings:** We used the same model architecture and parameters as in GAT for our GKAT to make the comparison as accurate as possible. The only difference is that we replaced the adjacency matrix masking in GAT by the normalized dot-product based similarity matrix generated from random walks, as described in Section 4.3. Both models used two-layer attention, with 8 attention heads in the first layer, and 1 head in the second layer. We used 8 hidden units in the first layer, and the number of output units in the second layer was the same as number of classes. Each layer was followed by an exponential linear unit (ELU) activation. We applied $L_2$-regularization with $\lambda = 0.0005$, dropout with $p = 0.6$ for inputs and normalized attention coefficients in both layers for all three datasets.

**Results:** The results are shown in Table 9. Our GKAT algorithm achieved lower accuracy on Cora dataset, but higher on the remaining two.

Table 9: Comparison of GAT and GKAT on citation networks datasets. For Cora and Citeseer, we reported the results for GAT from [35]. For GAT and Pubmed dataset, we reported the results averaged over 15 runs with the same parameter settings as in Cora and Citeseer. GKAT was run 15 times over multiple random walk lengths up to 7, and the best was reported.

|  | Cora | Citeseer | Pubmed |
|---|---|---|---|
| **GAT** | **83.0±0.7%** | $72.5 \pm 0.7\%$ | $77.2 \pm 0.6\%$ |
| **GKAT** | $82.1 \pm 0.7\%$ | **73.0±0.7%** | **78.0 ±0.7%** |

**Dynamic Generator of Random Walks:** We also tried the so-called *dynamic*-GKAT. The dynamic variant generated random walks from scratch in each training epoch, thus requiring additional compute. However, one advantage of the dynamic version is that we could assign different transition probabilities for adjacent nodes (rather than sampling next point of the walk uniformly at random). The transition probability matrix can be an attention matrix from Eq. 2 in Section. 4.1. An intuition behind that particular variant is that we assign higher transition probabilities for neighbors with higher

16

attention coefficients. The dynamic variant enabled us to improve accuracy of GKAT on Citeseer to **73.3**% (with reduced **0.6**% standard deviation).

### 7.5.3 Ablation Tests on Random Walk Length for GKAT

Figure 9 compares the effect of random walk path length of GKAT algorithms on training for Cora, Citeseer and Pubmed datasets. We run GKAT with multiple random walk lengths up to 7. The results show that a small path length no longer than 4 is enough for GKAT and dynamic-GKAT, which supports our claim that short walks are sufficient for GKAT.
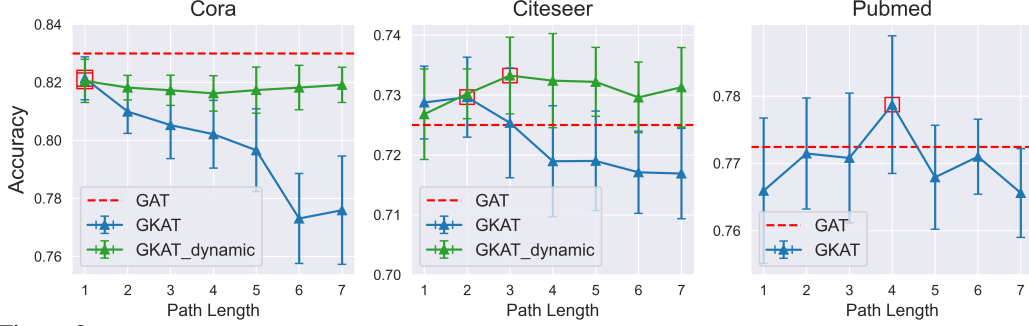


Figure 9: Ablation tests over random walk path lengths of GKAT and dynamic-GKAT on Cora, Citeseer and Pubmed datasets. The errorbar represents 1 standard deviation. Test-accuracies for the optimized GAT were shown as horizontal red dotted lines in each subplot. Best test accuracies of GKAT algorithms were highlighted with red squares. The dynamic-GKAT was tested only on datasets, where accuracy advantage of the regular optimized GKAT over optimized GAT was $\leq 0.5\%$.

### 7.6 Proof of Theorem 1

*Proof.* Note first that since $\omega(k)$ and $\omega(l)$ are chosen independently, we have:

$$\mathrm{K}^{\lambda,0,p}_{\mathrm{RWGNK}}(k,l) = \mathbb{E}_{\omega(k)}[f^{\omega(k),\lambda}_k] \cdot (\mathbb{E}_{\omega(l)}[f^{\omega(l),\lambda}_l])^\top = \mathbb{E}_{\omega(k),\omega(l)}[f^{\omega(k),\lambda}_k (f^{\omega(l),\lambda}_l)^\top] \quad (6)$$

Denote: $X = f^{\omega(k),\lambda}_k (f^{\omega(l),\lambda}_l)^\top$. The key observation is that $X$ can be rewritten as:

$$X = \sum_{\substack{u \in V(G) \\ }} \sum_{\substack{(j_1=k,\dots,j_{a+1}=u)=\mathrm{pref}(\omega(k)), \\ (j'_1=l,\dots,j'_{b+1}=u)=\mathrm{pref}(\omega(l))}} \lambda^a \lambda^b = \sum_{(j_1=k,\dots,j_{a+b+1}=l) \in \Omega(k,l)} \lambda^{a+b} \quad (7)$$

where $\Omega(k,l)$ is the multi-set of walks from $k$ to $l$ that are built from some prefix of $\omega(k)$ concatenated with some prefix of $\omega(l)$. Therefore we can write $X$ as:

$$X = \sum_{r \in R(k,l)} \sum_{i=0}^{\mathrm{len}(r)} \lambda^{\mathrm{len}(r)} 1[\mathcal{E}(r,i)], \quad (8)$$

where $R(k,l)$ is the set of walks from $k$ to $l$, $\mathrm{len}(r)$ stands for the length (number of edges) of walk $r$ and $\mathcal{E}(r,i)$ is an event that first $i$ edges of the walk $r$ (counting from $k$) form the prefix sub-walk of $\omega(k)$ and the remaining ones form the prefix sub-walk of $\omega(l)$. Therefore we have:

$$\mathrm{K}^{\lambda,0,p}_{\mathrm{RWGNK}}(k,l) = \mathbb{E}_{\omega(k),\omega(l)}\left[ \sum_{r \in R(k,l)} \sum_{i=0}^{\mathrm{len}(r)} \lambda^{\mathrm{len}(r)} 1[\mathcal{E}(r,i)] \right] =$$

$$\sum_{r \in R(k,l)} \sum_{i=0}^{\mathrm{len}(r)} \lambda^{\mathrm{len}(r)} \mathbb{P}_{\omega(k),\omega(l)}[\mathcal{E}(r,i)] = \sum_{r \in R(k,l)} \sum_{i=0}^{\mathrm{len}(r)} \lambda^{\mathrm{len}(r)} \prod_{j=0}^{i-1} \frac{1-p}{\deg(r^j)} \prod_{t=0}^{\mathrm{len}(r)-i-1} \frac{1-p}{\deg(r^{\mathrm{len}(r)-1-t})}, \quad (9)$$

where $r^y$ stands for the $y^{th}$ vertex of the walk $r$ starting from $k$ and $\deg(v)$ denotes the degree of a vertex $v$.

17

Therefore we obtain:

$$\sum_{r \in R(k,l)} \sum_{i=0}^{\text{len}(r)} \left( \frac{(1-p)\lambda}{d_{\max}(\text{G})} \right)^{\text{len}(r)} \le \text{K}_{\text{RWGNK}}^{\lambda,0,p}(k,l) \le \sum_{r \in R(k,l)} \sum_{i=0}^{\text{len}(r)} \left( \frac{(1-p)\lambda}{d_{\min}(\text{G})} \right)^{\text{len}(r)} \qquad (10)$$

We conclude that:

$$\sum_{i=0}^{\infty} r_{k,l}(i) \left( \frac{(1-p)\lambda}{d_{\max}(\text{G})} \right)^{i} (i+1) \le \text{K}_{\text{RWGNK}}^{\lambda,0,p}(k,l) \le \sum_{i=0}^{\infty} r_{k,l}(i) \left( \frac{(1-p)\lambda}{d_{\min}(\text{G})} \right)^{i} (i+1) \qquad (11)$$

To complete the proof, it suffices to notice that matrix $\text{Adj}^i(\text{G})$ encodes the number of walks of length $i$ between pairs of vertices in G. $\qquad \square$