

Graph Perceiver IO: A General Architecture for Graph-Structured Data

Seyun Bae^a (bsu1313@uos.ac.kr), Hoyoon Byun^b (hoyunb@gmail.com),
Changdae Oh^b (changdae.oh@uos.ac.kr), Yoon-Sik Cho ^c
(yoonsik@cau.ac.kr), Kyungwoo Song^{de} (kyungwoo.song@gmail.com)

^a Department of Computer Science and Engineering, University of Seoul

^b Department of Artificial Intelligence, University of Seoul

^c School of Computer Science and Engineering, Chung-Ang University

^d Department of Applied Statistics, Yonsei University

^e Department of Statistics and Data Science, Yonsei University

Corresponding Author:

Kyungwoo Song

Department of Applied Statistics, Yonsei University, 50, Yonsei-ro, Seodaemun-gu, Seoul, Republic of Korea

Tel: +82-2-2123-2473

Email: kyungwoo.song@gmail.com

Yoon-Sik Cho

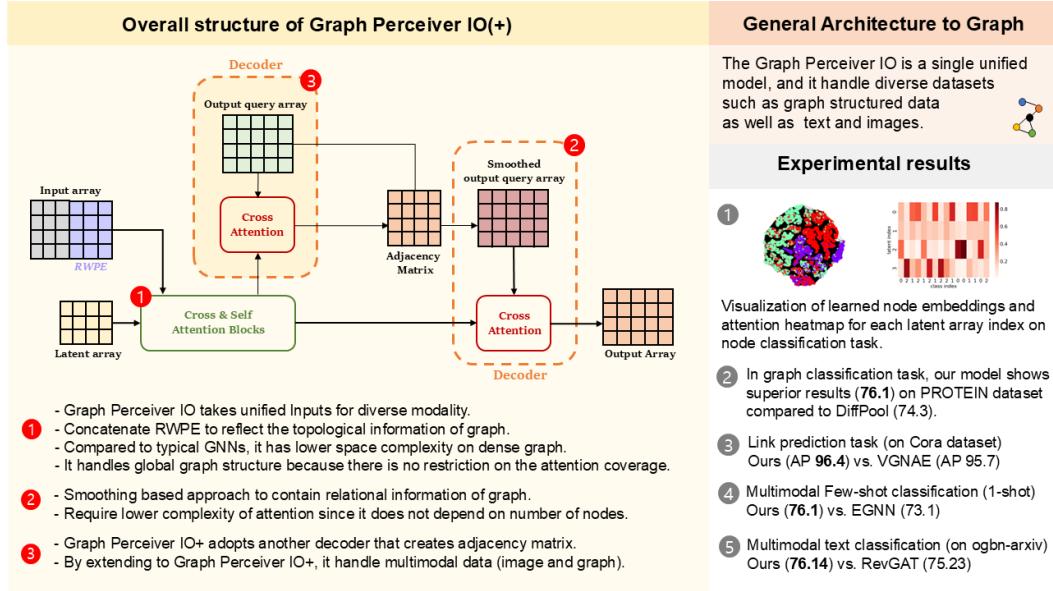
School of Computer Science and Engineering, Chung-Ang University, 84, Heukseok-ro, Dongjak-gu, Seoul, Republic of Korea

Email: yoonsik@cau.ac.kr

Graphical Abstract

Graph Perceiver IO: A General Architecture for Graph-Structured Data

Seyun Bae, Hoyoon Byun, Changdae Oh, Yoon-Sik Cho¹, Kyungwoo Song²



¹Corresponding author

²Corresponding author; Work partly done at University of Seoul

Highlights

Graph Perceiver IO: A General Architecture for Graph-Structured Data

Seyun Bae, Hoyoon Byun, Changdae Oh, Yoon-Sik Cho³, Kyungwoo Song⁴

- Graph Perceiver IO generalizes a multimodal method to the graph domain by designing specific graph features to reflect the topological information.
- Compared to the graph neural networks, Graph Perceiver IO requires a lower complexity, and it can handle the global and local information efficiently.
- By extending to two separated decoders, Graph Perceiver IO+ incorporates both images and graphs simultaneously for the multimodal few-shot classification.
- The Graph Perceiver IO(+) is the single unified model that takes input and output flexibly, and it handles diverse datasets such as graph structured data as well as text and images.

³Corresponding author

⁴Corresponding author; Work partly done at University of Seoul

Graph Perceiver IO: A General Architecture for Graph-Structured Data

Seyun Bae^a, Hoyoon Byun^b, Changdae Oh^b, Yoon-Sik Cho^{1c}, Kyungwoo Song^{2d,e}

^a*Department of Computer Science and Engineering University of Seoul*

^b*Department of Artificial Intelligence University of Seoul*

^c*School of Computer Science and Engineering Chung-Ang University*

^d*Department of Applied Statistics Yonsei University*

^e*Department of Statistics and Data Science Yonsei University*

Abstract

Multimodal machine learning has been widely studied for the development of general intelligence. Recently, the Perceiver and Perceiver IO, show competitive results for diverse dataset domains and tasks. However, recent works, Perceiver and Perceiver IO, have focused on heterogeneous modalities, including image, text, and there are few research works for graph structured datasets. A graph has an adjacency matrix different from other datasets such as text and image, and it is not trivial to handle the topological information. In this study, we provide a Graph Perceiver IO (GPIO), the Perceiver IO for the graph structured dataset. We keep the main structure of the GPIO as the Perceiver IO because the Perceiver IO already handles the diverse dataset well, except for the graph structured dataset. The GPIO is a general method that handles diverse datasets, such as graph-structured data, text, and images, by leveraging positional encoding and output query smooth-

¹Corresponding author

²Corresponding author; Work partly done at University of Seoul

ing. Compared to graph neural networks (GNNs), GPIO requires lower complexity and can efficiently incorporate global and local information, which is also empirically validated through experiments. Furthermore, we propose GPIO+ for the multimodal few-shot classification that incorporates both images and graphs simultaneously. GPIO achieves higher benchmark accuracy than GNNs across multiple tasks, including graph classification, node classification, and multimodal text classification, while also attaining superior AP and AUC in link prediction. Additionally, GPIO+ outperforms GNNs in multimodal few-shot classification. Our GPIO(+) can serve as a general architecture for handling various modalities and tasks.

Keywords: Graph Perceiver IO, Graph Neural Network, Multimodality

1. Introduction

Humans who have general intelligence has the capability to handle multiple datasets from different source simultaneously. One of the required building blocks for general intelligence is multimodal learning for a general perception, and there has been much interest in multimodal learning. With the development of the Transformer [1] and the Vision Transformer (ViT) [2], there have been multimodal learning approaches that adopts multiple Transformer encoders for each modality [3], and fuses the each information later [4]. Another remarkable algorithms for multimodal learning are the Perceiver [5] and its extension, the Perceiver IO [6], general architecture that has structured inputs and outputs. Perceiver IO handles various sizes of inputs and outputs, and they unify the inputs and outputs structures for diverse types of datasets such as images and text. Unified inputs and outputs structures

make the single model, Perceiver IO, handle multimodal datasets simultaneously. With Perceiver IO, the specific neural networks for each modality are not required.

Parallel to multimodal learning, there have been many works for the graph structured dataset. A graph has topological information and a graph is known to be a general type that can represent diverse types of the dataset, including image and text [7]. From its general property, the graph has diverse application areas such as recommender systems [8], drug discovery [9], knowledge graph [10] and pedestrian trajectory prediction [11]. The graph has meaningful knowledge information, and the intelligence learned from a graph structured dataset might be important for general intelligence as well as other datasets. The topological information can be divided by relational information and canonical positional information. Especially, the relational information, described as an adjacency matrix, is a distinct features different from the other dataset. The Perceiver and Perceiver IO flatten all inputs, and they adopt the positional encodings to describe the spatial information or the sequence order. However, a graph has no such sequence order but explicit relationships, requiring representation distinct from others. Besides, the proper input and output query array for multiple graph-related tasks and multimodal data, including text and images, has not been explored yet.

In this study, we propose the Graph Perceiver IO (GPIO), a general method to handle the graph structured dataset as well as diverse modality. The GPIO takes the raw node features and graph positional encoding to reflect node features and to distinguish various shape of graphs, canonical positional information. Besides, The GPIO can handle relational information

of graph through smoothing based output query array design. For the GPIO, we propose proper input array and output query array design to enhance the graph task performance while keeping the original model structure of Perceiver IO to preserve the multimodal ablity. Also, we propose Graph Perceiver IO+ to enhance the multimodal few-shot learning tasks on image and graph data.

The GPIO has three advantages over the traditional GNNs. First, the GPIO has a lower space complexity. For the given input array, latent array, and output query array, the computation of Perceiver IO does not depend on an adjacency matrix. While the complexity of traditional GNNs is proportional to the squared number of nodes, the GPIO has linear complexity. Second, the GPIO handles information on both the global and local structure. The GPIO does not have a restriction on the attention areas, and all nodes can interact with attention. It is noted that the locality inductive bias can be injected into the GPIO with the positional encoding and output query array. Third, the GPIO can handle multimodal data, including the graph, simultaneously. The single unified model, the GPIO, is necessary to consider the multimodal data such as image and graph, while the traditional methods need two networks for image and graph, respectively. Our research can contribute to the diagnosis of human diseases. MLP have been studied for diagnosing Alzheimer’s Disease (AD) [12] by classifying fMRI images and RNN-LSTM for early diagnosis of Diabetic retinopathy (DR) [13]. Following this, our multimodal research, including graphs, can be applied to diagnose AD [14] and DR [15]. For all experiments, here are the URLs of the source code you can reference : <https://github.com/MLAI-Yonsei/GPIO>

2. Related Works

2.1. Multimodal representation learning

Based on a tremendous amount of image-text pair datasets, contrastive pre-trained multimodal models such as CLIP [3], BLIP-2 [16] and Alpha-CLIP [17] recently show impressive progress w.r.t generalizability of representation, zero-shot transfer and OOD [18]. A key factor to their great success is the language-guided visual representation learning scheme. Given a set of image-caption pairs $(x_1, c_1), \dots, (x_N, c_N)$, they train an image encoder $E_{img}(\cdot)$ and text encoder $E_{txt}(\cdot)$ such that the similarity $\langle E_{img}(x_i), E_{txt}(c_i) \rangle$ between positively aligned pair is maximized relative to negatively unaligned pair through InfoNCE [19] objective function. In addition to image-text learning, video-text learning for segmentation [20] and video-audio-text sentiment analysis [21] and image-speech-text learning [22] have also studied actively.

2.2. General architecture

In another line of work, there are attempts that aim to unify the architecture or learning objective for training in different modalities. Most of these studies mainly utilize the Transformer [1] architecture which has less data-specific and task-specific inductive bias as the backbone. Instead of relying on modality-specific targets, Data2vec [23] uses the same learning objective for data from different modalities. Based on the combination objective of latent target prediction task [24] and Masked prediction [25], Data2vec

shows promising results on three different modalities including text, image and speech. [5] propose a general framework Perceiver to handle arbitrary configurations of different modalities. Without making domain-specific assumptions, Perceiver encodes the data point through iterative self-attention and cross-attention. And they mitigate the quadratic scaling problem of self-attention blocks in Transformer by processing a small set of latent units instead of high-dimensional inputs. Recently, Perceiver IO [6] appeared to extend the capable task that the original Perceiver can not perform, and Hierarchical Perceiver [26] improves the performance and efficiency of the Perceiver by introducing locality into the model architecture while preserving its modality-independence property. Our proposed model, GPIO, further extends the accommodatable modality of Perceiver IO to the graph structured dataset.

2.3. Graph Neural Networks

The representative method of Graph Neural Networks is GCN [27] which provides spatial domain graph convolution by approximating a spectral graph convolution. GCN propagates a node representation by using a normalized adjacency matrix. APPNP proposes propagation method based on personalized PageRank, which is one of the GCN variants. The memory efficient RevGAT [28] for the deeper and wider models, Dir-GNN [29] which adopts directionality for effective homophily and HiGCN [30] for the higher-order interaction have improved classification ability on node-level or graph-level tasks. Furthermore, various modified models of GNNs in link prediction tasks ensure reliable performance. [31] predict a link using an inner product between the latent embeddings of target nodes encoded based on a varia-

tional autoencoder [32]. [33] propose VGNAE that applies L_2 -normalization before propagation to alleviate the issue that an isolated node embed got close to zero by VGAE [31]. There are recent works of TokenGT [34] that applies the transformer to the graph-structured dataset. In addition, TAPE [35] that utilizes the vast knowledge of LLM to enrich text information and use it as a GNN feature, and GPT4Graph [36] that applies LLM itself to various graph-related tasks.

3. Preliminary

3.1. Perceiver

Jaegle et al. [5] proposes the Perceiver, a general perception module with iterative attention for handling the diverse and high-dimensional multimodal inputs. The Perceiver has two components of arrays, input arrays x , and latent arrays z . For given data instances, the Perceiver transforms the given data instances as the inputs arrays $x \in \mathbb{R}^{M \times C}$ where M and C denote the dataset properties such as the image size and the number of channels. The latent arrays $z \in \mathbb{R}^{N \times D}$ are learnable latent representations, where N and D are the number of latent and the latent dimension, respectively. The Perceiver adopts the two types of attention, cross-attention and self-attention, and both attentions are query-key-value attention, similar to the Transformer [1]. First, the cross-attention computes the attention between the latent arrays and the input arrays. The latent arrays are queries, and they attend to the important information from the key, the input arrays. The relevant input features for the given tasks are incorporated into the latent array with cross-attention. Second, the Perceiver adopts the self-attention for the latent

arrays. Attention between latent arrays and input arrays reduces the space and time complexity compared to the self-attention for the input arrays only. The attention complexity of the Perceiver is $O(NM)$ while the complexity of the self-attention for input-image is $O(N^2)$. In summary, the Perceiver is composed of the multiple blocks of cross-attention and the latent Transformer that adopts the self-attention for the latent arrays. The input arrays are fed into the Transformer iteratively across the layers. After the multiple blocks, the Perceiver adopts the average pooling over the index dimension and utilizes it to predict the target label.

3.2. Perceiver IO

The Perceiver is a general model that handles diverse datasets with flexible input structures. However, the output structure of the Perceiver is not flexible, and the Perceiver is limited to the classification tasks. To extend the Perceiver into diverse tasks such as language modeling and autoencoding, flexible output structures are necessary. Perceiver IO proposes an additional output query array components $y \in \mathbb{R}^{M \times E}$ for the flexible output structures. For classification task on image, the M denotes the number of images and E denotes the number of classes, respectively. Perceiver IO adds the additional query-key-value attention block, i.e. decoder, on the top of the Perceiver. The additional attention block computes the attention between the output query array and the latent arrays extracted by the Perceiver. The shape of the output query array is a controllable hyper-parameter, and it induces flexible output structures that produce any size of outputs.

The Perceiver and Perceiver IO attend to the relationship between the latent arrays and input arrays. They adopt query-key-value attention that

is similar to the Transformer. Different from the convolutional neural nets (CNN) [37] or the Vision Transformer (ViT) [38], the Perceiver operates on the individual pixels independently without patches or convolution operations. However, attention computation is permutation invariant, and the attention-only mechanism is hard to treat the sequence order. To incorporate the order information for sequence and spatial information for images, the Perceiver and the Perceiver IO introduce the learned positional encoding or Fourier-based positional encoding with sinusoid functions.

3.3. Perceiver IO and Graph Structured Data

The Perceiver IO is a general architecture, and it has high flexibility for input and output structure. From its flexibility and generality, The Perceiver and its variants handle the diverse tasks for image, text, speech, and 3D point clouds. They extend their works to incorporate the heterogeneous multimodal dataset such as audiovisual sequences, and the Perceiver shows competitive results. However, the application of Perceiver is limited to the non-graph structured dataset, and they do not handle the graph structured dataset that has topological information.

Extending the Perceiver to handle the graph structured dataset is non-trivial. First, the Perceiver is a general architecture, and it already shows the competitive results for diverse tasks. Therefore, simple extensions that do not modify the overall model structures are required. Second, the Perceiver receives the flattened inputs, and incorporating the topological information represented by the adjacency matrix into the Perceiver is challenging.

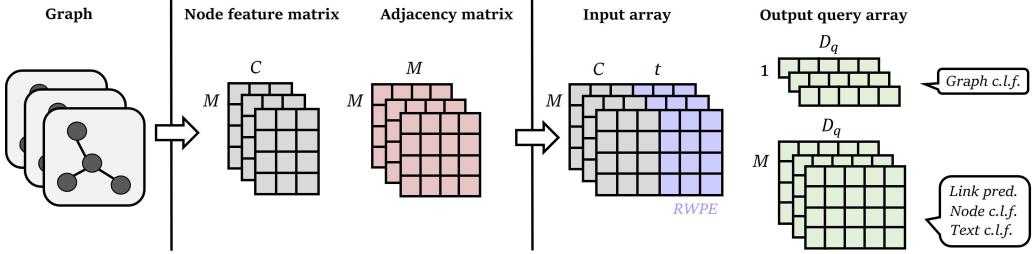


Figure 1: A graph-structured dataset splits into a node feature matrix and an adjacency matrix. GPIO construct the input array and the output query array to integrates graph data. Note that the node feature matrix can be text and image feature in multimodal learning. (i) Input array: Concatenation of node feature and random walk positional embedding (*RWPE*), where t denotes the dimension of *RWPE*. (ii) Output query array: A random initialized D_q dimensional vector used for the graph classification task, and (M, D_q) for node classification, link prediction and multimodal text classification. The output query array can be learnable or fixed, and *c.l.f.* denotes the classifier.

4. Methodology

4.1. Overall Structure

Graph data usually consists of $X \in \mathbb{R}^{M \times C}$ matrix, a set of nodes features, and adjacency matrix $A \in \mathbb{R}^{M \times M}$, which represents a set of relational information between nodes. In order to train graphs using Perceiver IO, the input array and the output query array should be constructed by using X and A properly. To handle the graph-structured data, topological information is essential. We can divide the topological information on the graph as relational information and canonical positional information. Both two information is essential for the graph-related tasks requiring the use of adjacency matrix. We propose the methods (i) output query array smoothing to incorporate relational information and (ii) positional embeddings to incorporate canonical

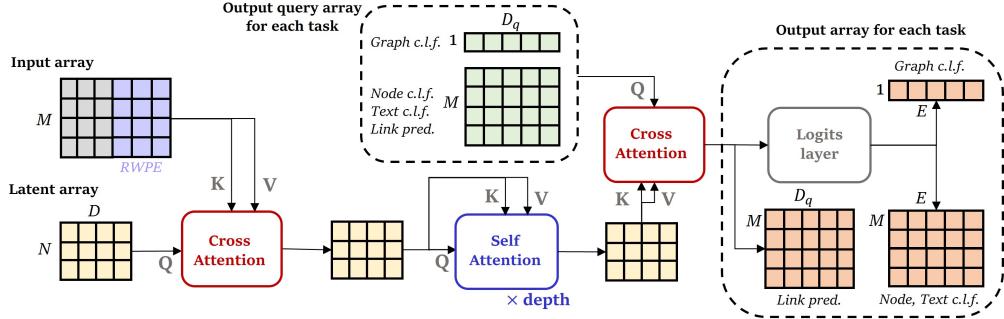


Figure 2: Overall structure of GPIO. First, the cross-attention between the initial latent and input query enables the latent to absorb the necessary information from the input. Then, the latent progressively encode the salient feature of a given data point through repeated self-attention blocks. Finally, the output query array and final latent communicate via cross-attention to make proper output for each task. D_q is an arbitrarily configurable output query array dimension, and E is the number of classes about a given task. The number of depth or layer of self attention block is a controllable hyperparameters.

positional information. As shown in Figure 1, the methods are implemented by concatenating the node features and positional embeddings and smoothing the node features. The shape of the output query array depends on the given task, and we will discuss it in Section *Output Query Array*. The overall structure of our proposed model is in Figure 2.

4.2. Output Query Array

One of the major problems of the GPIO is to incorporate the adjacency matrix. The adjacency matrix that is widely utilized in GNNs is not adaptable for the GPIO. Also, the GPIO requires a flexible output structure to handle the diverse graph-related tasks such as link prediction, graph classification, and node classification. For the flexible output array structure, we utilize the output query array of Perceiver IO. The output query shape can

be determined by the task specifically. It is noted that the output query array can be used as a pretrained embeddings, or it can be learnable parameters. For graph classification task that requires a single label per graph, we set the output query array size as $1 \times D_q$, where D_q denotes the output query array dimension. For node-specific tasks such as node classification and link prediction, we set the output query array shape as $M \times D_q$ where M is the number of nodes. The simple choice of the output query array is raw nodes features $X \in \mathbb{R}^{M \times C}$ that is the same as the input array. However, we find that the output query array based on the raw node features tends to show poor performance. Appendix 8.2. presents a detailed result. Furthermore, raw node features do not contain any relational information of adjacency matrix. Instead of raw node features, we propose a smoothing based output query array that is smoothed features of raw node features L times [39]. The following Eq. 1 denotes the L times smoothed node features, the output query array of the GPIO,

$$\hat{X} = \left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} \right)^L X \quad (1)$$

where $\tilde{A} = A + I$ denotes a adjacency matrix with self-loop and $\tilde{D}_{ii} = \sum_{j=0} \tilde{A}_{ij}$ denotes node degree matrix of \tilde{A} . SGC [39] removes nonlinear function between GCN layers and collapse repeated multiplication of matrix into L -th power of matrix like Eq. 1 which extracts smoothed features of X . As a result, extracted features \hat{X} tends to have similar representation to nearby nodes and help to make similar predictions. That is, the SGC incorporates adjacency matrix into its input to get L -hops neighbor information. In a similar way, we utilize adjacency matrix to obtain \hat{X} as an output query array which contain relational information.

Besides, we adopt the APPNP smoothing [40] as an output query array as shown in Eq. 2.

$$\begin{aligned} X^{(0)} &= X, \\ X^{(l)} &= (1 - \alpha)\tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}X^{(l-1)} + \alpha X^{(0)}, \\ \hat{X}_{APPNP} &= X^{(L)} \end{aligned} \quad (2)$$

APPNP adopts personalized propagation, and it provides more data specific smoothing methods. We empirically find that the output query array with a smoothing method is one of the key factors for the GPIO. Section 4.5. presents complexity analysis of output query array.

4.3. Input Array

The GPIO requires the flexible input structure as the Perceiver or Perceiver IO. For image and text dataset, *jaegle et al.* propose Fourier embedding to consider the spatial or sequential information, but it shows poor performance for graph classification [6]. Appendix 8.2.3. presents a detailed result. To reflect the topological information, specifically canonical positional information, with the original Perceiver input array structure, we adopt the random walk positional embedding (*RWPE*) [41]. Random walk operators are defined as $R = AD^{-1}$ where $A \in \mathbb{R}^{M \times M}$ is the adjacency matrix, and D is the degree matrix. *RWPE* of node i is defined as

$$RWPE_i = [R_{ii}, R_{ii}^2, \dots, R_{ii}^t] \in \mathbb{R}^t, \quad (3)$$

where R_{ii}^t is the probability of starting from the i node, walking t times, and returning to the i node. *dwivedi et al.* adopt *RWPE* as positional encoding initialization for distinguishing several cases of non-isomorphic graphs and

structurally different nodes which message-passing GNNs and 1-WL [42] fail to distinguish. The choice of sufficiently large t gives a unique node representation to distinguish each node and graph. Besides the node features, we concatenate the $RWPE$ to obtain nodes representation distinguishing non-isomorphic graphs and structurally different nodes. There is no restriction on the attention coverage, and it makes the GPIO handles the global structure information efficiently. It is noted that the injection of $RWPE$ into the GPIO can be interpreted as an inductive bias to consider the local information as well as the global information. Compared to typical GNNs, the Perceiver has lower space complexity in terms of the input array. GNNs requires $O(M \times C + M \times M)$ on dense graph, while the GPIO needs $O(M(C + t))$, where t is the controllable hyper-parameters. Consequently, GPIO is not proportional to the number of edges, making it more efficient on dense graphs. Experiments on Barabási-Albert (BA) graphs of varying densities confirm that GNN memory usage grows with density, whereas GPIO depends only on the number of nodes M . This highlights GPIO’s space efficiency for large, dense graphs. The specific experimental analysis is in the Appendix 8.5.2.

4.4. Model Analysis

We demonstrate how the GPIO leverages graph structures via adjacency matrix A . The decoder in the GPIO computes the attention scores between latent array z and output query array X . The attention score matrix $\mathbf{S} \in \mathbb{R}^{M \times N}$ can be represented as follows, where W_Q and W_K denote query and key projection matrix with hidden size d .

$$\mathbf{S} = \frac{(XW_Q)(zW_K)^\top}{\sqrt{d}} \quad (4)$$

Each element in \mathbf{S} is a score indicating how much each latent is focusing on which node features. This means that z attends globally, regardless of whether the node feature is connected to any of its neighbors. By taking feature smoothing \widehat{X} , the GPIO exploits A to propagate the global score to its local neighbors. The propagated score matrix $\widehat{\mathbf{S}}$ can be derived as follows, where $\ddot{A} = \widetilde{D}^{-\frac{1}{2}} \widetilde{A} \widetilde{D}^{-\frac{1}{2}}$.

$$\begin{aligned}\widehat{\mathbf{S}} &= \frac{(\widehat{X}W_Q)(zW_K)^\top}{\sqrt{d}} = \frac{(\ddot{A}^L X W_Q)(zW_K)^\top}{\sqrt{d}} \\ &= \ddot{A}^L \mathbf{S},\end{aligned}\tag{5}$$

The relational information \ddot{A} is responsible for propagating the weighted score to its neighbor that are L -hops away. The global focus of a latent spreads its impact across local neighbor nodes, ensuring that nodes at least L -hops away have similar scores. Detailed experimental results are covered in *Discussion*.

An alternative is APPNP smoothing, which prevents oversmoothing compared to Eq. 5 to propagate to more distant neighbors.

$$\widehat{\mathbf{S}}_{APPNP} = (\alpha \sum_{l=0}^{L-1} (1-\alpha)^l \ddot{A} + (1-\alpha)^L \ddot{A}^L) \mathbf{S}\tag{6}$$

The global attention score that a latent assigns to a particular node is propagated to nodes within an L -hops with different weights based on \ddot{A} , allowing the GPIO to handle graph structures.

4.5. Complexity Analysis

The GPIO requires a single adjacency matrix-related operation to make the smoothed output query array before the model training. The smoothing

operation is only once conducted, and it never appears during the model training. The mechanism is similar to the Simplifying Graph Convolutional (SGC) [39] that only performs single adjacency matrix-related operation. Therefore, the GPIO has an efficient computation time compared to the other GNNs. Besides, there is no attention between inputs, but there is attention between the latent array and the input array, or the output query array and the input array. The complexity of attention between latent arrays is independent of the number of nodes, and the complexity of attention between latent and input arrays depends on the number of nodes linearly. However, the current Transformer variants, such as Graph Transformer architecture [43] adopt the input-wise attention, and the computation time is proportional to the squared number of nodes. Therefore, GPIO has advantages in the time complexity as well as its simplicity and multimodal learning.

4.6. Graph Perceiver IO+

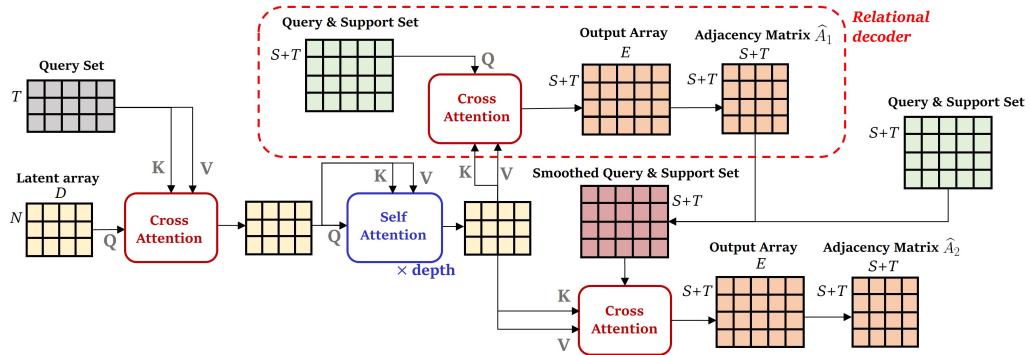


Figure 3: Overall structure of GPIO+. The red dashed area is the relational decoder introduced for the extension to GPIO+. \hat{A}_1 gives rich information to classify \mathcal{Q} . The second decoder performs multimodal few-shot image classification.

The few-shot classification task is to classify query set \mathcal{Q} which contains unlabeled T samples, when support set \mathcal{S} which contains labeled S samples is given. n -way k -shot classification problem means that support set \mathcal{S} contains number of n classes and each class has number of k labeled samples. Constructing graph-structures between \mathcal{Q} and \mathcal{S} provides rich relational information to classify which class a sample from \mathcal{Q} belongs to among the \mathcal{S} . To enhance the multimodal few-shot task performance, we additionally propose Graph Perceiver IO+ (GPIO+), that handle both image features and relational information of each set by adopting *Relational decoder*. As seen in Figure 3, the *Relational decoder* is extension of the original GPIO’s decoder that decodes relational information between \mathcal{Q} and \mathcal{S} [44]. We configure input array as samples of \mathcal{Q} and output query array as samples of \mathcal{Q} and \mathcal{S} in the both decoders. The *Relational decoder* creates adjacency matrix $\widehat{A}_1 \in \mathbb{R}^{(S+T) \times (S+T)}$, and the second decoder that is original decoder of GPIO creates $\widehat{A}_2 \in \mathbb{R}^{(S+T) \times (S+T)}$. The *Relational decoder* takes the unsmoothed samples of \mathcal{Q} and \mathcal{S} as an output query array, and creates adjacency matrix \widehat{A}_1 which is obtained by inner product of output array itself. \widehat{A}_1 represents relational information of each pair of \mathcal{Q} and \mathcal{S} . The second decoder takes a smoothed output query array which is smoothed by \widehat{A}_1 , and it generates an adjacency matrix \widehat{A}_2 . \widehat{A}_2 is final output to classify \mathcal{Q} . The edge label matrix $Y_e \in \mathbb{R}^{(S+T) \times (S+T)}$ with a value of 1 if the class is the same and 0 if the class is different for all query set and support set pairs. GPIO+ follows episodic training, which mimics few-shot learning setting of test time, and the each label in train and test time is sampled from mutually exclusive set of all classes. In each episodic, the final loss \mathcal{L} is obtained by computing loss

between Y_e and weighted average of \hat{A}_1 and \hat{A}_2 such that

$$\hat{A} = \beta\hat{A}_1 + (1 - \beta)\hat{A}_2, \text{ and } \mathcal{L} = \mathcal{L}_{bce}(Y_e, \hat{A}) \quad (7)$$

where β is controllable hyperparameter and \mathcal{L}_{bce} is binary cross entropy loss. The final class of query set selects one of the class of all support set with the highest edge predicted value of \hat{A} [44].

5. Results

To validate our proposed model, the GPIO, we conduct comprehensive experiments, node classification, graph classification, link prediction, multi-modal few-shot image classification, and multimodal text classification. Appendix 8.1. provides the more detailed experimental settings.

5.1. Link Prediction

| Model | Cora | | CiteSeer | | PubMed | |
|----------------------|-----------------|-----------------|-----------------|-----------------|------------------|------------------|
| | AUC | AP | AUC | AP | AUC | AP |
| Spectral Clustering* | 84.6±0.01 | 88.5±0.0 | 80.5±0.01 | 85.0±0.01 | 84.2±0.02 | 87.8±0.01 |
| DeepWalk* | 83.1±0.01 | 85.0±0.0 | 80.5±0.02 | 83.6±0.01 | 84.4±0.0 | 84.1±0.0 |
| DGI* | 89.8±0.8 | 89.7±1.0 | 95.5±1.0 | 95.7±1.0 | 91.2±0.6 | 92.2±0.5 |
| ARGVA* | 92.4±0.004 | 93.2±0.003 | 92.4±0.003 | 93.0±0.003 | 96.8±0.001 | 97.1±0.001 |
| GIC [45] | 90.0±1.0 | 89.9±1.3 | 95.8±0.7 | 95.8±0.9 | 90.9±1.0 | 91.6±0.9 |
| VGAE [31] | 95.2±0.5 | 94.7±0.6 | 92.0±1.7 | 91.6±1.7 | 95.6±0.7 | 95.3±0.6 |
| GNNAE [33] | 95.6±0.7 | <u>96.0±0.8</u> | 97.2±0.5 | 97.3±0.4 | 97.7±0.2 | <u>97.6±0.2</u> |
| VGNNAE [33] | <u>95.8±0.6</u> | 95.7±0.8 | 96.8±0.6 | 96.7±0.6 | 97.3±0.1 | 97.2±0.2 |
| GPIO | 95.9±0.3 | 96.4±0.2 | <u>96.8±0.3</u> | <u>97.2±0.2</u> | <u>97.6±0.08</u> | 97.8±0.07 |

Table 1: Performance of link predictions. GPIO shows the competitive results with the well-known link prediction model VGAE, which adopts the graph convolutional neural networks. * denotes the reported performance in [45].

Besides the classification task, we validate our models on the link prediction task. Table 1 denotes the results of the link prediction task. We utilize the features extracted by the cross attention between the output query array and latent array for the link prediction task. The GPIO shows the competitive performance across all dataset. We provide detailed experimental settings and t-SNE plot to compare the performance of the three models in Appendix 8.2.1.

5.2. Node Classification

| Model | Cora | | CiteSeer | | PubMed | |
|-------|------------------|------------------|------------------|------------------|------------------|------------------|
| | Fixed | Random | Fixed | Random | Fixed | Random |
| GCN | 81.4±0.7 | 78.7±1.7 | 71.1±0.7 | 68.1±1.7 | 78.9±0.6 | 77.3±2.4 |
| GAT | 83.0±0.5 | 80.9±1.5 | 70.9±0.5 | <u>68.8</u> ±1.7 | 78.9±0.4 | 77.6±2.4 |
| Cheb | 80.5±1.0 | 77.0±2.7 | 69.8±1.2 | 67.2±2.1 | 78.2±0.6 | 75.4±2.5 |
| SGC | 81.7±0.1 | 80.1±1.9 | <u>71.3</u> ±0.2 | 68.5±2.0 | 78.9±0.0 | 76.6±2.4 |
| ARMA | 82.2±0.9 | 79.8±1.7 | 71.0±0.6 | 67.9±1.9 | 78.8±0.3 | 77.6±2.2 |
| APPNP | <u>83.3</u> ±0.5 | 82.1 ±1.5 | 71.7 ±0.5 | 69.8 ±1.8 | 80.1 ±0.2 | <u>79.1</u> ±2.3 |
| GPIO | 83.9 ±0.6 | 81.6±1.8 | 70.1±1.0 | 68.1±1.7 | 79.9±0.4 | 79.6 ±2.2 |

Table 2: Accuracy for the node classification. The GPIO shows the competitive results on PubMed. Fixed and Random denotes the dataset split methods from [46].

For node classification experiments, we adopt the three benchmark data, Cora [49], CiteSeer [50], and PubMed [51]. We follow all experimental settings and GNN baseline from [46], and we repeat the experiment 100 times. Table 2 denotes the accuracy of node classification task. The GPIO shows the competitive results on the Cora and PubMed dataset.

Besides, for the PubMed dataset, which has 3 classes and 500-dimensional features, we extracted high-dimensional embeddings from the penultimate

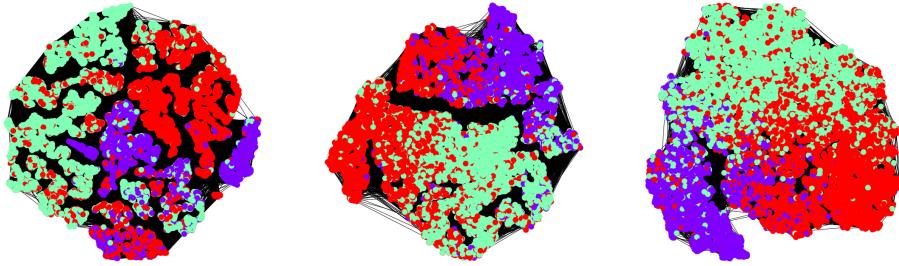


Figure 4: t-SNE [47] visualization of learned nodes embedded by the GPIO (left), APPNP (middle), GAT (right) on PubMed dataset. The learned embedding of GPIO shows a relatively large uniformity compared to the APPNP and GAT. Large uniformity denotes that a feature distribution utilizes maximal information, and it has a positive correlation with the downstream task performance [48].

layer of each model and visualized their 2D projection in Figure 4. Each color represents a class. The embedding visualization confirms the GPIO has a relatively high uniformity embedding space.

5.3. Graph Classification

We validate our models on the graph classification task. Graph classification tasks require the understanding of both the global and local graph structures. We repeat the experiments ten times, and we follow the same experimental protocols with the [46]. Table 3 denotes the accuracy of the graph classification on diverse benchmark dataset. For MUTAG [52] and PROTEINS [52], GPIO shows superior performance compared to the recent GNNs. In short, GPIO has a high capacity to incorporate both the node features and link information simultaneously. It is noted that MUTAG and PROTEINS have node features while IMDB, REDDIT, and COLLAB [53] have no node features. Even though there are no node features, the GPIO

| Model | MUTAG | PROTEINS | IMDB-BINARY | REDDIT-BINARY | COLLAB |
|------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| GCNWithJK | 72.9±12.0 | 72.6±3.6 | 73.2±5.0 | 89.4±2.9 | 81.5±2.1 |
| SAGPool | 74.0±8.7 | 72.3±2.8 | 72.3±4.7 | 89.0±2.1 | 78.9±1.0 |
| DiffPool | <u>84.6±8.7</u> | <u>74.3±6.4</u> | 74.8±4.8 | 92.7±2.0 | 79.4±1.9 |
| GCN | 70.7±11.0 | 72.2±2.4 | 74.2±4.4 | 89.1±2.0 | <u>81.0±1.4</u> |
| GraphSAGE | 75.1±11.4 | 74.1±2.3 | 73.2±4.4 | 90.7±2.3 | 79.8±1.1 |
| GIN0 | 81.9±8.0 | 73.1±3.8 | <u>73.7±4.1</u> | 90.9±2.1 | 80.5±1.9 |
| Set2SetNet | 74.5±11.9 | 74.1±3.8 | 72.7±5.0 | 90.3±2.4 | 79.4±1.7 |
| SortPool | 83.0±9.0 | 73.9±4.5 | 71.8±3.0 | 84.3±5.0 | 77.8±1.6 |
| ASAP | 78.7±11.8 | 74.0±3.0 | 72.2±4.3 | OOM | 79.4±1.7 |
| GPIO | 86.1±6.9 | 76.1±3.0 | 72.9±4.4 | <u>91.6±2.3</u> | 79.1±2.5 |

Table 3: Accuracy of the graph classification tasks. GPIO shows the superior performance on MUTAG and PROTEINS dataset that has node features. The topological information is significant because IMDB, REDDIT, and COLLAB have no node features. GPIO shows the competitive results even though there are no node features.

has competitive results with the recent GNNs. It denotes that the GPIO handles the topological information of graph-structured data with the positional embeddings.

5.4. Multimodal Few-shot classification

We perform a few-shot image classification task to evaluate GPIO+ for image-graph multimodal learning. We follow public source code of [44] as an experimental settings. Also, we employ pretrained VGG16 as feature extractor, and change the number of training and test iterations to 10K and 1K, respectively. Table 4 shows the results on 5-way setting on miniImageNet dataset. GPIO+ shows superior performance, compared to EGNN in 1-shot and 3-shot settings. The *relational decoder* of GPIO+ takes raw image

| Model | 1-shot | 3-shot | 5-shot |
|-------------------|-------------|-------------|-------------|
| MAML+Transduction | 50.8* | - | 66.2** |
| TPN | 53.8* | - | 69.4** |
| EGNN (original) | 59.2 | 71.2 | 75.7 |
| EGNN† | <u>73.1</u> | <u>83.5</u> | 87.9 |
| GPIO+† | 76.1 | 83.7 | <u>86.9</u> |

Table 4: The 5-way, 1-shot, 3-shot and 5-shot few-shot classification accuracies on miniImageNet dataset. * and ** denotes the reported performance in [54] and [44], respectively. † denotes result to use pre-trained feature extractor.

features as the output query array to generate the relational information of the images (nodes). The graph structure of images (nodes) is generated in this step. The second decoder takes smoothed image features as output query array to incorporate graph structure of nodes. Using two separated decoders, GPIO+ handle multimodal data such as image and graph.

5.5. Multimodal Text Classification

We conduct evaluation on OGB datasets [55] which contains large graphs for text-graph multimodal learning. We follow [35] as GNNs baseline and experimental settings and employ fine-tuned language model as text feature

| Model | ogbn-products | ogbn-arxiv |
|--------|---------------------------------------|---------------------------------------|
| MLP | 0.7429 ± 0.0075 | 0.7265 ± 0.0014 |
| GCN | 0.7561 ± 0.0202 | 0.7409 ± 0.0035 |
| SAGE | 0.7428 ± 0.0115 | 0.7355 ± 0.0037 |
| RevGAT | 0.7663 ± 0.0064 | 0.7523 ± 0.0025 |
| GPIO | 0.7763 ± 0.0007 | 0.7614 ± 0.0006 |

Table 5: Text classification accuracy for the ogbn-products and ogbn-arxiv datasets [55]. For ogbn-products, we used the same subset data as in [35].

extractor which is size 129M [35]. The nodes of the graph consist of the title, summary, or product description of the article, which are composed of the citation network or the concurrent purchase network, respectively. The model should classify these text nodes into the correct class using topological information. Dataset statistics and are in the Appendix 8.4. In table 5, GPIO outperforms GNNs on all datasets and handles large structured graph datasets effectively. The result illustrates that our model performs well in text classification by utilizing the text and its topological information.

6. Discussion

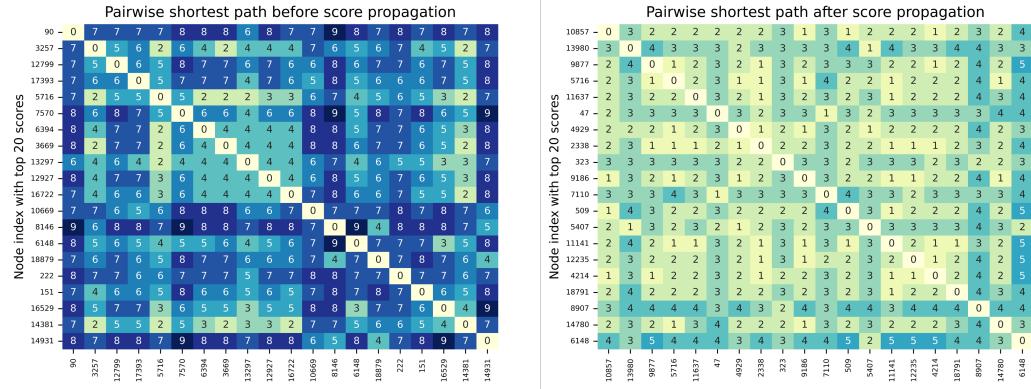


Figure 5: The pairwise shortest path heatmap on PubMed [51] dataset. The left side represents before 2-hops away score propagation and the right side represents after.

We explore how GPIO propagates global information locally with feature smoothing through a visualization. We take one latent that has absorbed the required global information from the input and compute the attention scores $\mathbf{s} \in \mathbb{R}^{M \times 1}$ with the output query array to select the top 20 nodes with the highest attention score out of M nodes. That is, these are the nodes that

the latent has chosen to focus on. We computed the pairwise shortest path of these nodes to see their distance. In Figure 5, the nodes are relatively farther apart before propagation than after. Feature smoothing has changed the global focus to a neighboring local focus. Also, although we used \hat{X} with $L = 2$, some nodes focused by latent still have distances greater than 2. Therefore, GPIO can incorporate the global and local graph structures efficiently.

In graph-related tasks, understanding complex graph structures requires the ability to distinguish between two graphs with different structures. The classic way to tell if two graphs are isomorphic is the 1-WL test [42]: if two graphs are not isomorphic using the 1-WL test, they have different structures. However, the 1-WL test does not guarantee isomorphism. In other words, there are cases where the 1-WL test fails to distinguish between graphs. The typical message-passing GNNs have an expressive power of up to 1-WL. The key to competing with GNN is to be as expressive as or more expressive than 1-WL. A simple choice would be to sequentially give each node a position embedding, like Perceiver IO, but that would be less expressive, as the only thing that distinguishes the graph is the number of nodes. Inspired by [41], we present a concrete analysis of the choice of *RWPE* as a treatment for learning complex graphs compared to GNNs. In order to learn a complex structure like PROTEINS [52], the nodes in the graph must be distinguished by the correct unique positions, so that they are distinct and form a different graph when viewed as a whole. The Figure 6 illustrates the number of uniquely distinct nodes in each graph of PROTEINS when distinguished by 1-WL and *RWPE*. A point on the plot indicates how many of the nodes are

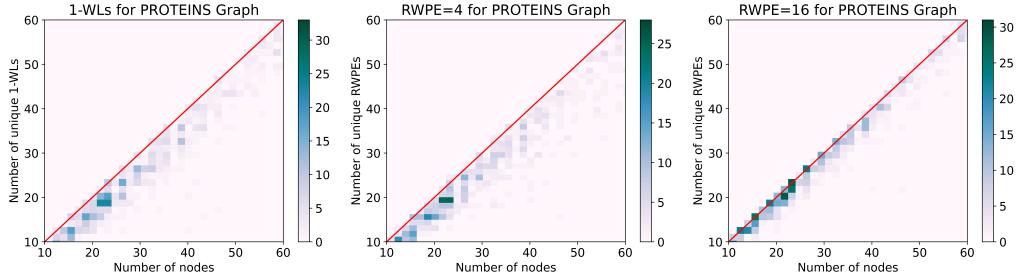


Figure 6: The figure indicates how many nodes are distinguished out of the total number of nodes in the PROTEINS graphs when 1-WL or $RWPE$ is used. The x-axis is the total number of nodes in each graph, and the color on the plot is the number of graphs. Note that having the same number of nodes and unique nodes does not mean that the graphs are isomorphic, only that they can have the same number of (unique) nodes.

uniquely distinguished relative to the total nodes, which means the number of graphs at that point. The more graphs that are close to the red axis, the more uniquely distinguished nodes they have, and the more expressive they are. On the left side of Figure 6, we can see that the graphs distinguished by 1-WL are close to the red axis. Also, the $RWPE$ with $t = 4$ in the middle is also close to the red axis, similar to 1-WL. This means that 1-WL and low-dimensional RWPE have similar expressive power in distinguishing graphs. When $RWPE$ is set to $t = 16$, it is much closer to the red axis. To summarize, we can see that the choice of RWPE is more expressive compared to GNN, which is as expressive as 1-WL. The choice of a suitable dimension of $RWPE$ is more expressive compared to GNNs, which is as expressive as 1-WL. Thus, by utilizing $RWPE$, GPIO has effective expressivity for graph-related tasks.

7. Conclusion

Multimodal learning graph structured dataset, as well as text, image, and speech, is necessary for general perception. In this paper, we propose GPIO(+) that handles the graph structured dataset as well as other datset such as image and text.

The strengths of our work is that we show that graph data, along with other modalities, can be effectively trained on a single architecture without the need for additional inductive bias. This means that GPIO(+) learns topological information from positional encoding and output query design in the same way as other modalities, without changing the structure of the model to specialize in a particular modality. These approaches demonstrates that GPIO is unaffected by the number of graph edges, achieving around 73% memory savings over a typical GNN on a BA graph (edge density 0.03) with around 8000 nodes. For further research, there is potential to extend to the study of canonical positional embeddings on domain specific graphs and various graph-related multimodal learning with output query design. However, our work does not scale well to domains where edge features are crucial, such as social graphs with multiple features like temporal dynamics and relationship types. It is necessary to propose a general architecture by adapting output query smoothing to incorporate edge features.

GPIO achieves higher accuracy on the PROTEIN dataset for graph classification (76.1%) compared to DiffPool (74.3%), and surpasses APPNP on PubMed (79.6% vs. 79.1%) for node classification. It also outperforms VG-NAE on Cora (96.4% vs. 95.7%) in link prediction and improves upon RevGAT on ogbn-arxiv (76.14% vs. 75.23%) for multimodal text classifi-

cation. Finally, GPIO+ achieves 76.1% accuracy in multimodal few-shot classification, exceeding EGNN’s 73.1%.

Overall, GPIO(+) is a general architecture for handling diverse modalities and tasks, including graph-related works, and can be applied to fields like recommendation systems and information retrieval as it accommodates graph, text, and image modalities.

Declaration of Competing Interest

There is no conflict of interest.

Data availability

All of the code and data will be available upon request.

Acknowledgments

This work was supported by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT) (No. 2022R1A4A3033874), and supported by the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT). (No.2021R1F1A1060117).

8. Appendix

8.1. Experimental Settings

In this section, we describe the hyperparameter settings for each task and dataset. We set the same hyperparameter search spaces for each model in order to compare results fairly.

8.1.1. Learning for Node Classification tasks

For node classification tasks, we set the size of the output query array to $(M \times D_q)$. M and D_q are number of nodes and size of the learnable embedding vector respectively. The output array from the last cross attention layer is passed through the logits layer to make the vector size of each node equal to the class size. The size of the output array is, then, $(M \times E)$. E should be the same size with number of classes. We evaluate the cross-entropy loss for node classification tasks.

$$\mathcal{L}_{C.E} = -\frac{1}{M} \sum_{m=1}^M \log\left(\frac{w_{my}}{\sum_{e=1}^E w_{me}}\right) \quad (8)$$

w_{my} and w_{me} denote the value of element of m th node's row vector corresponding to the label index in output array and each element value of m th node's row vector respectively.

8.1.2. Learning for Link Prediction tasks

We use the GPIO as an encoder to embed the raw node features to the latent embeddings and, we adopt an inner product to predict edges. Specifically, the output array generated by our model becomes a latent embedding matrix. In the same manner, as GAE [31], we conduct the inner product between each pair of node features.

$$P(\mathbf{A}|\mathbf{W}) = \prod_{i=1}^N \prod_{j=1}^N P(A_{ij}|w_i, w_j), \quad (9)$$

with $P(A_{ij} = 1|w_i, w_j) = \sigma(w_i^T \cdot w_j)$

Matrix \mathbf{A} denotes the adjacency matrix, and \mathbf{W} is the output array. σ denotes a sigmoid function. We use reconstruction loss for training the model

for link prediction.

$$\begin{aligned}\mathcal{L}_{recon} = & -\log\left(\frac{1}{|N_{tr}^{pos}|} \sum_{(i,j) \in N_{tr}^{pos}} \sigma(w_i^T \cdot w_j)\right) \\ & -\log\left(1 - \frac{1}{|N_{tr}^{neg}|} \sum_{(k,t) \in N_{tr}^{neg}} \sigma(w_k^T \cdot w_t)\right)\end{aligned}\quad (10)$$

N_{tr}^{pos} denotes a set of positive train edges. N_{tr}^{neg} denotes a set of random sampled negative train edges. A σ function maps the inner product of embeddings of node pairs to a proper value between 0 and 1. The reconstruction loss Eq. 10 is binary cross-entropy loss, which assigns positive edges a label of one and negative edges a label of zero.

8.1.3. Learning for Graph Classification tasks

Lastly, for graph classification tasks, the output query array is 1-dimensional array. After through the logits layer, It still be a 1-dimensional array but the size of array should match the number of classes. The characteristic of this procedure is that we can conduct graph-level tasks without the read-out layer (e.g., global pooling layer). We use the cross-entropy loss generally used in graph classifications.

$$\mathcal{L}_{C.E} = -\frac{1}{H} \sum_{h=1}^H \log\left(\frac{w_{hy}}{\sum_{e=1}^E w_{he}}\right) \quad (11)$$

H denotes the number of graphs. w_{hy} and w_{he} denote the value of element of h th graph's row vector corresponding to the label index in the output array and each element value of h th graph's row vector respectively.

8.1.4. Hyperparameter Setting for Node Classification

Table 6 describes the hyperparameter search spaces for the node classification task.

| Hyperparameter | Candidate values | | |
|---------------------|------------------|---------------|--------------------|
| | Cora | CiteSeer | PubMed |
| learning rate | {5e-4, 1e-3} | {1e-3} | {5e-3, 1e-2} |
| weight decay | {1e-2} | {5e-2} | {5e-2, 1e-1} |
| latent length | {16, 32, 64} | {4} | {4, 8} |
| latent dimension | {32, 64} | {64, 128} | {32, 64, 128, 256} |
| # of MHCA heads | {32, 64} | {32, 64} | {1} |
| # of MHSA heads | {8, 16} | {8, 16, 32} | {1, 8, 16} |
| MHCA head dimension | {64, 128} | {32, 64, 128} | {64, 128} |
| MHSA head dimension | {64, 128} | {32, 64, 128} | {64, 128} |
| depth | {1} | {1} | {1} |

Table 6: Hyperparameters for the node classification task. MHCA and MHSA denote the multi-head cross attention and multi-head self attention, respectively.

8.1.5. Hyperparameter Setting for Link Prediction

Table 7 describes the hyperparameter search spaces for the link prediction task.

| Hyperparameter | Candidate values | | |
|---------------------|------------------|--------------|--------------|
| | Cora | CiteSeer | PubMed |
| learning rate | {5e-5} | {1e-5} | {5e-4} |
| weight decay | {5e-4, 1e-3} | {5e-3, 1e-2} | {5e-4, 1e-3} |
| latent length | {16, 32} | {32, 64} | {32, 64} |
| latent dimension | {32, 64} | {256, 512} | {32, 64} |
| # of MHCA heads | {64, 128} | {128, 256} | {16, 32} |
| # of MHSA heads | {8, 16} | {4, 8} | {4, 8} |
| MHCA head dimension | {128, 256} | {128, 256} | {128, 256} |
| MHSA head dimension | {128, 256} | {128, 256} | {128, 256} |
| depth | {1} | {1} | {1} |

Table 7: Hyperparameters for the link prediction task. MHCA and MHSA denote the multi-head cross attention and multi-head self attention, respectively.

8.1.6. Hyperparameter Setting for Graph Classification

Table 8 describes the hyperparameter search spaces for the graph classification task.

| Hyperparameter | Candidate values | | | | |
|---------------------|--------------------|--------------------|--------------------|---------------|--------------|
| | MUTAG | PROTEINS | IMDB-BINARY | REDDIT-BINARY | COLLAB |
| learning rate | {5e-3, 1e-3, 5e-4} | {5e-3, 1e-3, 5e-4} | {5e-3, 1e-3} | {5e-3} | {5e-3, 1e-3} |
| weight decay | {5e-4, 1e-4} | {5e-4, 1e-4} | {5e-4, 1e-4} | {5e-4, 1e-4} | {5e-4, 1e-4} |
| latent length | {8, 16, 32} | {16, 32, 64} | {32, 64, 128} | {32, 64} | {64} |
| latent dimension | {32, 64, 128} | {32, 64, 128} | {32, 64, 128, 256} | {32, 64} | {256} |
| # of MHCA heads | {1} | {1} | {1} | {8, 16} | {8, 16} |
| # of MHSA heads | {1} | {1} | {1} | {4} | {4} |
| MHCA head dimension | {4} | {4} | {4} | {4} | {4} |
| MHSA head dimension | {4} | {4} | {4} | {4} | {4} |
| depth | {1} | {1} | {1} | {1} | {1} |

Table 8: Hyperparameters for the graph classification task. MHCA and MHSA denote the multi-head cross attention and multi-head self attention, respectively.

8.1.7. Hyperparameter setting for Multimodal Few-shot Classification

Table 9 describes the hyperparameter search spaces for the multimodal few-shot classification task.

| Hyperparameter | Candidate values |
|---------------------|------------------|
| learning rate | {5e-6, 1e-5} |
| weight decay | {0, 1e-6, 1e-5} |
| latent length | {32} |
| latent dimension | {32, 64} |
| # of MHCA heads | {16, 64} |
| # of MHSA heads | {4} |
| MHCA head dimension | {256, 512} |
| MHSA head dimension | {256, 512} |
| depth | {1, 2} |

Table 9: Hyperparameters for the multimodal few-shot classification task. MHCA and MHSA denote the multi-head cross attention and multi-head self attention, respectively.

8.1.8. Hyperparameter setting for Multimodal Text Classification

Table 10 describes the hyperparameter search spaces for the multimodal text classification task.

| Hyperparameter | Candidate values | |
|---------------------|------------------|--------------|
| | ogbn-products | ogbn-arxiv |
| learning rate | {1e-3} | {5e-4} |
| weight decay | {5e-4} | {5e-4} |
| latent length | {16, 32, 64} | {16, 32, 64} |
| latent dimension | {128, 256} | {128, 256} |
| # of MHCA heads | {8, 16} | {8, 16} |
| # of MHSA heads | {8, 16, 32} | {8, 16, 32} |
| MHCA head dimension | {64, 128} | {64, 128} |
| MHSA head dimension | {64, 128} | {64, 128} |
| depth | {1} | {1} |

Table 10: Hyperparameters for the multimodal text classification task. MHCA and MHSA denote the multi-head cross attention and multi-head self attention, respectively.

8.1.9. Additional Experimental Settings

For all classification tasks except link prediction, we used the accuracy metric to measure whether the node, graph, text, or image is correctly classified. For link prediction, a binary classification task, we adopt average precision (AP) and area under the ROC curve (AUC) to effectively handle class imbalance, as non-existent edges are more frequent than existing ones.

We validate our models on RTX A6000, RTX3090, and A100 GPU. We adopt random seed value 2025 for all experiments. We use PyTorch 1.11.0 and PyTorch-Geometric 2.0.1 with CUDA 11.2.

8.2. Additional Experimental Results

In this section, we provide additional quantitative results.

8.2.1. Link Prediction

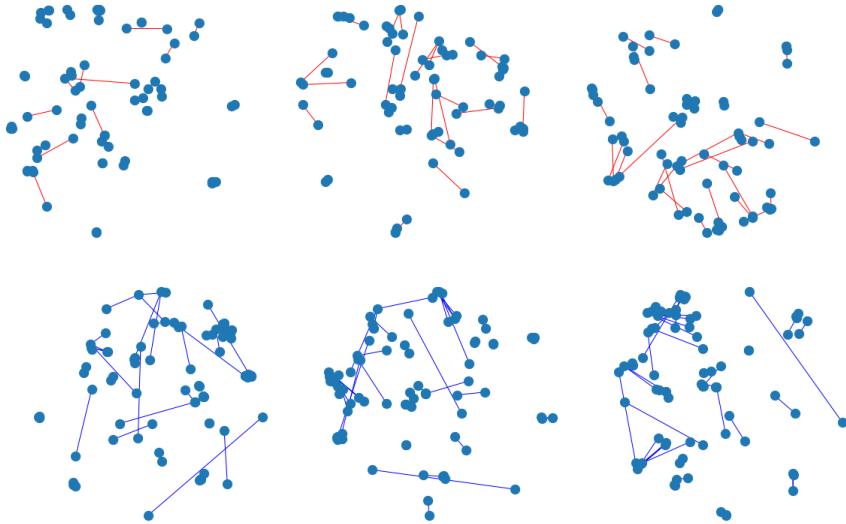


Figure 7: The first row denotes the t-SNE of positive edges (red color) learned by GPIO (left), GNAE (middle), VGNAE (right). The second row denotes the t-SNE of negative edges (blue color) learned by GPIO (left), GNAE (middle), VGNAE (right). Compared to other models, our model embeds the node pairs of the positive edge set closer together, indicating that the node representation is better in the edge level task. Contrary, in the negative edge case, the model is more expressive when the distance between two nodes increases.

We analyze the impact of the smoothing times L and the effect of $RWPE$ on link prediction. Table 11 compares performance to the baseline, showing how our model performs as L varies from 0 to 2 and based on the presence or absence of $RWPE$. Feature smoothing that reflects relational information performs better than without it. In the Link prediction task, we can see that the presence of $RWPE$ does not have a significant impact on performance, as enough relational information is captured by feature smoothing.

| Model | <i>RWPE</i> | Output query smoothing | Cora | | CiteSeer | | PubMed | |
|---------------------------|-------------|---------------------------|-----------------|------------------|-----------------|-----------------|------------------|------------------|
| | | | AUC | AP | AUC | AP | AUC | AP |
| Spectral Clustering [56]* | | | 84.6±0.01 | 88.5±0.0 | 80.5±0.01 | 85.0±0.01 | 84.2±0.02 | 87.8±0.01 |
| DeepWalk [57]* | | | 83.1±0.01 | 85.0±0.0 | 80.5±0.02 | 83.6±0.01 | 84.4±0.0 | 84.1±0.0 |
| DGI [58]* | | | 89.8±0.8 | 89.7±1.0 | 95.5±1.0 | 95.7±1.0 | 91.2±0.6 | 92.2±0.5 |
| ARGVA [59]* | | | 92.4±0.004 | 93.2±0.003 | 92.4±0.003 | 93.0±0.003 | 96.8±0.001 | 97.1±0.001 |
| GIC [45] | | | 90.0±1.0 | 89.9±1.3 | 95.8±0.7 | 95.8±0.9 | 90.9±1.0 | 91.6±0.9 |
| VGAE [31] | | | 95.2±0.5 | 94.7±0.6 | 92.0±1.7 | 91.6±1.7 | 95.6±0.7 | 95.3±0.6 |
| GNAE [33] | | | 95.6±0.7 | 96.0±0.8 | 97.2±0.5 | 97.3±0.4 | 97.7±0.2 | <u>97.6±0.2</u> |
| VGNAE [33] | | | 95.8±0.6 | 95.7±0.8 | 96.8±0.6 | 96.7±0.6 | 97.3±0.1 | 97.2±0.2 |
| GPIO | - | $L = 0$ | 94.4±0.3 | 95.4 ±0.3 | 91.5±0.3 | 91.1±0.7 | 95.9±0.1 | 95.9±0.09 |
| | ✓ | $L = 0$ | 94.3±0.4 | 95.2 ±0.3 | 89.3±0.7 | 91.0±0.5 | 95.8±0.1 | 95.7±0.1 |
| | - | $L = 1$ | 95.6±0.2 | <u>96.3±0.1</u> | 96.2±0.3 | 96.7±0.2 | <u>97.6±0.08</u> | 97.8±0.07 |
| | ✓ | $L = 1$ | 95.9±0.3 | 96.4 ±0.2 | 95.6±0.5 | 96.1±0.3 | 97.6±0.06 | 97.7±0.08 |
| | - | $L = 2$ | <u>95.8±0.2</u> | 96.3 ±0.1 | <u>96.8±0.3</u> | <u>97.2±0.2</u> | 97.0±0.09 | 97.2±0.08 |

Table 11: The performance of link prediction task. The GPIO shows competitive performance compared with the state-of-the-art model, [33]. Also, the GPIO shows the superior results with the well-known link prediction model VGAE, which adopts the graph convolutional neural networks. * denotes the reported performance in [45]. We use 4 dimensional *RWPE* in link prediction. $L = 0, 1, 2$ denotes number of smoothing.

8.2.2. Node Classification

We analyzed the impact of the output query array smoothing type and its hyperparameters on node classification and the effect on *RWPE*. We find that the output query array smoothed by APPNP shows comparable performance with the GNNs. Since \hat{X}_{APPNP} has less oversmoothing, we can set the L value to be larger than \hat{X} , and it has the effect of propagating the score \hat{S}_{APPNP} to more distant nodes and is effective in node classification. On the other hand, as with link prediction, *RWPE* does not have a significant impact, as the relational information is well reflected by \hat{X}_{APPNP} .

To analyze the attention properties of the GPIO, we visualize the attention weights of cross-attention and the attention between the latent arrays

| Model | <i>RWPE</i> | Output query smoothing | Cora | | CiteSeer | | PubMed | |
|------------|-------------|---------------------------|------------------|------------------|-----------------|-----------------|-----------------|-----------------|
| | | | Fixed | Random | Fixed | Random | Fixed | Random |
| GCN [27] | | | 81.4±0.7 | 78.7±1.7 | 71.1±0.7 | 68.1±1.7 | 78.9±0.6 | 77.3±2.4 |
| GAT [60] | | | 83.0±0.5 | 80.9±1.5 | 70.9±0.5 | 68.8±1.7 | 78.9±0.4 | 77.6±2.4 |
| Cheb [61] | | | 80.5±1.0 | 77.0±2.7 | 69.8±1.2 | 67.2±2.1 | 78.2±0.6 | 75.4±2.5 |
| SGC [39] | | | 81.7±0.05 | 80.1±1.9 | 71.3±0.2 | 68.5±2.0 | 78.9±0.0 | 76.6±2.4 |
| ARMA [62] | | | 82.2±0.9 | 79.8±1.7 | 71.0±0.6 | 67.9±1.9 | 78.8±0.3 | 77.6±2.2 |
| APPNP [40] | | | 83.3±0.5 | 82.1±1.5 | 71.7±0.5 | 69.8±1.8 | 80.1±0.2 | 79.1±2.3 |
| GPIO | - | $L = 0$ | 59.3 ±1.2 | 57.9 ±2.6 | 57.4 ±1.3 | 55.0 ±2.3 | 72.3 ±1.7 | 70.0 ±3.0 |
| | ✓ | $L = 0$ | 58.6 ±1.9 | 57.4 ±2.3 | 57.4 ±1.5 | 55.5 ±2.5 | 73.0 ±1.1 | 70.0 ±2.6 |
| | - | $L = 2$ | 81.8 ±0.7 | 79.9 ±1.8 | 69.3 ±1.2 | 67.6 ±1.8 | 78.6 ±0.5 | 77.9 ±2.4 |
| | - | $L = 3$ | 82.9 ±0.7 | 80.7 ±1.6 | 68.9 ±1.1 | 67.7 ±1.9 | 78.2 ±0.4 | 78.0 ±2.4 |
| | - | $L = 10, \alpha = 0.1$ | 83.9 ±0.6 | <u>81.6 ±1.8</u> | 70.1 ±1.0 | 68.1 ±1.7 | 79.8 ±0.5 | 79.6±2.2 |
| | ✓ | $L = 10, \alpha = 0.1$ | 83.7 ±0.8 | 81.5 ±1.5 | 69.2 ±1.0 | 67.9 ±1.7 | <u>79.9±0.4</u> | 79.0 ±2.5 |

Table 12: Accuracy for the node classification. The GPIO shows the competitive results on the PubMed. Fixed and Random denotes the dataset split methods from [46]. We use 4 or 8 dimensional *RWPE* in node classification. $L = 0, 2, 3$ denotes number of smoothing, $L = 10, \alpha = 0.1$ denotes hyperparameters of APPNP.

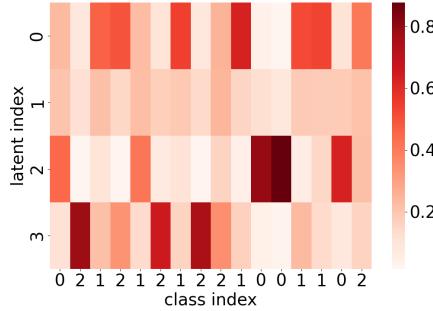


Figure 8: GPIO attention heatmap for each latent index on PubMed. The x-axis and y-axis indicate the node class and latent index, respectively. Each latent of the GPIO captures the different nodes with diversity. Besides, each latent prone to focus on the specific class.

and the input arrays as shown in Figure 8. Each latent index in the GPIO captures the different nodes with high diversity. Additionally, each latent

focuses on the specific class group of nodes relatively.

8.2.3. Graph Classification

| Model | MUTAG | PROTEINS | IMDB-BINARY | REDDIT-BINARY | COLLAB |
|--------------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| GCNWithJK | 72.9±12.0 | 72.6±3.6 | 73.2±5.0 | 89.4±2.9 | 81.5±2.1 |
| SAGPool | 74.0±8.7 | 72.3±2.8 | 72.3±4.7 | 89.0±2.1 | 78.9±1.0 |
| DiffPool | <u>84.6±8.7</u> | <u>74.3±6.4</u> | 74.8±4.8 | 92.7±2.0 | 79.4±1.9 |
| EdgePool | 72.3±13.4 | 72.6±3.4 | 73.3±5.3 | 89.2±3.8 | 79.3±1.2 |
| GCN | 70.7±11.0 | 72.2±2.4 | 74.2±4.4 | 89.1±2.0 | <u>81.0±1.4</u> |
| GraphSAGE | 75.1±11.4 | 74.1±2.3 | 73.2±4.4 | 90.7±2.3 | 79.8±1.1 |
| GIN0 | 81.9±8.0 | 73.1±3.8 | <u>73.7±4.1</u> | 90.9±2.1 | 80.5±1.9 |
| GIN | 82.9±11.3 | 72.0±3.2 | 73.5±4.9 | 90.5±2.5 | 80.7±2.0 |
| GlobalAttentionNet | 77.7±12.1 | 72.6±2.6 | 72.9±3.7 | 88.6±3.2 | 79.1±0.7 |
| Set2SetNet | 74.5±11.9 | 74.1±3.8 | 72.7±5.0 | 90.3±2.4 | 79.4±1.7 |
| SortPool | 83.0±9.0 | 73.9±4.5 | 71.8±3.0 | 84.3±5.0 | 77.8±1.6 |
| ASAP | 78.7±11.8 | 74.0±3.0 | 72.2±4.3 | OOM | 79.4±1.7 |
| Graph Percevier IO | | | | | |
| + None PE | 69.6±6.6 | 71.1±3.7 | 72.0±5.2 | 86.4±2.9 | 77.7±2.4 |
| + Fourier PE | 83.6±9.5 | 72.5±4.3 | 71.7±4.3 | 76.1±3.5 | 75.5±2.0 |
| + <i>RWPE</i> | 86.1±6.9 | 76.1±3.0 | 72.9±4.4 | <u>91.6±2.3</u> | 79.1±2.5 |

Table 13: Accuracy of the graph classification tasks. The GPIO shows the superior performance on MUTAG and PROTEINS dataset that has node features. IMDB, REDDIT, and COLLAB have no node features, and an understanding of the topological information of the graph is required. The GPIO shows the competitive results even though there are no node features. We use 64 dimensional *RWPE* in graph classification. For None PE in the REDDIT dataset, we remove first layer normalization.

Table 13 denotes how the presence and type of positional encoding affects graph classification performance. Unlike node classification, GPIO uses an

output query array of $1 \times D_q (= E)$ because the model classifies the entire graph rather than each node in graph classification. This means that unlike other tasks, it reflects topological information via *RWPE* without output query array smoothing. The results show that on relatively small and simple MUTAG and PROTEINS datasets, Fourier positional encoding that incorporates only sequential information outperforms no positional encoding. However, in large, complex graphs in other datasets, sequential information caused overfitting. This is because a graph with a same structure seen at training time may have been encoded as a completely different graph at test time, depending on how the Fourier PE were started. In contrast, *RWPE* outperformed None and Fourier PE on all datasets. This is because *RWPE* effectively reflects the graph structure that the GPIO saw at training time to the same or similar graph structure at test time. This demonstrates how effective *RWPE* is at providing canonical positional information in graph classification.

8.3. Graph Related Tasks

| Tasks | input array | Output query array | Logits layer |
|----------------------|--------------------|----------------------|--------------|
| Node classification | $M \times (C + t)$ | $M \times D_q (= C)$ | O |
| Graph classification | $M \times (C + t)$ | $1 \times D_q (= E)$ | O |
| Link prediction | $M \times (C + t)$ | $M \times D_q (= C)$ | X |

Table 14: Configuration of input and output query array used in the actual training process. We adopt a smoothing based output query array for node classification and link prediction. We apply learnable output query array in graph classification. The logits layer serves to map the final output to the class dimension.

The GPIO has the capability to perform diverse graph-related tasks such as prediction and classification. In node classification tasks, the output array obtained after going through the logits layer contains each node’s predicted score. We use the cross-entropy loss to learn node embedding. For graph classification tasks, we set a size of output query array to 1-dimension vector. Then we pass the array to the logits layer (e.g., a single linear layer) instead of the read-out layer usually used for global pooling in graph-level tasks. We evaluate the cross-entropy loss with target label. For link prediction tasks, we conduct the inner product between each pair of node features which is the output of the final cross attention layer. In the same manner, as GAE [31], We employ the reconstruction loss for training our model. Table 14 provides configuration of each array and whether or not the logits layer is used for each task.

8.4. Dataset

Table 15 describes the dataset statistics for graph-related tasks. For the preprocessing and train/test splits, we follow the benchmark experimental settings for each task and dataset. We follow the same settings with the [46] for node classification, graph classification and link prediction, respectively. We provide the dataset as well as the source code as supplementary files. Note that the graph-related dataset will be automatically downloaded and normalized during preprocessing when running the source code.

Cora, CiteSeer, and PubMed are all citation networks. Cora consists of scientific articles classified into seven distinct research areas, each of which is represented by a word vector based on a dictionary of 1,433 words. CiteSeer contains papers classified into six distinct research areas, represented by a

| Dataset | Graphs | Nodes | Edges | Features | Classes |
|---------------|--------|--------|----------|----------|---------|
| Cora [49] | 1 | 2,708 | 5,278 | 1,433 | 7 |
| CiteSeer [50] | 1 | 3,327 | 4,552 | 3,703 | 6 |
| PubMed [51] | 1 | 19,717 | 44,324 | 500 | 3 |
| MUTAG | 188 | 17.93 | 19.79 | 7 | 2 |
| PROTEINS | 1,113 | 39.06 | 72.82 | 3 | 2 |
| IMDB-BINARY | 1,000 | 19.77 | 96.53 | - | 2 |
| REDDIT-BINARY | 2,00 | 429.63 | 497.754 | - | 2 |
| COLLAB | 5,000 | 74.49 | 2,457.22 | - | 3 |

Table 15: Dataset statistics for the graph-related tasks, link prediction, graph classification, and node classification.

word vector using a dictionary of 3,703 words. PubMed includes diabetes-related papers classified into three distinct topics, with each paper represented by a TF-IDF weighted word vector based on a dictionary of 500 words. MUTAG is nitroaromatic compounds dataset aimed at predicting their mutagenicity on *Salmonella typhimurium*. Chemical compounds are represented as graphs, where nodes correspond to (one-hot encoded) atoms and edges represent chemical bonds. PROTEINS is a protein dataset classified as enzymes or non-enzymes, where nodes represent amino acids and edges connect those within a distance of less than 6 Angstroms. IMDB-BINARY is a movie collaboration containing the ego-networks of 1,000 actors and actresses from IMDB. In each graph, nodes represent actors or actresses, and edges connect those who appeared in the same movie. The class of IMDB-BINARY shows reviews as negative or positive. REDDIT-BINARY is a dataset of Reddit discussion graphs, where nodes represent users and edges indicate replies. The graph is classified as either question/answer-based or discussion-based. COLLAB is a dataset of scientific collaboration where each graph represents

a researcher’s ego network, with nodes as the researcher and collaborators, and edges indicating collaborations. Each graph is labeled based on the researcher’s field: High Energy Physics, Condensed Matter Physics, or Astro Physics. In our experimental setup, we did not use edge features for any of the graphs, all of which are undirected.

Also, we provide dataset statistics for the multimodal text classification. The ogbn-products is a graph modeling Amazon’s product co-purchasing network, which is undirected and unweighted. Nodes represent products, and edges indicate items bought together. Node features are derived from product descriptions using a bag-of-words approach [63] with a 100-dimensional feature vector. The goal is to classify products into one of 47 top-level classes. The ogbn-arxiv dataset is a directed graph modeling citation relationships among Computer Science papers on ARXIV. Nodes represent papers, with edges indicating citations. Each paper is characterized by a 128-dimensional feature vector, computed as the average of word embeddings from its title and abstract [64]. For the ogbn-products and ogbn-arxiv datasets, which don’t contain edge features, we followed the standard train/validation/test split provided by OGB [55].

| Dataset | #Nodes | #Edges | #Classes |
|------------------------|---------|-----------|----------|
| ogbn-products (subset) | 54,025 | 74,420 | 47 |
| ogbn-arxiv | 169,343 | 1,166,243 | 40 |

Table 16: Statistics of OGB datasets [55]. They contain a large amount of nodes and edges. ogbn-products is subset data as in [35].

8.5. Additional Material

In this section, we provide additional explanatory material to help you understand the paper.

8.5.1. Additional Illustration Material

To make it easier to understand at a glance, we have shown in Figure 9 how a unified model can take different modalities as input and perform various graph-related and multimodal tasks in a unified way.

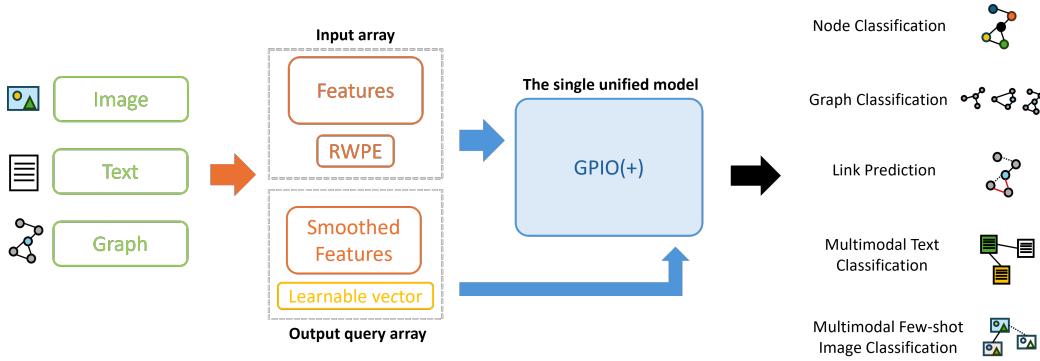


Figure 9: $\text{GPIO}(+)$, a single architecture, takes multiple modalities such as graph text images as input and output queries (e.g. RWPE , Smoothed Features) all in the same way and performs various tasks. The tasks include node classification, graph classification, link prediction, multimodal text classification, and multimodal few-shot image classification.

8.5.2. Space Complexity Analysis

In Section *Input Array*, GPIO does not use the adjacency matrix $A \in \mathbb{R}^{M \times M}$ in training and evaluation directly, so it offers lower space complexity. In practice, most GNNs use a sparse representation and have the space complexity of $O(M \times C + V)$, where V is the number of edges. GPIO , on the other hand, is $O(M \times C)$ since size of the positional encoding t is relatively

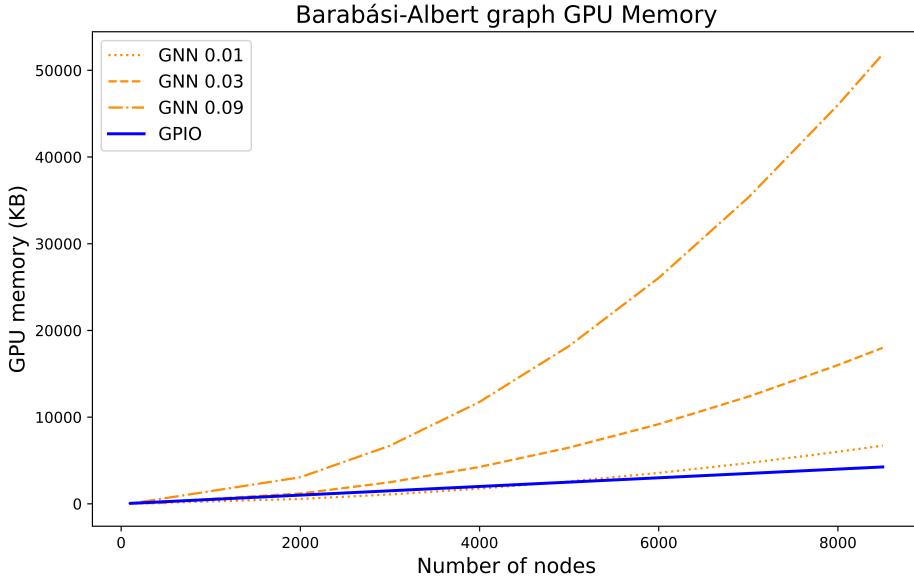


Figure 10: GPU memory usage for Barabási-Albert graph densities 0.01, 0.03, and 0.09. GNN 0.09 in the legend means density 0.09 graph. GPIO is invariant to graph density changes and has lower GPU memory usage compared to GNN at high densities.

small. Therefore, the space complexity is all proportional to the number of nodes M . However, GPIO is not proportional to V . This means that GPIO offer lower space complexity in graphs with a higher density of edges. Figure 10 illustrates the GPU memory usage of GNN and GPIO for a given density and number of nodes. We generated graphs with densities of 0.01, 0.03, and 0.09 using the Barabási-Albert model [65], which captures graph shape similar to real networks. We increased M from about 500 to 8000 and set $C = 64$ and $t = 32$. Both increase GPU memory usage as M increases, but there is a sharp difference at higher density (e.g. GNN 0.03, GNN 0.09 in the legend). Note taht the GPU memory usage of GPIO is not affected by changes in density, and is only affected by M . This result emphasizes that

GPIO is efficient in terms of the space complexity required as inputs compared to a typical GNNs when the number of nodes is large, especially for dense graphs. There are many cases that transform the tasks and problems as a graph-structured dataset, such as traveling salesman problems [66]; the transformed graph is fully connected.

References

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
- [2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, N. Houlsby, An image is worth 16x16 words: Transformers for image recognition at scale, in: *International Conference on Learning Representations*, 2021.
- [3] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al., Learning transferable visual models from natural language supervision, in: *International Conference on Machine Learning*, PMLR, 2021, pp. 8748–8763.
- [4] J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, S. C. H. Hoi, Align before fuse: Vision and language representation learning with momentum distillation, *Advances in Neural Information Processing Systems* 34 (2021).

- [5] A. Jaegle, F. Gimeno, A. Brock, O. Vinyals, A. Zisserman, J. Carreira, Perceiver: General perception with iterative attention, in: International Conference on Machine Learning, PMLR, 2021, pp. 4651–4664.
- [6] A. Jaegle, S. Borgeaud, J.-B. Alayrac, C. Doersch, C. Ionescu, D. Ding, S. Koppula, D. Zoran, A. Brock, E. Shelhamer, O. J. Henaff, M. Botvinick, A. Zisserman, O. Vinyals, J. Carreira, Perceiver IO: A general architecture for structured inputs & outputs, in: International Conference on Learning Representations, 2022.
- [7] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zam baldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner, et al., Relational inductive biases, deep learning, and graph networks, arXiv preprint arXiv:1806.01261 (2018).
- [8] J. Gong, Y. Zhao, J. Zhao, J. Zhang, G. Ma, S. Zheng, S. Du, J. Tang, Personalized recommendation via inductive spatiotemporal graph neural network, Pattern Recognition 145 (2024) 109884.
- [9] X.-b. Ye, Q. Guan, W. Luo, L. Fang, Z.-R. Lai, J. Wang, Molecular substructure graph attention network for molecular property identification in drug discovery, Pattern Recognition 128 (2022) 108659.
- [10] Y. Zhang, Y. Li, X. Wei, Y. Yang, L. Liu, Y. L. Murphey, Graph matching for knowledge graph alignment using edge-coloring propagation, Pattern Recognition 144 (2023) 109851.
- [11] X. Zhang, P. Angeloudis, Y. Demiris, Dual-branch spatio-temporal

graph neural networks for pedestrian trajectory prediction, Pattern Recognition 142 (2023) 109633.

- [12] A. Chelladurai, D. L. Narayan, P. B. Divakarachari, U. Loganathan, fmri-based alzheimer's disease detection using the sas method with multi-layer perceptron network, Brain Sciences 13 (2023) 893.
- [13] Y. B. Özçelik, A. Altan, Overcoming nonlinear dynamics in diabetic retinopathy classification: a robust ai-based model with chaotic swarm intelligence optimization and recurrent long short-term memory, Fractal and Fractional 7 (2023) 598.
- [14] J. Xu, C. Yuan, X. Ma, H. Shang, X. Shi, X. Zhu, Interpretable medical deep framework by logits-constraint attention guiding graph-based multi-scale fusion for alzheimer's disease analysis, Pattern Recognition 152 (2024) 110450.
- [15] C. Yu, H. Pei, Dynamic graph clustering learning for unsupervised diabetic retinopathy classification, Diagnostics 13 (2023) 3251.
- [16] J. Li, D. Li, S. Savarese, S. Hoi, Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, in: International conference on machine learning, PMLR, 2023, pp. 19730–19742.
- [17] Z. Sun, Y. Fang, T. Wu, P. Zhang, Y. Zang, S. Kong, Y. Xiong, D. Lin, J. Wang, Alpha-clip: A clip model focusing on wherever you want, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024, pp. 13019–13029.

- [18] P. Mayilvahanan, T. Wiedemer, E. Rusak, M. Bethge, W. Brendel, Does clip’s generalization performance mainly stem from high train-test similarity?, arXiv preprint arXiv:2310.09562 (2023).
- [19] A. v. d. Oord, Y. Li, O. Vinyals, Representation learning with contrastive predictive coding, arXiv preprint arXiv:1807.03748 (2018).
- [20] M. Lan, F. Rong, Z. Li, W. Yu, L. Zhang, Bidirectional correlation-driven inter-frame interaction transformer for referring video object segmentation, Pattern Recognition 153 (2024) 110535.
- [21] Z. Liu, L. Cai, W. Yang, J. Liu, Sentiment analysis based on text information enhancement and multimodal feature fusion, Pattern Recognition 156 (2024) 110847.
- [22] D. Wang, X. Guo, Y. Tian, J. Liu, L. He, X. Luo, Tetfn: A text enhanced transformer fusion network for multimodal sentiment analysis, Pattern Recognition 136 (2023) 109259.
- [23] A. Baevski, W.-N. Hsu, Q. Xu, A. Babu, J. Gu, M. Auli, data2vec: A general framework for self-supervised learning in speech, vision and language, in: K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, S. Sabato (Eds.), Proceedings of the 39th International Conference on Machine Learning, volume 162 of *Proceedings of Machine Learning Research*, PMLR, 2022, pp. 1298–1312.
- [24] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, A. Joulin, Emerging properties in self-supervised vision transformers, in:

Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 9650–9660.

- [25] H. Bao, L. Dong, S. Piao, F. Wei, BEit: BERT pre-training of image transformers, in: International Conference on Learning Representations, 2022.
- [26] J. Carreira, S. Koppula, D. Zoran, A. Recasens, C. Ionescu, O. J. Hénaff, E. Shelhamer, R. Arandjelovic, M. Botvinick, O. Vinyals, et al., Hierarchical perceiver (2022).
- [27] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, in: International Conference on Learning Representations, 2017.
- [28] G. Li, M. Müller, B. Ghanem, V. Koltun, Training graph neural networks with 1000 layers, in: International conference on machine learning, PMLR, 2021, pp. 6437–6449.
- [29] E. Rossi, B. Charpentier, F. Di Giovanni, F. Frasca, S. Günnemann, M. M. Bronstein, Edge directionality improves learning on heterophilic graphs, in: Learning on Graphs Conference, PMLR, 2024, pp. 25–1.
- [30] Y. Huang, Y. Zeng, Q. Wu, L. Lü, Higher-order graph convolutional network with flower-petals laplacians on simplicial complexes, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, 2024, pp. 12653–12661.
- [31] T. N. Kipf, M. Welling, Variational graph auto-encoders, arXiv preprint arXiv:1611.07308 (2016).

- [32] D. P. Kingma, M. Welling, Auto-encoding variational bayes, CoRR abs/1312.6114 (2013).
- [33] S. J. Ahn, M. Kim, Variational graph normalized autoencoders, in: Proceedings of the 30th ACM International Conference on Information & Knowledge Management, 2021, pp. 2827–2831.
- [34] J. Kim, D. Nguyen, S. Min, S. Cho, M. Lee, H. Lee, S. Hong, Pure transformers are powerful graph learners, Advances in Neural Information Processing Systems 35 (2022) 14582–14595.
- [35] X. He, X. Bresson, T. Laurent, A. Perold, Y. LeCun, B. Hooi, Harnessing explanations: Llm-to-lm interpreter for enhanced text-attributed graph representation learning, arXiv preprint arXiv:2305.19523 (2023).
- [36] J. Guo, L. Du, H. Liu, M. Zhou, X. He, S. Han, Gpt4graph: Can large language models understand graph structured data? an empirical evaluation and benchmarking, arXiv preprint arXiv:2305.15066 (2023).
- [37] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [38] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, H. Jégou, Training data-efficient image transformers & distillation through attention, in: International Conference on Machine Learning, PMLR, 2021, pp. 10347–10357.
- [39] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, K. Weinberger, Simplifying

graph convolutional networks, in: International conference on machine learning, PMLR, 2019, pp. 6861–6871.

- [40] J. Gasteiger, A. Bojchevski, S. Günnemann, Combining neural networks with personalized pagerank for classification on graphs, in: International Conference on Learning Representations, 2019.
- [41] V. P. Dwivedi, A. T. Luu, T. Laurent, Y. Bengio, X. Bresson, Graph neural networks with learnable structural and positional representations, in: International Conference on Learning Representations, 2022.
- [42] B. Weisfeiler, A. Leman, The reduction of a graph to canonical form and the algebra which appears therein, nti, Series 2 (1968) 12–16.
- [43] V. P. Dwivedi, X. Bresson, A generalization of transformer networks to graphs, arXiv preprint arXiv:2012.09699 (2020).
- [44] J. Kim, T. Kim, S. Kim, C. D. Yoo, Edge-labeling graph neural network for few-shot learning, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 11–20.
- [45] C. Mavromatis, G. Karypis, Graph infoclust: Maximizing coarse-grain mutual information in graphs, in: K. Karlapalem, H. Cheng, N. Ramakrishnan, R. K. Agrawal, P. K. Reddy, J. Srivastava, T. Chakraborty (Eds.), Advances in Knowledge Discovery and Data Mining - 25th Pacific-Asia Conference, PAKDD 2021, Virtual Event, May 11-14, 2021, Proceedings, Part I, volume 12712 of *Lecture Notes in Computer Science*, Springer, 2021, pp. 541–553.

- [46] M. Fey, J. E. Lenssen, Fast graph representation learning with PyTorch Geometric, in: ICLR 2019 Workshop on Representation Learning on Graphs and Manifolds, 2019.
- [47] L. Van der Maaten, G. Hinton, Visualizing data using t-sne., Journal of machine learning research 9 (2008).
- [48] T. Wang, P. Isola, Understanding contrastive representation learning through alignment and uniformity on the hypersphere, in: International Conference on Machine Learning, PMLR, 2020, pp. 9929–9939.
- [49] A. K. McCallum, K. Nigam, J. Rennie, K. Seymore, Automating the construction of internet portals with machine learning, Information Retrieval 3 (2000) 127–163.
- [50] C. L. Giles, K. D. Bollacker, S. Lawrence, CiteSeer: An automatic citation indexing system, in: Proceedings of the Third ACM Conference on Digital Libraries, Pittsburgh, PA, USA, ACM, New York, NY, USA, 1998, pp. 89–98.
- [51] G. M. Namata, B. London, L. Getoor, B. Huang, Query-driven active surveying for collective classification, in: Workshop on Mining and Learning with Graphs, 2012.
- [52] A. K. Debnath, R. L. Lopez de Compadre, G. Debnath, A. J. Shusterman, C. Hansch, Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity, Journal of medicinal chemistry 34 (1991) 786–797.

- [53] P. Yanardag, S. Vishwanathan, Deep graph kernels, in: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, 2015, pp. 1365–1374.
- [54] Y. Liu, J. Lee, M. Park, S. Kim, E. Yang, S. Hwang, Y. Yang, LEARNING TO PROPAGATE LABELS: TRANSDUCTIVE PROPAGATION NETWORK FOR FEW-SHOT LEARNING, in: International Conference on Learning Representations, 2019.
- [55] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, J. Leskovec, Open graph benchmark: Datasets for machine learning on graphs, *Advances in neural information processing systems* 33 (2020) 22118–22133.
- [56] L. Tang, H. Liu, Leveraging social media networks for classification, *Data Mining and Knowledge Discovery* 23 (2011) 447–478.
- [57] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: Online learning of social representations, in: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, 2014, pp. 701–710.
- [58] P. Velickovic, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, R. D. Hjelm, Deep graph infomax., *ICLR (Poster)* 2 (2019) 4.
- [59] S. Pan, R. Hu, S.-f. Fung, G. Long, J. Jiang, C. Zhang, Learning graph embedding with adversarial training methods, *IEEE transactions on cybernetics* 50 (2019) 2475–2487.

- [60] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, Y. Bengio, Graph attention networks, arXiv preprint arXiv:1710.10903 (2017).
- [61] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, Advances in neural information processing systems 29 (2016).
- [62] F. M. Bianchi, D. Grattarola, L. Livi, C. Alippi, Graph neural networks with convolutional arma filters, IEEE transactions on pattern analysis and machine intelligence (2021).
- [63] W.-L. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, C.-J. Hsieh, Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks, in: Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining, 2019, pp. 257–266.
- [64] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, Advances in neural information processing systems 26 (2013).
- [65] R. Albert, A.-L. Barabási, Statistical mechanics of complex networks, Reviews of modern physics 74 (2002) 47.
- [66] W. Kool, H. van Hoof, M. Welling, Attention, learn to solve routing problems!, in: International Conference on Learning Representations, 2019.