

Fast Attributed Graph Embedding via Density of States

Saurabh Sawlani
Carnegie Mellon University
saurabh.sawlani@gmail.com

Lingxiao Zhao
Carnegie Mellon University
lingxia1@andrew.cmu.edu

Leman Akoglu
Carnegie Mellon University
lakoglu@andrew.cmu.edu

Abstract—Given a node-attributed graph, how can we efficiently represent it with few numerical features that expressively reflect its topology and attribute information? We propose A-DOGE, for Attributed DOS-based Graph Embedding, based on density of states (DOS, a.k.a. spectral density) to tackle this problem. A-DOGE is designed to fulfill a long desiderata of desirable characteristics. Most notably, it capitalizes on efficient approximation algorithms for DOS, that we extend to blend in node labels and attributes for the first time, making it fast and scalable for large attributed graphs and graph databases. Being based on the entire eigenspectrum of a graph, A-DOGE can capture structural and attribute properties at multiple (“glocal”) scales. Moreover, it is unsupervised (i.e. agnostic to any specific objective) and lends itself to various interpretations, which makes it suitable for exploratory graph mining tasks. Finally, it processes each graph independent of others, making it amenable for streaming settings as well as parallelization. Through extensive experiments, we show the efficacy and efficiency of A-DOGE on exploratory graph analysis and graph classification tasks, where it significantly outperforms unsupervised baselines and achieves competitive performance with modern supervised GNNs, while achieving the best trade-off between accuracy and runtime.

Index Terms—attributed graphs, spectral embedding, graph filters, band-pass, density of states

I. INTRODUCTION

Graphs are widely used to model structured data from different domains such as chemistry [1], biology [2], cybersecurity [3], finance [4], etc. The effectiveness and popularity of data-driven machine learning algorithms has necessitated expressive vector representations of different kinds of complex data, and graphs are no exception. Different from images or text, graphs pose novel challenges in finding effective representations as graph databases may contain graphs that vary in size and structure, and do not necessarily exhibit alignment (i.e. correspondence) between the nodes of different graphs.

Formally, we want to design a function $R : G \mapsto \mathbf{z}_G \in \mathbb{R}^D$, where D is a fixed embedding size that does not depend on the input graph size. Ideally, given a graph database with N graphs (with n nodes and m edges per graph on average), we want R to be (i) *permutation and size invariant*; where graphs with similar structure and label/attribute distribution have similar embeddings irrespective of node ordering and number of nodes, (ii) *flexible*; that leverages information from node labels and/or multi-attributes as well as edge weights, (iii) *multi-scale/glocal*; that can capture local/microscopic, mesoscopic, as well as global/macroscale properties of a graph, and (iv) *task-agnostic/unsupervised*; that can produce embeddings independent of any downstream task or related

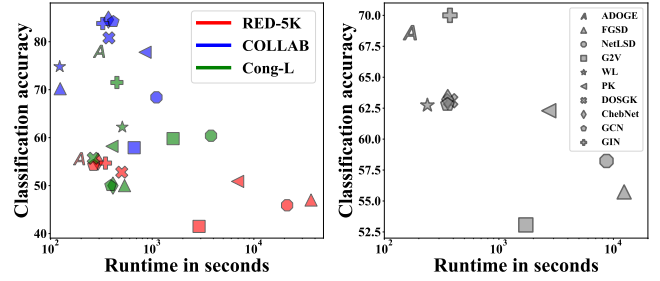


Figure 1. A-DOGE achieves superior runtime-performance trade-off. (left) Runtime (log-scale) vs. accuracy on three individual datasets (with largest N , $\max(m)$ and the largest synthetic dataset). (right) Average runtime vs. graph classification accuracy across datasets.

class labels, where not being tied to a specific task allows embeddings to be general-purpose for use e.g. in graph mining and exploratory data analysis. In addition, as with any algorithm, we want R to be (v) *efficient* and *scalable* to large graphs (large n , m) as well as large databases (large N). Finally, R that can produce one embedding at a time (vi) *independently per graph* (as opposed to “collective processing”) may be desirable, which allows on-the-fly embedding per incoming graph in streaming settings, as well as embarrassing parallelization for speed.

Spectrally-designed embeddings are a popular class of techniques based on the graph eigenspectrum [5], as it captures key structural graph properties, such as cuts [6], random walk stationarity [7], dynamical processes and epidemic thresholds [8], diameter, connectedness, clustering, etc. [9]. However, the high complexity of computing the eigenspectrum exactly has proven to be a barrier for creating spectrum-based graph embeddings. Moreover, while the eigenspectrum can capture important topological properties, blending in node attributes/labels into spectrally-designed embeddings is nontrivial.

In this paper, we leverage fast algorithms for approximating the spectral density of a graph [10], and use it to independently construct unsupervised graph embeddings that are permutation and size invariant, flexible and multi-scale. Here, the focus is on representing the entire spectrum of the graph, which helps capture any arbitrary “band” of eigenvalues (band-pass), rather than only the extremal eigenpairs (low/high pass).

A. Prior Work

Table I gives a comparison with three categories of relevant prior work in the context of desired properties for a graph embedding. These existing work do not satisfy one or more of the aforementioned properties (ii)–(vi) as we discuss next.

Unsupervised explicit graph embedding (UEGE): Several unsupervised methods construct an explicit vector representa-

Table 1
A-DOGE SATISFIES ALL PROPERTIES, WHILE PRIOR WORK MISS ONE OR MORE OF THE INPUT GRAPH OR EMBEDDING PROPERTIES.

	Method	Inp. graph			Embedding			
		Node labels	Node attributes	Edge weights	Band-pass	Task-agnostic	Scalable	Indpt. per graph
UEGE	FGSD [11], NetLSD [12]			✓		✓	✓	✓
	G2VEC [13]	✓				✓	✓	✓
GK	WL [14], WL-OA [15]	✓				✓		
	SAGE [16], PK [17]	✓	✓	✓		✓		
	RetGK [18]			✓		✓		
	DOSGK [19]			✓	✓	✓		
GNN	GCN [20], GIN [21]	✓	✓	✓				✓
	ChebNet [22], CALEYNet [23]	✓	✓	✓	✓		✓	✓
	A-DOGE [this paper]	✓	✓	✓	✓	✓	✓	✓

tion for each graph. Among those, spectrum-based methods have gained popularity in recent years. FGSD [11] treats a graph as a collection of spectral distances between its vertices. NetLSD [12] represents a graph as a collection of heat traces of the graph at several time points. Both methods are effective at capturing local and global structural properties of a graph, however, they ignore node labels and attributes. graph2vec [13] creates Weisfeiler-Lehman (WL) subtree based features, and learns an embedding of the graph trained to predict the existence of subtrees in the graph. It admits node labels, but ignores node attributes as well as edge weights.

Graph kernels (GK): Due to the existence of many effective distance measures between graphs, graph kernels are a more widely studied method of graph representations [24]. While most popular kernels are effective at capturing characteristics of the graph structure, only a few, including the Propagation Kernel (PK) [17] are able to factor in edge weights, node labels and continuous node attributes (see Table 1 in [24]).

Several graph kernels which use spectral properties have been developed in recent years. RetGK [18] represents each graph as a collection of node embeddings, where the node features are the return-probabilities of random walks of varying lengths. SAGE [16] extends this idea to graphs with labeled and attributed nodes by appending each node embedding with its one-hot encoded label and/or attributes. However, both these methods do not scale well for large graphs. Moreover, computing return probabilities of random walks tends to over-represent local features near a node, and often fails to capture global properties of the graph [19]. These issues are addressed by the Density Of States (DOS) GK, and its point-wise (i.e. node-level) extension (PDOS)¹, which uses Chebyshev polynomials to efficiently capture global properties of random walks, and uses fast approximation techniques [10]. However, despite their efficiency, they are limited to plain graphs and do not admit node labels or attributes.

Moreover, although graph kernels have proven effective at

¹This is called LDOS in their paper, but we use PDOS to avoid confusion with our definition of LDOS in this paper.

modelling graph structure, and in some cases node labels and attributes, for many kernel methods, computing an $N \times N$ -sized kernel matrix can be restrictive in terms of both time and space, which do not scale to large databases with many graphs.

Graph Neural Networks (GNNs): While existing unsupervised embedding and kernel methods are ill-equipped to handle continuous node attributes, GNNs are able to leverage such data to a great extent. However, deep parameterized models come with their own drawbacks. They are resource-hungry, not task-agnostic, and can be slow to train. Moreover, when viewed through a spectral lens [25], most neighborhood-aggregation based GNNs such as GCN [20] and GIN [26] can only act as low-pass or high-pass filters on a graph spectrum. Only spectrally designed GNNs such as ChebNet [22] and CaleyNet [23] can act as band-pass filters.

A perhaps subtle characteristic of graph embedding methods is *independent* versus dependent/collective processing of the graphs. By design, all GNN-based methods including graph2vec require collective processing due to end-to-end training. WL and PK respectively obtain the compressed labels and histogram bins based on all graphs which makes them dependent. RetGK, DOSGK, and SAGE obtain graph-level embeddings through kernelizing the set of node-level embeddings, which is of different sizes across graphs, and hence they are inherently bound to create $N \times N$ pairwise kernel values rather than an explicit/independent embedding for each graph.

B. Our Contributions

We propose A-DOGE (Attributed DOS-based Graph Embedding), for extremely fast unsupervised embedding of attributed graphs that is permutation and size invariant, flexible, and multi-scale, which is produced independently per graph.

Our main technical contributions are as follows:

- **New graph-level embedding algorithm:** We introduce a new spectrally-designed graph embedding approach, called A-DOGE, that leverages the whole (eigen)spectrum of a graph. A-DOGE capitalizes on recent algorithms that can efficiently approximate the (local) density of states (L)DOS [10], extending to attributed graphs for the first time.
- **Desired characteristics:** Thanks to efficient approximations, A-DOGE is extremely fast. It can handle node labels, continuous multi-attributes, and edge weights. Leveraging the whole spectrum, it enables variable band-pass filtering as well as features that capture multi-scale properties. Further, it processes each graph independently of others, which makes it amenable for streaming scenarios as well as parallelization.
- **Exploratory graph analysis:** A-DOGE is not tied to any specific objective, which makes it suitable for both un/supervised tasks. In fact, our embedding features lend themselves to various interpretations, related to graph signal convolution, random walks, and band-filters, which prove useful in data mining and exploratory analysis of real-world graph datasets as we show through experiments.
- **Efficacy and Efficiency:** Extensive experiments show that A-DOGE is on par with or superior to all unsupervised baselines, and competitive against modern supervised GNNs

on graph classification tasks. Notably, it achieves the best runtime-accuracy trade-off. (See Fig. 1.)

Reproducibility: We share all datasets and source code at <https://github.com/sawhani/A-DOGE>.

II. PROBLEM STATEMENT & PRELIMINARIES

Notation. We denote scalars, vectors, matrices and sets by lowercase (x), lowercase boldface (\mathbf{x}), uppercase boldface (\mathbf{X}), and calligraphic (\mathcal{X}) letters, respectively. $\mathbf{X}_{:,j}$ and \mathbf{X}_{ij} refer to the j -th column and the (i, j) -th entry of a matrix.

We consider undirected, weighted node-attributed graphs $G = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathcal{A})$ where $\mathcal{V} = \{v_1, \dots, v_n\}$ denotes the set of n nodes, and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ denotes the set of m edges. \mathbf{W} depicts the weighted adjacency matrix where $\mathbf{W}_{ij} > 0$ if $(v_i, v_j) \in \mathcal{E}$, and 0 otherwise. \mathbf{X} is the $n \times d$ node-attribute matrix, where $\mathcal{A} = \{a_1, \dots, a_d\}$ denotes the set of d attributes, with $\text{dom}(a_j)$ depicting the domain of attribute a_j . In terms of Graph Signal Processing (GSP) terminology, any $\mathbf{x} = \mathbf{X}_{:,j}$ can be thought as a *graph signal* on the nodes, with one scalar per node.

Problem 1 (Unsupervised Graph-level Embedding). *Given a set of undirected, weighted and node-attributed/labeled graphs $\mathcal{G} = \{G_1, \dots, G_N\}$, for $G_i = (\mathcal{V}_i, \mathcal{E}_i, \mathbf{X}_i, \mathcal{A})$, where*

- (i) *graphs in \mathcal{G} can be of varying sizes, (ii) there exists no particular correspondence between the nodes of different graphs, and (iii) the (categorical and/or continuous) attributes and their domain are shared among all graphs,*

Find *D -dimensional graph-level embedding $\mathbf{z}_G \in \mathbb{R}^D$ for each $G \in \mathcal{G}$ that captures both structural and attribute information.*

Let $\widetilde{\mathbf{W}} = \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$ denote the symmetrically normalized adjacency matrix, where \mathbf{D} is the diagonal degree matrix with $\mathbf{D}_{i,i} = \sum_j \mathbf{W}_{i,j}$. Let $\widetilde{\mathbf{L}} = \mathbf{I} - \widetilde{\mathbf{W}}$ denote the Laplacian matrix, and $\mathbf{P} = \mathbf{D}^{-1} \mathbf{W}$ the random walk matrix. For a connected graph, $\widetilde{\mathbf{W}}$ has eigenvalues $-1 = \lambda_0 < \lambda_1 \leq \dots \leq \lambda_{n-1} = 1$ with corresponding eigenvectors $\{\mathbf{u}_k\}_{k=0}^{n-1}$. $\widetilde{\mathbf{W}}$ has the same set of eigenvectors as $\widetilde{\mathbf{L}}$ whose eigenvalues are the shifted set $\{\mu_k = 1 - \lambda_k\}_{k=0}^{n-1} \in [0, 2]$. $\widetilde{\mathbf{W}}$ also shares the same eigenvalues as \mathbf{P} . As such, the spectral density function has bounded support for these graph matrices. Following GSP convention, we refer to the eigenvalues as the *graph frequencies*.

In this work we use $\widetilde{\mathbf{W}}$ as the so-called graph shift operator \mathbf{S} which generalizes to any symmetric matrix of a graph. Let $\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ depict the eigendecomposition, where $\mathbf{\Lambda} := \text{diag}([\lambda_1 \dots \lambda_n])$ and $\mathbf{U} = [\mathbf{u}_1 \dots \mathbf{u}_n]$.

Definition 1 (Graph spectrum). *The **spectrum** of a graph is composed of the set of graph eigenvalues together with their multiplicities of the (normalized) adjacency matrix.*

Graph Fourier transform. The graph Fourier transform (GFT) of a graph signal $\mathbf{x} \in \mathbb{R}^n$ is defined as the projection

$$\widehat{\mathbf{y}} = \mathcal{F}(\mathbf{x}) = \mathbf{U}^T \mathbf{x}$$

and the inverse GFT of $\widehat{\mathbf{y}} \in \mathbb{R}^n$ is given as

$$\mathbf{x} = \mathcal{F}^{-1}(\widehat{\mathbf{y}}) = \mathbf{U} \widehat{\mathbf{y}}$$

Graph filtering. A graph filter is an operation on a graph signal with output in the graph frequency domain, that is,

$$\widehat{\mathbf{y}}_{flt} = \phi(\mathbf{\Lambda}) \widehat{\mathbf{y}}, \quad (1)$$

where $\phi(\mathbf{\Lambda})$ is a diagonal matrix with filter frequency response values as its diagonal elements.

Definition 2 (Frequency Response Function (FRF)). *The **frequency response function** of a graph filter is written as*

$$\phi : \mathbb{C} \mapsto \mathbb{R}, \quad \lambda_i \rightarrow \phi(\lambda_i), \quad (2)$$

which, simply put, assigns a scalar value $\phi(\lambda_i)$ to each graph frequency (i.e., eigenvalue) λ_i .

By applying the inverse GFT on both sides of Eq. (1), we can get the filter output in the node domain as

$$\mathbf{x}_{flt} = \mathbf{U} \phi(\mathbf{\Lambda}) \widehat{\mathbf{y}} = \mathbf{U} \phi(\mathbf{\Lambda}) \mathbf{U}^T \mathbf{x} = \phi(\mathbf{S}) \mathbf{x}.$$

Signal convolution. Graph convolution of two signals, say \mathbf{x} and \mathbf{x}' , each in \mathbb{R}^n , yields another signal $\mathbf{c} \in \mathbb{R}^n$ as

$$\mathbf{c}_{\mathbf{x}, \mathbf{x}'} = \mathbf{x} *_G \mathbf{x}' = \mathbf{U} (\mathbf{U}^T \mathbf{x} \odot \mathbf{U}^T \mathbf{x}') = \sum_{i=1}^n \mathbf{u}_i (\mathbf{u}_i^T \mathbf{x}) (\mathbf{u}_i^T \mathbf{x}')$$

where \odot depicts the Hadamard product. We can write the Fourier transform of the convolution as

$$\mathcal{F}(\mathbf{c}_{\mathbf{x}, \mathbf{x}'}) = \widehat{\mathbf{c}}_{\mathbf{x}, \mathbf{x}'} = \{(\mathbf{u}_i^T \mathbf{x})(\mathbf{u}_i^T \mathbf{x}')\}_{i=1}^n. \quad (3)$$

Density of States. Spectral density is the overall distribution of the eigenvalues as induced by any symmetric $n \times n$ graph matrix $\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$. It is also referred to as the density of states (DOS) in the physics literature, reflecting the number of states at different energy levels [27]. Formally,

Definition 3 (Density of States (DOS)). *DOS or the spectral density induced by \mathbf{S} is the density function*

$$f(\lambda) = \frac{1}{n} \sum_{i=1}^n \delta(\lambda - \lambda_i), \quad (4)$$

where $\delta(\cdot)$ is the Dirac delta function.

Definition 4 (Local Density of States (LDOS)). *Likewise, for any input vector $\mathbf{v} \in \mathbb{R}^n$, LDOS is given as*

$$f(\lambda; \mathbf{v}) = \sum_{i=1}^n |\mathbf{v}^T \mathbf{u}_i|^2 \delta(\lambda - \lambda_i). \quad (5)$$

The following related equalities can be derived easily respectively for DOS and LDOS.

$$\int \phi(\lambda) f(\lambda) = \sum_{i=1}^n \phi(\lambda_i) = \text{trace}(\phi(\mathbf{S})) \quad (6)$$

$$\int \phi(\lambda) f(\lambda; \mathbf{v}) = \sum_{i=1}^n \phi(\lambda_i) (\mathbf{v}^T \mathbf{u}_i) (\mathbf{u}_i^T \mathbf{v}) = \mathbf{v}^T \phi(\mathbf{S}) \mathbf{v} \quad (7)$$

Scaling (L)DOS. The extremal (i.e., a few top largest or smallest) eigenpairs of various graph matrices have been associated with important graph characteristics, such as small-cut partitions [6], convergence rate of random walks to stationarity [7], unfolding of dynamical processes and epidemic thresholds [8], etc. Obtaining those few eigenpairs is also computationally easy. On the other hand, (L)DOS provides the distribution of the *entire* spectrum, which opens the door for the analysis of graph properties that are not evident from only the extremal eigenpairs. However, computing all n eigenvalues

and eigenvectors of a graph with n nodes is considerably more demanding. Therefore, analyzing large graphs through their density of states has been obstructed by the lack of scalable algorithms, until recently.

In their award-winning work, Dong *et al.* [10] introduced highly-efficient approximation algorithms to compute spectral densities, scalable to graphs with as large as tens of millions of nodes and billions of edges. Their main focus has been scaling the computation of these functions, with approximation-error analysis on plain graphs. In this paper, we capitalize on their work for speed and extend it to leverage node attributes for the first time for fast, attributed graph-level embedding.

III. GRAPH-LEVEL EMBEDDING WITH A-DOGE

A. Motivation

Our spectrally-designed A-DOGE derives graph-level features based on the node attributes and the *entire* spectrum of $\widetilde{\mathbf{W}}$ (can be other symmetric graph matrix, w.l.o.g. referred as \mathbf{S} , see §II), where the spectrum is composed of *all* the eigenvalues. Before delving into details, we discuss the motive for using the full spectrum and present an illustrative example.

Why the entire spectrum? We design graph-level features based on all of the eigenpairs of a graph matrix for two primary reasons. First, a large number of studies have found that the full eigenvalue spectra of different classes of real-world networks differ considerably [9], [28], [29], [30]. This suggests that the spectra can play a key discriminative role. Second, real-world networks are observed to exhibit localization on low-order eigenvectors, which are those eigenvectors associated with the non-extremal eigenvalues (in the sense of being the largest or smallest), but that are “buried” further down in the eigenvalue spectrum [31]. Notably, they capture mesoscopic inhomogeneity in networks which is defined as topologically distinct groupings of nodes, from few nodes to large modules, communities, or different interconnected subnetworks [32].

Illustrative example: To illustrate the valuable information that non-extremal eigenpairs carry, we present a visual analysis of low-order eigenvector localization using the MIG graph (See §IV-A). It consists of the counties across 49 mainland U.S. states as nodes, and an edge depicts the total number of people that migrated between two counties during 1995-2000 [31].

Eigenvector localization arises when most of the entries of an eigenvector are zero or near-zero, and implies that the nonzero components of the eigenvector coincide with a particular set of geometrically distinguished nodes in the graph. Extremal eigenvectors typically exhibit low localization; as shown in Fig. 2(a), the 2nd eigenvector has many non-zeros and mainly captures macroscopic properties, in this case, the graph cut depicting relatively fewer migrations between west- and east-coasts. A lower-order eigenvector, namely the 41st in (b) reflects mesoscopic structure in terms of migration patterns in and around South Dakota. 128th eigenvector in (c) has even larger localization, narrowing in a few counties within Texas near Austin, reflecting microscopic patterns. It is remarkable that the low-order eigenvectors align with geographical and political boundaries, carrying useful information at multiple scales.

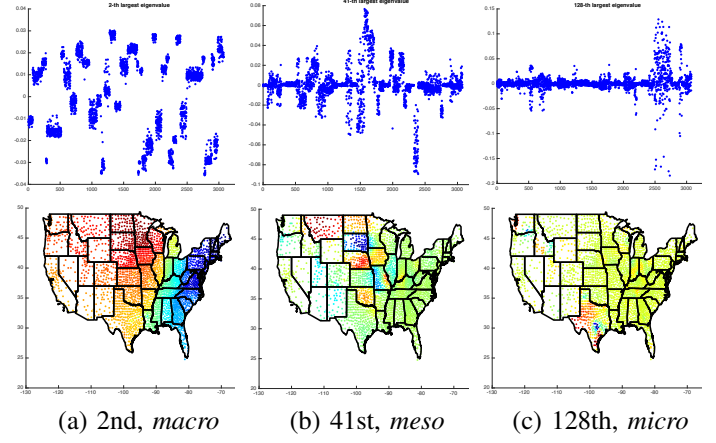


Figure 2. (top) Eigenvector entries (y-axis) versus node (i.e. U.S. county) index (x-axis) for the top (from left to right) 2nd, 41st, and 128th eigenvectors of the MIG graph (See §IV-A); (bottom) Eigenvector entries as heatmap (red:high to blue:low) for nodes (U.S. counties) shown in 2-d coordinates, solid black lines depict official U.S. state borders (best in color).

B. Spectrum as Histogram

Density of states (DOS) or spectral density as given in Eq. (4) is a continuous probability density function $f(\lambda)$ of the eigenvalues. We represent it with a histogram density estimator, denoted $h^{DOS}(\lambda)$ that partitions the eigenvalue range $[-1, 1]$ for $\widetilde{\mathbf{W}}$ into $B = 2/\Delta$ disjoint bins of equal width Δ . Let us denote their centers by $\tilde{\lambda}_b$ for $b \in \{1, \dots, B\}$. For any $i \in \{1, \dots, n\}$, let $\text{Bin}(\lambda_i)$ denote the bin that λ_i belongs to. We define our DOS histogram features for a graph as follows.

Definition 5 (DOS histogram features). *DOS histogram is a B -dimensional vector, denoted $\mathbf{h}^{DOS} \in \mathbb{R}^B$, where*

$$\mathbf{h}^{DOS}(\tilde{\lambda}_b) = \frac{1}{\Delta} \frac{\sum_{i=1}^n \mathbb{I}(\lambda_i \in \text{Bin}(\tilde{\lambda}_b))}{n}, b \in \{1, \dots, B\} \quad (8)$$

We also represent the local density of states (LDOS) in Eq. (5) similarly and define LDOS histogram features.

Definition 6 (LDOS histogram features). *For a given vector $\mathbf{v} \in \mathbb{R}^n$, the LDOS histogram is a B -dimensional vector, denoted $\mathbf{h}_v^{LDOS} \in \mathbb{R}^B$, where*

$$\mathbf{h}_v^{LDOS}(\tilde{\lambda}_b) = \frac{1}{\Delta} \frac{\sum_{i=1}^n |\mathbf{v}^T \mathbf{u}_i|^2 \mathbb{I}(\lambda_i \in \text{Bin}(\tilde{\lambda}_b))}{n}, \forall b \quad (9)$$

Note that by abusing convention slightly, we use the word histogram to refer to Eq. (9) although it is not a normalized density mass function. Fig. 3 shows examples to DOS (top) and LDOS (bottom) histograms with $B = 40$ each.

Computing both histograms requires all of the eigenvalues $\lambda_i, i = \{1 \dots n\}$ for a graph with n nodes. Further, LDOS requires all the corresponding eigenvectors \mathbf{u}_i ’s. For even moderate size graphs, computing the complete set of eigenpairs is prohibitive. Most recently, Dong *et al.* [10] introduced fast and scalable approximation algorithms to estimate these spectral densities. Our work is inspired by and builds on their work to efficiently obtain both \mathbf{h}^{DOS} and \mathbf{h}^{LDOS} based on the Gauss Quadrature and Lanczos (GQL) algorithm [33].

On the other hand, both in [10] and their follow-up work [19], $\mathbf{v} = \mathbf{e}_i$ is used in Eq. (5) to capture the spectral information about each particular node $i = \{1, \dots, n\}$, called point-wise

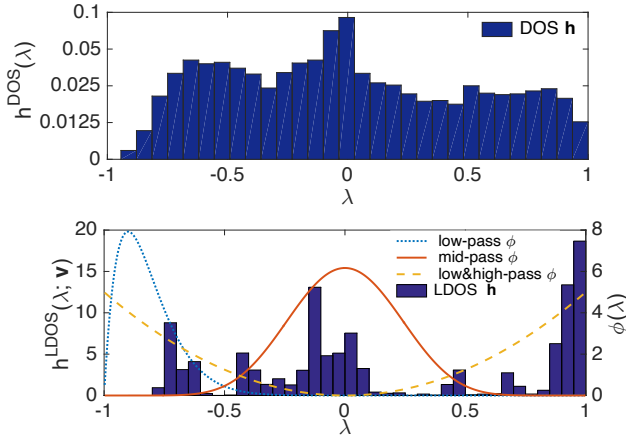


Figure 3. Example (top) histogram density estimator of the spectral density, a.k.a. density of states (DOS), of a graph for symmetrically-normalized $\tilde{\mathbf{W}}$ with eigenvalues (i.e. frequencies) in $[-1, 1]$, (bottom) local density of states (LDOS) histogram for a given vector \mathbf{v} . Also shown (in color) are three different frequency response functions depicted by curves over the spectrum.

density of states (PDOS), where \mathbf{e}_i is the i -th standard basis vector with i -th entry equal to 1 and 0 elsewhere. As such, both work are limited to plain graphs without node labels/attributes. We extend the use of LDOS to attributed graphs for the first time, by setting $\mathbf{v} \in \mathbb{R}^n$ in Eq. (9) to capture a graph signal on the nodes associated with an attribute.

Specifically, given a categorical or binary attribute a_j , we create a separate \mathbf{v} for each unique value $val \in \text{dom}(a_j)$ where $\mathbf{v}_i := 1$ if $\mathbf{X}_{ij} = val$ and 0 otherwise. For numerical attributes, we set $\mathbf{v} := \bar{\mathbf{X}}_{\cdot j}$ where $\bar{\mathbf{X}}$ denotes the column-wise standardized attribute matrix. Notably, LDOS can be extended to structural node-level attributes, such as degree or other node centrality measures, eccentricity, etc.

Interpreting LDOS. There is an intuitive interpretation of a LDOS feature in Eq. (9). The term $\mathbf{v}^T \mathbf{u}_i$, that is the dot product between an attribute vector and a graph eigenvector, is to reflect the *alignment* between attribute values and the structurally distinct group of nodes that the eigenvector captures. The better the alignment, the larger is the LDOS feature value for the bin that the corresponding eigenvalue falls into.

In addition to original LDOS, we also create *interaction* features between pairs of node attributes. Accordingly, the coupled-LDOS histogram features are defined as follows.

Definition 7 (cLDOS histogram features). *For two input vectors $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^n$, the coupled-LDOS histogram is a B -dimensional vector, denoted $\mathbf{h}_{\mathbf{v}, \mathbf{v}'}^{cLDOS} \in \mathbb{R}^B$, where*

$$\mathbf{h}_{\mathbf{v}, \mathbf{v}'}^{cLDOS}(\tilde{\lambda}_b) = \frac{1}{\Delta} \frac{\sum_{i=1}^n (\mathbf{v}^T \mathbf{u}_i)(\mathbf{u}_i^T \mathbf{v}') \mathbb{I}(\lambda_i \in \text{Bin}(\tilde{\lambda}_b))}{n}, \forall b \quad (10)$$

Note that $(\mathbf{v}^T \mathbf{u}_i)(\mathbf{u}_i^T \mathbf{v}')$ is the i -th entry of $\hat{\mathbf{c}}_{\mathbf{v}, \mathbf{v}'}$ (from Eq. 3). Hence $\mathbf{h}_{\mathbf{v}, \mathbf{v}'}^{cLDOS}$ can simply be viewed as $\hat{\mathbf{c}}_{\mathbf{v}, \mathbf{v}'}$, binned according to the corresponding eigenvalues of each entry.

Moreover, recall that we use the GQL algorithm to approximate the LDOS features, where the terms $\mathbf{v}^T \mathbf{u}_i$ or $\mathbf{u}_i^T \mathbf{v}'$ are not computed using the individual eigenvectors explicitly. Nevertheless it is easy to acquire cLDOS features in Eq. (10) using the LDOS features in Eq. 9 and simple algebra. Given

the separate LDOS features for \mathbf{v} and \mathbf{v}' , we also create those for $(\mathbf{v} + \mathbf{v}')$. Then,

$$\mathbf{h}_{\mathbf{v}, \mathbf{v}'}^{cLDOS} = [\mathbf{h}_{\mathbf{v} + \mathbf{v}'}^{LDOS} - \mathbf{h}_{\mathbf{v}}^{LDOS} - \mathbf{h}_{\mathbf{v}'}^{LDOS}] / 2. \quad (11)$$

C. Functions over the Spectrum

DOS, LDOS and cLDOS histograms provide “raw” information about the graph spectrum and the attributes. In addition, we define aggregate scalar features by specifying various frequency response functions (FRF) [25] (Eq. (2)) over these histograms.

Definition 8 ((cL)DOS aggregate features). *Given a DOS, LDOS or cLDOS histogram $\mathbf{h} \in \mathbb{R}^B$, and a frequency response func. $\phi(\cdot)$, a (cL)DOS aggregate feature $g_\phi \in \mathbb{R}$ is written as*

$$g_\phi = \sum_{b=1}^B \mathbf{h}(\tilde{\lambda}_b) \phi(\tilde{\lambda}_b) \quad (12)$$

Each FRF $\phi(\cdot)$ focuses on a different part of the spectrum, inducing a variety of graph filters. In Fig. 3 (bottom) we show three example FRFs; a low-pass one (blue) that has high values for smaller eigenvalues, a mid-pass one (red) as well as one that is both low-and-high pass (orange). To extract graph connectivity and attribute information broadly, we construct a *portfolio* of these graph filters, i.e. associated FRF’s $\{\phi(\cdot)\}$, called a **filterbank**.

Before delving into the details of our filterbank, we make a few remarks. First, note that the sum in Eq. (12) is an approximation of the integral in Eq. (6) for \mathbf{h}^{DOS} , that of Eq. (7) for \mathbf{h}^{LDOS} , and accordingly an approximation of $\mathbf{v}^T \phi(\mathbf{S}) \mathbf{v}'$ for \mathbf{h}^{cLDOS} . Second, given the efficiently-computed histograms thanks to the GQL algorithm, computing the aggregate features by Eq. (12) is extremely fast and simply involves a weighted sum. This allows us to employ a large filterbank containing many different FRF’s almost for “free”. Finally, we have seen that our cLDOS aggregate features relate to graph signal convolution. Denoting the vector of frequency responses by $\phi := \{\phi(\lambda_i)\}_{i=1}^n$, based on Eq.s (7) and (3),

$$\int \phi(\lambda) f(\lambda; \mathbf{v}, \mathbf{v}') = \sum_{i=1}^n \phi(\lambda_i) (\mathbf{v}^T \mathbf{u}_i) (\mathbf{u}_i^T \mathbf{v}') = \phi^T \hat{\mathbf{c}}_{\mathbf{v}, \mathbf{v}'} \quad (13)$$

In the following, we present two classes of FRF’s that A-DOGE employs to extract (cL)DOS aggregate features.

1) **Chebyshev polynomials:** We use the series of Chebyshev polynomials as a set of FRF’s defined by the recurrence $\phi_1(\lambda) = 1$, $\phi_2(\lambda) = 2(\lambda/\lambda_{\max}) - 1$, and $\phi_k(\lambda) = 2\phi_2(\lambda)\phi_{k-1}(\lambda) - \phi_{k-2}(\lambda)$, where λ_{\max} is the maximum eigenvalue.

Interpretation. As shown in Fig. 4(a), Chebyshev polynomials provide frequency profiles that cover various parts of the spectrum. For example, the 2nd one is mostly a low- and high-pass filter and stops the middle band, while the 3rd one passes the middle bands as well as very high and very low bands of the spectrum, and so on. Given a number of these FRF’s, emphasis can be put on passing/stopping specific bands by a weighted combination of them.

The flexibility of any-band filtering by A-DOGE is favorable over several modern graph neural networks (GNNs). GCN [20], for instance, works as a low-pass-only filter and hence does

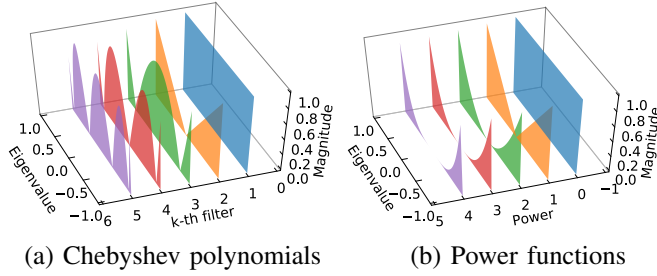


Figure 4. Frequency response functions $\phi_k(\lambda)$ (magnitude in absolute val.s)

Table II

SUMMARY OF A-DOGE GRAPH-LEVEL FEATURES BASED ON SPECTRAL DENSITIES (cL)DOS, ORGANIZED AS HISTOGRAM AND AGGREGATE FEATURES USING CHEBYSHEV POLYNOMIALS AND POWER FUNCTIONS.

DOS			LDOS			cLDOS		
hist	agg. g_ϕ		hist	agg. g_ϕ		hist	agg. g_ϕ	
	Cheb.	Pow.		Cheb.	Pow.		Cheb.	Pow.
B	K	K	BD	KD	KD	$B\binom{D}{2}$	$K\binom{D}{2}$	$K\binom{D}{2}$

not cover the whole spectrum. GIN [21] has a learnable scalar parameter ϵ that determines which band to stop, however its FRF is a linearly decreasing function, which is not a strong low-pass or high-pass filter. (See Fig. 2 in [25].) In contrast, spectrally-designed ChebNet [22] is more expressive and also employs the Chebyshev polynomials. We compare to these modern GNNs in the experiments on graph classification tasks.

2) **Power functions:** The second class of FRF’s in our filterbank uses (both positive and negative) powers of the spectrum, that is,

$$\phi_k(\lambda) = \lambda^k, \quad k = \pm\{1, \dots, K/2\}$$

Interpretation. Our aggregate features using the power functions relate to random walks on the graph. Consider positive values of k and $\mathbf{S} = \tilde{\mathbf{W}}$. Recall that for \mathbf{h}^{DOS} , Eq. (12) is an approximation of $\text{trace}(\phi(\mathbf{S}))$, which is equal to the total return-probability of a k -step random walk to a node. For \mathbf{h}^{LDOS} , aggregate features approximate $\mathbf{v}^T \phi(\mathbf{S}) \mathbf{v}$. For a binary/categorical attribute where \mathbf{v} depicts a certain value, e.g. $\text{val} := (\text{party_affiliation} : \text{democrat})$, it corresponds to the probability that a k -step random walk starting at any node with value val “hits”/reaches another node with the same value. For $\mathbf{h}^{\text{cLDOS}}$, similarly, it is the probability that such a walk starting at any node with a certain val will reach another node with a different val' . Moreover, for two continuous attributes \mathbf{v} and \mathbf{v}' , approximating $\mathbf{v}^T \phi(\mathbf{S}) \mathbf{v}'$ via $\mathbf{h}^{\text{cLDOS}}$ would capture the covariance of the attributes over “ k -hop connected” pairs of nodes that can reach each other within k -steps.

The interpretations extend to the negative powers as well, which correspond to many walks of different lengths in the limit. In that respect, aggregate features using power functions depict *multi-scale* properties, where increasingly positive values of k capture microscopic to mesoscopic properties related to short/local random walks, whereas negative powers relate to the long-range walks and thereby macroscopic structure.

We conclude with an overview of all the graph-level features described in this section. Table II gives the number of features by category, where B is the number of histogram bins, K is

Input: Graph G (with D node attributes), parameters B, K .

- (Preprocess) Compute $2K$ agg. functions on B bin centers.
- $\tilde{\mathbf{W}} \leftarrow$ normalized adjacency matrix of G
- $\mathbf{v}_d \leftarrow$ attribute vector for each $d \in \{1 \dots D\}$
- Compute \mathbf{h}^{DOS} using GQL [10] on $\tilde{\mathbf{W}}$
- For each d , compute $\mathbf{h}_{\mathbf{v}_d}^{\text{LDOS}}$ using GQL on $\tilde{\mathbf{W}}, \mathbf{v}_d$.
- For each pair $d, d' \in \{1 \dots D\}$, compute $\mathbf{h}_{\mathbf{v}_d + \mathbf{v}_{d'}}^{\text{LDOS}}$ using GQL and then $\mathbf{h}_{\mathbf{v}_d + \mathbf{v}_{d'}}^{\text{cLDOS}}$ using Eq. (11).
- For each histogram computed above, dot product with all aggregate functions to produce aggregate features.

Figure 5. Steps to generate all A-DOGE features (See Table II)

the number of Chebyshev/power FRF’s, and $D \geq d$ is the total number of attributes upon one-hot-encoding the categorical and binary attributes. A-DOGE yields $(B + 2K)(1 + D + \binom{D}{2})$ features in total for an attributed graph, which are permutation- and size-invariant, task-agnostic, variable band-pass, multi-scale, and extremely efficient to compute. We outline the steps in Fig. 5 and give detailed complexity analysis next.

D. Computational Complexity

Since A-DOGE computes an embedding for each graph independently, it scales linearly with the number of graphs in the dataset, i.e., N .

We analyze the asymptotic runtime of A-DOGE on a single graph G with n nodes, m edges, and D total node attributes (including one-hot encoded labels and categorical attributes). We use the Gauss Quadrature and Lanczos algorithm described by Dong *et al.* [10] to compute a (cL)DOS *histogram*. This involves (i) running η_L Lanczos iterations, each requiring $O(m)$ operations, followed by (ii) the eigendecomposition of a tridiagonal $n \times n$ matrix, with $O(n^2)$ operations. Note that although a tridiagonal matrix eigendecomposition has a quadratic worst-case complexity theoretically, this operation is extremely fast in practice – especially for real-world matrices. Each *aggregate feature* requires a dot product of two vectors of size B for $O(B)$, where we use $2K$ different frequency response functions (i.e. $\phi(\cdot)$ ’s) in total (K each for Chebyshev and powers). Then, the total complexity of computing one histogram and its related aggregate features is $O(n^2 + \eta_L m + KB)$. This gives a total runtime of $O((n^2 + \eta_L m + KB) \cdot \alpha)$, where α denotes the number of desired graph-level features (i.e. embedding size) in A-DOGE.

Notably, A-DOGE is modular and can include any subset of the features in Table II. For datasets with a large number of node attributes, one can skip cLDOS features, or only choose important attribute-pairs to ensure $\alpha = O(D)$. Also note that each aggregate feature for a given $\phi(\cdot)$ can be computed independently, and hence can be easily parallelized.

IV. EXPERIMENTS

To evaluate A-DOGE we design both quantitative and qualitative experiments to answer the following questions.

Q1. Graph Classification How does A-DOGE (unsupervised) compare to the modern GNNs and graph kernels (un/supervised) on benchmark graph classification tasks?

Q2. Exploratory Graph Analysis Can A-DOGE provide insights for mining real-world attributed graphs?

Q3. Efficiency How fast and scalable is A-DOGE?

A. Experiment Setup

Datasets. List of datasets and summary are in Table III.

For graph classification, we use eight benchmark datasets from TUDataset repository². Five are commonly used social network datasets, REDDIT-B, REDDIT-5K, COLLAB, IMDB-BIN and IMDB-MUL. These contain only plain graphs—a setting with which all the baselines are compatible. The other three are biochemistry datasets, PROTEINS, DD and AIDS, which have node labels and/or attributes.

We also use four other graph datasets to specifically showcase the strengths of A-DOGE in leveraging the full graph spectrum.

- **BandPass** is a synthetic dataset consisting of images generated via sinusoidal patterns from two frequency ranges [25].
- **Congress** is based on the voting patterns in 41 U.S. Senates (1927–2008) [31], where nodes represent senators (labeled by party affiliation) and edge weights represent voting agreement. To create separate classes of graphs, we add noise to edge weights between same-party senators (class 1), and randomly picked senators (class 2) in one randomly picked congress.
- **Congress-1** is generated by picking one congress at random and shuffling the labels of senators via random swaps; 50 swaps in class 1, and 300 in class 2.
- **MIG** is based on the county-to-county migration in the U.S. [31]. To create separate classes of graphs, we add noise to edges between a pair of bordering states (class 1) or within a state (class 2).

In addition to the above datasets, we perform graph exploratory analysis using A-DOGE on two more datasets:

- **Facebook100** consists of Facebook college social networks from 100 American institutions [34], with student demographic information (major, dorm, status, class-year, etc.) as node attributes.
- **BorderStates** is built from the MIG dataset, by inducing 49 separate graphs - one for each mainland state and its bordering states. We label counties of the selected state as 1, and the counties of the neighbors as 2.

Baselines. We compare A-DOGE quantitatively to various unsupervised and supervised graph embedding, graph kernel, and graph neural network methods on graph classification tasks.

Unsupervised explicit graph embeddings are in the same category as A-DOGE and hence most comparable. As baselines from this category, we use • FGSD [11], • NETLSD [12] and • G2VEC [13], which we described briefly in §I-A.

Graph kernels are also unsupervised; here we use three of the best performing kernels on classification benchmarks, and a recent DOS-based graph kernel. • WL [14]: the Weisfeiler-Lehman graph kernel, • WL-OA [15]: the Weisfeiler-Lehman Optimal Assignment kernel, • PK [17]: the Propagation Kernel, and • DOSGK [19], the Density of States Graph Kernel.

GNN baselines include state-of-the-art supervised models, such as • CHEBNET [22], • GCN [20], and • GIN [21].

Note that FGSD, NETLSD, and DOSGK are for plain graphs only. G2VEC, WL, and WL-OA admit node labels

Table III
DATASET SUMMARY STATISTICS.

	N	Cls.	Avg. n	Avg. m	Lbl.	Attr.
REDDIT-B	2000	2	429.6	497.7	-	-
REDDIT-5K	5000	5	508.5	594.8	-	-
COLLAB	5000	3	74.5	2457.8	-	-
IMDB-BIN	1000	2	19.8	96.5	-	-
IMDB-MUL	1500	3	13.0	65.9	-	-
DD	1178	2	284.3	715.7	78	-
PROTEINS	1113	2	39.1	72.8	-	1
AIDS	2000	2	15.7	16.2	38	4
BandPass	5000	2	200	1072.6	-	1
Congress	200	2	4196	450662.5	3	-
MIG	200	2	3075	1092282.0	19	-
Facebook100	100	n/a	12083.2	469845.4	-	7
BorderStates	49	n/a	367.4	21633.9	2	-

but not (continuous) attributes. Therefore, they input only the admissible parts of a graph dataset for classification.

Model configuration. In our experiments with A-DOGE, we set $\eta_L=100$, $B=200$ and $K=100$ (see Table II). For plain datasets, we use node degree as a continuous attribute. For FGSD, we use L^{-1} as the distance function and 0.001 as the binwidth. For NETLSD, we use heat trace signatures at 250 different values of t logarithmically spaced in $[10^{-2}, 10^2]$. For G2VEC, we set the WL iteration count to 5 and output dimension to 1024. For the kernels WL, WL-OA and PK, we use the implementation from the GraKel package³, and the default parameters suggested. For DOSGK, same as with A-DOGE, we use 200 bins and 100 Chebyshev moments. For all the GNNs, we use mean-pooling as the readout function.

We run all non-GNN experiments on one core of Intel(R) Xeon(R) CPU E5-2667 v3 CPU @3.20GHz. GNN experiments are run on a server with NVIDIA Tesla V100 GPU and one core of Intel(R) Xeon(R) Gold 6248 CPU @2.50GHz.

B. Graph Classification

Classifier configurations. For classification with the embeddings produced by unsupervised methods, we use the kernel-SVM⁴ classifier with the regularization parameter C chosen from $\{10^{-3}, 10^{-2}, \dots, 10^3\}$ via 10-fold cross-validation. We perform this experiment 10 times using random splits. For explicit embeddings, we normalize each feature, and set γ to be the inverse of the median of pairwise ℓ_2 distances between all embeddings. For A-DOGE, we also set the option of using LDOS, cLDOS features, and the option of using aggregate FRFs as hyperparameters. We normalize all kernels symmetrically. For GNNs, we train them end-to-end using cross-entropy loss, and hyperparameters (learning-rate at 0.005, layers in $\{2, 3, 5, 7\}$, hidden sizes from $\{32, 64, 128\}$ and epochs up to 200) selected via 10-fold cross-validation. For each of the above methods, we report the mean test accuracy for the best choice of hyperparameters, and the corresponding standard deviation on every dataset except BandPass, for which we use the single train-validation-test split as specified in [25].

³<https://ysig.github.io/GraKel/>

⁴SVM facilitates comparable results between implicit and explicit kernels.

²<https://chrsmrrs.github.io/datasets/docs/datasets/>

Table IV

GRAPH CLASSIFICATION PERFORMANCE BY A-DOGE AND ITS DOS-ONLY (I.E. NO ATTRIBUTES) VARIANT DOGE, COMPARED WITH THREE TYPES OF BASELINES. THE HIGHEST PERFORMANCE PER DATASET IS IN **BOLD** AND THE HIGHEST AMONG UNSUPERVISED METHODS IS UNDERLINED. E DENOTES THE CODE OUTPUTTING AN ERROR; THE NUMBERS WITH SYMBOLS DENOTE THE PAPER FROM WHICH THE NUMBERS ARE TAKEN: [‡][24], * [19], [†][25].

	Graph Embedding (Unsupervised)					Graph Kernels (Unsupervised)				GNNs (Supervised)		
	A-DOGE	DOGE	FGSD	NETLSD	G2VEC	WL	WL-OA	PK	DOSGK	CHEBNET	GCN	GIN
RED-B	91.6 (1.5)	90.3 (1.8)	82.4 (2.6)	85.6 (2.2)	74.2 (2.7)	83.9 (0.5) [‡]	88.9 (0.1) [‡]	85.5 (0.3) [‡]	88.8 (0.3)*	90.2 (2.0)	89.9 (2.0)	91.7 (1.6)
RED-5K	55.6 (2.2)	53.8 (2.1)	47.0 (1.8)	45.9 (2.1)	41.5 (1.6)	51.2 (0.3)*	E	E	52.8 (0.2)*	55.0 (2.2)	54.2 (1.7)	54.7 (2.0)
COLLAB	72.2 (2.0)	72.2 (2.0)	70.2 (1.8)	68.4 (1.9)	57.9 (1.5)	74.8 (0.2)*	79.8 (1.6)	77.8 (1.7)	<u>80.8</u> (0.2)*	84.6 (1.1)	84.2 (1.2)	83.8 (1.6)
IMDB-B	72.6 (4.3)	71.6 (4.3)	70.6 (4.1)	69.7 (4.1)	56.0 (4.1)	71.3 (1.0) [‡]	<u>73.5</u> (0.6)	71.2 (0.7) [‡]	72.8 (0.9)*	80.2 (3.9)	79.9 (3.7)	80.8 (4.5)
IMDB-M	47.8 (3.5)	47.6 (3.7)	48.6 (3.4)	47.9 (3.7)	44.4 (3.8)	50.7 (0.6) [‡]	50.7 (0.5) [‡]	<u>51.0</u> (0.7) [‡]	49.4 (0.5)*	55.6 (2.7)	55.2 (2.7)	56.3 (3.1)
DD	80.1 (3.5)	76.2 (3.4)	76.5 (3.5)	76.6 (3.5)	76.2 (3.5)	80.9 (0.3)	79.9 (0.5)	81.6 (0.5)	73.4 (3.7)	78.9 (1.9)	78.0 (1.8)	79.3 (1.9)
PROTN	74.9 (3.5)	74.9 (3.5)	74.2 (3.3)	74.5 (4.0)	72.1 (3.1)	73.9 (0.7) [‡]	<u>75.9</u> (0.6) [‡]	74.6 (0.5) [‡]	72.1 (3.9)	78.3 (2.7)	76.7 (3.5)	78.4 (3.9)
AIDS	99.8 (0.3)	99.8 (0.3)	99.6 (0.4)	99.6 (0.5)	98.8 (0.7)	99.7 (0.0) [‡]	99.7 (0.0) [‡]	99.7 (0.0) [‡]	99.1 (0.7)	96.9 (1.6)	95.5 (1.3)	98.6 (0.6)
Cong	99.5 (1.5)	54.7 (11.0)	95.1 (4.3)	99.5 (1.5)	86.8 (7.4)	84.8 (7.3)	81.1 (7.7)	68.6 (8.3)	60.0 (10)	50.0 (0.0)	50.0 (0.0)	57.0 (5.9)
Cong-1	78.0 (8.6)	58.9 (10.0)	50.0 (0.0)	60.4 (9.7)	59.8 (11)	62.2 (10)	62.3 (10)	58.2 (10)	55.7 (9.7)	50.0 (0.0)	50.0 (0.0)	71.5 (9.4)
MIG	100.0 (0)	62.3 (9.7)	99.5 (1.5)	99.9 (1.1)	50.0 (0.0)	99.8 (1.4)	99.8 (1.4)	100 (0.0)	53.5 (12)	100.0 (0.0)	78.5 (1.7)	100.0 (0.0)
BPass	<u>90.8</u>	51.9	47.9	51.4	50	50	51.6	70.4	48.5	98.2 [†]	77.9 [†]	87.6 [†]
Avg.	82.5	69.1	74.1	75.8	66.0	75.6	76.6	76.3	68.6	78.4	74.2	80.5

Results. Table IV contains all the performance results of our classification experiments. Among the benchmark datasets, A-DOGE achieves on par performance with the most competitive unsupervised baselines and is often comparable to (supervised) GNNs, while being considerably more resource-frugal.

On the other four datasets, A-DOGE significantly outperforms all baseline methods due to its ability to capture the alignment of labels and attributes with graph structure at a multi-scale level, even in databases with as few as 200 graphs. Provided A-DOGE uses considerably lower resources in comparison with kernels and GNNs, and considering that the latter are trained end-to-end, we do not expect A-DOGE to exhibit state-of-the-art performance on every dataset. Still, A-DOGE outperforms/equals all baselines on 7 of the datasets. Moreover, A-DOGE stands out as the top choice based on average performance across all datasets.

On the BandPass dataset, only the spectrally-designed CHEBNET is able to outperform A-DOGE. This can be attributed to the way that BandPass is created, wherein graph classes are formed based on the frequency band used to generate the underlying image. Fig. 6 depicts the LDOS histograms of the graphs in the BandPass dataset. We can clearly see that capturing specific bands of the eigenspectrum suffices to characterize the disparity between the two graph classes.

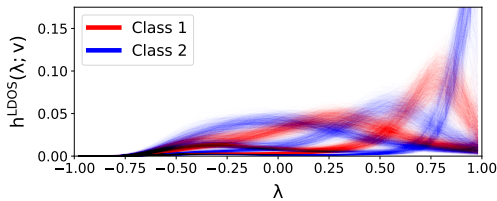


Figure 6. LDOS histograms for all graphs in BandPass, plotted as lines. Each class can be characterized by specific bands of eigenvalues.

Feature ablation. Table IV also shows the DOS-only version of A-DOGE without using node labels and attributes, called DOGE. We observe that in the benchmark datasets, graph structure seems to hold most of the useful information needed for classification, and hence there is only a small improvement

in performance from using node attributes. In the rest of the datasets, node attributes play an important role, causing significant improvements in results for A-DOGE by using LDOS and cLDOS features.

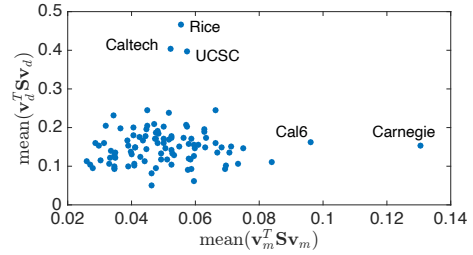


Figure 7. Average homophily w.r.t. *major* vs. *dorm* in 100 Facebook college social networks in the U.S., where \mathbf{v}_m and \mathbf{v}_d respectively refer to an attribute vector corresponding to a particular *major* m and *dorm* d .

C. Graph Data Mining

To demonstrate the interpretability of the A-DOGE features, we perform exploratory graph analysis on three real-world datasets, Facebook100, Congress and BorderStates. **Facebook100.** In Facebook100, we denote each categorical feature (e.g. *major*) with its one-hot encoding, and hence, each particular value (e.g. Computer Science) has its own (binary) attribute vector. We first visualize the Facebook100 graphs via LDOS aggregate features using these attribute vectors, with small positive power functions as FRF to capture the assortativity (homophily) of different attributes across different college networks. In each graph, we compute the aggregate feature that estimates $\mathbf{v}_m^T \mathbf{S} \mathbf{v}_m$ for every *major* captured by \mathbf{v}_m , and similarly $\mathbf{v}_d^T \mathbf{S} \mathbf{v}_d$ for every *dormitory* captured by \mathbf{v}_d . Fig. 7 plots the mean homophily with respect to *major* and *dorm* for each of the 100 colleges.

While Carnegie pops up as having the highest correlation between edges and students with the same *major*, comparing the ranges of both axes suggests that *dorm* is a much stronger indicator of students within a college being friends. Moreover, this tendency seems to be more pronounced in Rice, Caltech and UCSC. This is also backed up by findings in [34] and

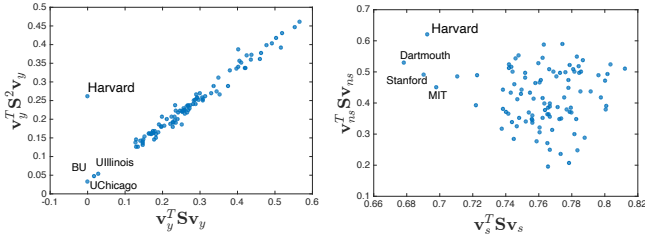


Figure 8. (left) Homophily w.r.t. *class_year* based on $k=1$ and $k=2$ -length paths over all 100 colleges. \mathbf{v}_y refers to continuous attribute vector with class years. (right) Homophily within student and non-student communities in all 100 colleges. Binary vector \mathbf{v}_s (\mathbf{v}_{ns}) depicts student (non-student) status.

the real-world knowledge that Rice and Caltech are organized predominantly by dorms and other on-campus housing.

We also analyze similar aggregate functions over the continuous attributes. Fig. 8(left) plots the assortativity with respect to *class_year* for $k=1$ and $k=2$ for the power functions, which capture 1- and 2-length paths. As we expect, these features are highly correlated in most colleges—with the striking exception of Harvard, where it appears that 2-length paths are common between individuals of similar *class_year*, but this is not the case with 1-length paths. To investigate further, we plot homophilies for student and non-student populations for all colleges in Fig. 8(right) and we learn that the Harvard network consists of a comparatively higher number of edges amongst non-student members, most of whom have empty or very disparate *class_year*. Even if edges between students are fewer, this is corrected when we look at 2-length paths instead.

Congress. Next, we want to explore scenarios where *interactions between attributes* prove important to understanding properties of a graph. To this end, we look at the Congress graph, where the two attribute vectors are binary vectors \mathbf{v}_d and \mathbf{v}_r , corresponding to Democrat and Republican senators respectively (ignoring the small minority of independents). We plot within-party agreement ($(\mathbf{v}_d^T \mathbf{S} \mathbf{v}_d + \mathbf{v}_r^T \mathbf{S} \mathbf{v}_r)/2$) and cross-party agreement ($\mathbf{v}_d^T \mathbf{S} \mathbf{v}_r$) over the years in Fig. 9.

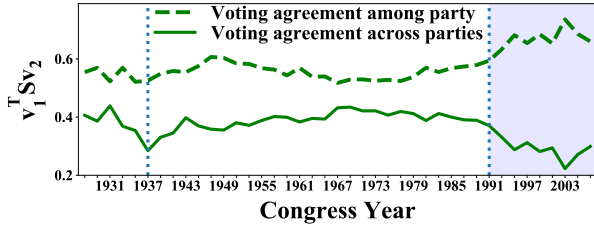


Figure 9. Voting agreement within (dashed curve, $\mathbf{v}_1=\mathbf{v}_2=\mathbf{v}_d$ or \mathbf{v}_r) and across (solid curve, $\mathbf{v}_1=\mathbf{v}_d$, $\mathbf{v}_2=\mathbf{v}_r$) political parties over the years, for 41 Senates during 1927–2008.

We can observe that beginning from the 1990s, senators tend to agree among their parties, and disagree with the opposite party to a higher extent, hinting at a growing polarization in politics. We note that agreement across parties is also low in 1937 (see the “dip”), however, this is better explained by the fact that this congress had overwhelmingly more number of democrats. There is no hint of polarization for that instance, since there is no corresponding rise in the dashed (within-party) curve. Fig. 9 shows that aggregate functions from A-DOGE not only help us observe such phenomenon but also help quantify them to a relative extent.

BorderStates. Lastly, we analyze BorderStates, comparing within-state migration against cross-border migration for each of the 49 mainland states in the U.S. We focus on LDOS aggregate features, this time using both positive and negative power functions, to analyze both short and long-range migration patterns. In other words, while small positive powers ($k=1$) capture local migration patterns, negative powers ($k=-1$) depicting paths of all lengths also reflect long-range migration behavior on a relatively global scale.

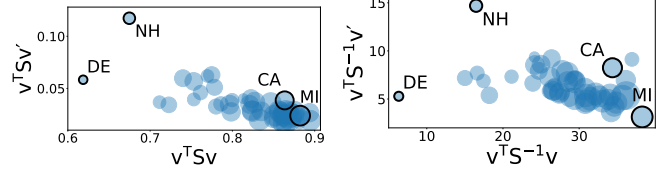


Figure 10. Comparison of migration patterns for each U.S. state – within its counties vs. across its borders; migration (left) over a local range, and (right) on a global scale. \mathbf{v}_w and \mathbf{v}_b refer to binary vectors denoting within and border-state counties, respectively. Node sizes correlate to size of state.

From Fig. 10(left), we observe that at the local scale, most states have greater within-state migration than cross-border migration. Comparatively, NH and DE, being the states with the least number of counties, exhibit lower within-state migration. Moreover, due to NH’s geographical and political similarity with its bordering states, it shows highest cross-border migration. On the other hand, larger states such as CA and MI exhibit mostly within-state migrations on the local scale. However, on the global scale (Fig. 10(right)), the difference between these is more pronounced, since CA is a more popular long-range migration destination than MI.

D. Scalability

A-DOGE is not directly comparable to all the baselines in terms of resource requirements. GNNs and G2VEC need GPU processing, which make them incomparable to CPU-based A-DOGE and the rest. Other differences, such as supervised training and collective processing of the graphs via multiple passes over the dataset (in contrast to one-by-one/independent processing by A-DOGE) put them in a different “league”.

On the other hand, kernel baselines need considerably more memory. WL, WL-OA and PK compute intermediate data (e.g. compressed labels) based on all the graphs in memory. These and DOSGK produce a $N \times N$ kernel matrix that is also memory-resident.

FGSD and NETLSD are comparable in the sense that, similar to A-DOGE, they process the graphs independently one-by-one. Likewise, they are also unsupervised. However, they cannot handle node labels/attributes. Nevertheless, we provide running time and scalability comparison in Fig. 11 that plots the runtime vs. size in number of nodes for individual graphs in the REDDIT-5K dataset. For any graph from the dataset (up to 9500 edges), A-DOGE does not take more than 1 second to compute. Fig. 1 compares this runtime for 3 of our largest datasets. We can see that A-DOGE achieves the best time-accuracy trade-off among competing baselines. For methods with comparable or better accuracy scores (e.g. GIN),

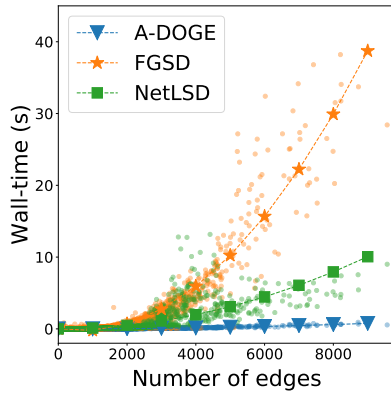


Figure 11. Runtimes per graph in the REDDIT-5K dataset for each of the A-DOGE and the two baselines which can compute embeddings independently.

A-DOGE is almost twice as fast on average. For baselines with similar runtime (e.g., WL), A-DOGE achieves significantly higher accuracy.

V. CONCLUSION

We propose A-DOGE, an unsupervised graph embedding technique designed to efficiently capture structural properties as well as node labels and attributes of a graph. To this end A-DOGE uses spectral density, or density of states (DOS), derived from the eigenspectrum of the graph, as a tool to capture both global and local properties of a graph. Further, we extend local density of states to leverage node labels and attributes, and capitalize on fast approximation algorithms making A-DOGE efficient and scalable to large graphs both in terms of time and space. Being unsupervised, it is not only suitable for downstream supervised graph classification tasks, but also applies well to exploratory graph analysis. Through both quantitative and qualitative experiments, we show the efficacy and efficiency of A-DOGE, where it outperforms unsupervised baselines and performs comparably to the supervised GNNs on graph classification tasks, and provides various insights into the analysis of real-world attributed graphs.

ACKNOWLEDGMENTS

This work is sponsored by NSF CAREER 1452425. We also thank PwC Risk and Regulatory Services Innovation Center at CMU. Any conclusions expressed in this material are those of the authors and do not necessarily reflect the views of the funding parties.

REFERENCES

- [1] N. Wale, I. A. Watson, and G. Karypis, "Comparison of descriptor spaces for chemical compound retrieval and classification," *Knowledge and Information Systems*, vol. 14, no. 3, pp. 347–375, 2008.
- [2] N. Przulj, "Biological network comparison using graphlet degree distribution," *Bioinform.*, vol. 26, no. 6, pp. 853–854, 2010.
- [3] H. Duen, N. Carey, W. Jeffrey, W. Adam, and F. Christos, "Polonium: tera-scale graph mining for malware detection," in *SIAM SDM*, 2011.
- [4] B. Ribeiro, N. Chen, and A. Kovacec, "Shaping graph pattern mining for financial risk," *Neurocomputing*, vol. 326, pp. 123–131, 2019.
- [5] P. Van Mieghem, *Graph spectra for complex networks*. Cambridge University Press, 2010.
- [6] A. Pothen, H. D. Simon, and K.-P. Liou, "Partitioning sparse matrices with eigenvectors of graphs," *SIAM J. on Matrix Anal. and Appl.*, vol. 11, no. 3, pp. 430–452, 1990.
- [7] J. A. Fill, "Eigenvalue bounds on convergence to stationarity for nonreversible markov chains, with an application to the exclusion process," *The Annals of Applied Probability*, pp. 62–87, 1991.
- [8] D. Chakrabarti, Y. Wang, C. Wang, J. Leskovec, and C. Faloutsos, "Epidemic thresholds in real networks," *ACM TISSEC*, vol. 10, no. 4, pp. 1–26, 2008.
- [9] S. Jin and R. Zafarani, "The spectral zoo of networks: Embedding and visualizing networks with spectral moments," in *KDD*, 2020, pp. 1426–1434.
- [10] K. Dong, A. R. Benson, and D. Bindel, "Network density of states," in *KDD*, 2019, pp. 1152–1161.
- [11] S. Verma and Z.-L. Zhang, "Hunt for the unique, stable, sparse and fast feature learning on graphs," in *NIPS*, 2017, pp. 88–98.
- [12] A. Tsitsulin, D. Mottin, P. Karras, A. Bronstein, and E. Müller, "Netlsd: hearing the shape of a graph," in *KDD*, 2018, pp. 2347–2356.
- [13] A. Narayanan, M. Chandramohan, R. Venkatesan, L. Chen, Y. Liu, and S. Jaiswal, "graph2vec: Learning distributed representations of graphs," *arXiv preprint arXiv:1707.05005*, 2017.
- [14] N. Shervashidze, P. Schweitzer, E. J. van Leeuwen, K. Mehlhorn, and K. M. Borgwardt, "Weisfeiler-lehman graph kernels," *J. Mach. Learn. Res.*, vol. 12, pp. 2539–2561, 2011.
- [15] N. M. Kriege, P.-L. Giscard, and R. C. Wilson, "On valid optimal assignment kernels and applications to graph classification," in *NIPS*, 2016, pp. 1615–1623.
- [16] L. Wu, Z. Zhang, A. Nehorai, L. Zhao, F. Xu, and A. S. Learning, "Sage: Scalable attributed graph embeddings for graph classification," in *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [17] M. Neumann, R. Garnett, C. Bauckhage, and K. Kersting, "Propagation kernels: efficient graph kernels from propagated information," *Mach. Learn.*, vol. 102, no. 2, pp. 209–245, 2016.
- [18] Z. Zhang, M. Wang, Y. Xiang, Y. Huang, and A. Nehorai, "RetGK: Graph kernels based on return probabilities of random walks," in *NeurIPS*, 2018, pp. 3968–3978.
- [19] L. Huang, A. J. Graven, and D. Bindel, "Density of states graph kernels," in *SDM*. SIAM, 2021, pp. 289–297.
- [20] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," in *ICLR*, 2017.
- [21] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How Powerful are Graph Neural Networks?" in *ICLR*, 2019, pp. 1–17.
- [22] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NIPS*, 2016, pp. 3837–3845.
- [23] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, "CayleyNets: Graph convolutional neural networks with complex rational spectral filters," *IEEE Trans. Sign. Process.*, vol. 67, no. 1, pp. 97–109, 2019.
- [24] N. M. Kriege, F. D. Johansson, and C. Morris, "A survey on graph kernels," *Appl. Netw. Sci.*, vol. 5, no. 1, p. 6, 2020.
- [25] M. Balcilar, R. Guillaume, P. Héroux, B. Gaüzère, S. Adam, and P. Honeine, "Analyzing the expressive power of graph neural networks in a spectral perspective," in *ICLR*, 2021.
- [26] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *ICLR*. OpenReview.net, 2019.
- [27] F. Wang and D. P. Landau, "Efficient, multiple-range random walk algorithm to calculate the density of states," *Physical Review Letters*, vol. 86, no. 10, p. 2050, 2001.
- [28] I. J. Farkas, I. Derényi, A.-L. Barabási, and T. Vicsek, "Spectra of real-world graphs: Beyond the semicircle law," *Physical Review E*, vol. 64, no. 2, 2001.
- [29] A. Banerjee and J. Jost, "Spectral plot properties: Towards a qualitative classification of networks," *Networks & Heterogeneous Media*, vol. 3, no. 2, p. 395, 2008.
- [30] P. N. McGraw and M. Menzinger, "Laplacian spectra as a diagnostic tool for network structure and dynamics," *Physical Review E*, vol. 77, no. 3, 2008.
- [31] M. Cucuringu and M. W. Mahoney, "Localization on low-order eigenvectors of data matrices," *arXiv preprint arXiv:1109.1355*, 2011.
- [32] M. Mitrović and B. Tadić, "Spectral and dynamical properties in classes of sparse networks with mesoscopic inhomogeneities," *Physical Review E*, vol. 80, no. 2, p. 026123, 2009.
- [33] G. H. Golub and G. Meurant, "Matrices, moments, and quadrature," in *Milestones in Matrix Computation*. Oxford U. Press, 2007, pp. 380–433.
- [34] A. L. Traud, P. J. Mucha, and M. A. Porter, "Social structure of facebook networks," *Physica A: Statistical Mechanics and its Applications*, vol. 391, no. 16, pp. 4165–4180, 2012.