

MECCH: Metapath Context Convolution-based Heterogeneous Graph Neural Networks

Xinyu Fu^a, Irwin King^a

^aDepartment of Computer Science and Engineering, The Chinese University of Hong Kong, Hong Kong, China

Abstract

Heterogeneous graph neural networks (HGNNs) were proposed for representation learning on structural data with multiple types of nodes and edges. To deal with the performance degradation issue when HGNNs become deep, researchers combine metapaths into HGNNs to associate nodes closely related in semantics but far apart in the graph. However, existing metapath-based models suffer from either information loss or high computation costs. To address these problems, we present a novel *Metapath Context Convolution-based Heterogeneous Graph Neural Network* (MECCH). MECCH leverages *metapath contexts*, a new kind of graph structure that facilitates lossless node information aggregation while avoiding any redundancy. Specifically, MECCH applies three novel components after feature preprocessing to extract comprehensive information from the input graph efficiently: (1) metapath context construction, (2) metapath context encoder, and (3) convolutional metapath fusion. Experiments on five real-world heterogeneous graph datasets for node classification and link prediction show that MECCH achieves superior prediction accuracy compared with state-of-the-art baselines with improved computational efficiency. The code is available at <https://github.com/cynricfu/MECCH>. The formal publication is available at <https://doi.org/10.1016/j.neunet.2023.11.030>.

Keywords: Graph neural networks, Heterogeneous information networks, Graph representation learning

1. Introduction

Many real-world networks are *heterogeneous graphs*, which contain multiple types of nodes and edges. As illustrated in Figure 1, a movie information network may consist of actors, movies, directors, and different types of relationships between them. The complex and irregular interactions among different types of nodes and edges make it challenging to extract knowledge from heterogeneous graphs efficiently. Therefore, heterogeneous graph representation learning, which aims to represent nodes using low-dimensional vectors, is a desirable way to automatically process and make inferences on such data.

Over the past decade, heterogeneous graph representation learning has drawn significant attention. Early attempts usually combine skip-gram model (Mikolov et al., 2013a) and metapath-guided random walks (Dong et al., 2017; Fu et al., 2017; Shi et al., 2019). With the rapid development of deep learning, graph neural networks (GNNs) (Kipf and Welling, 2017; Hamilton et al., 2017; Velickovic et al., 2018) are proposed to incorporate node features and benefit from neural network architectures. Initially, GNNs focused on homogeneous graphs. But it is straightforward for researchers to generalize GNNs to the heterogeneous scenario, where multiple types of nodes and edges introduce another layer of complexity into the GNN design.

Recent efforts in heterogeneous GNNs (HGNNs) can be divided into two categories: *relation-based HGNNs* and *metapath-based HGNNs*. Relation-based HGNNs consider message passing parameterized by edge types and aggregate information from direct neighbors (Schlichtkrull et al., 2018; Zhang et al., 2019a; Hu et al., 2020; Lv et al., 2021; Yang et al., 2023d). Models of this kind are generally simple and fast. Still, they usu-

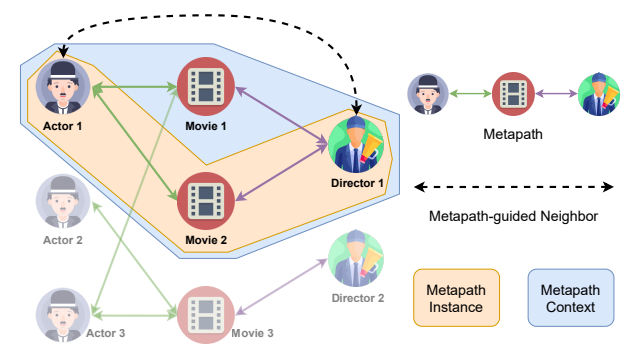


Figure 1: An illustration of the heterogeneous graph and related concepts.

ally require stacking many GNN layers to leverage information multiple hops away, potentially deteriorating model performance (Li et al., 2018; Zhou et al., 2021). Another line of research takes advantage of *metapaths*, which are ordered sequences of node and edge types describing composite relationships between nodes (as shown in Figure 1). Metapath-based HGNNs consider message passing through metapaths and aggregate information from metapath-guided neighbors. Through this way, the models can comfortably obtain information multiple hops away with very few layers and capture high-level semantic information embedded in the graph.

However, existing metapath-based HGNNs can hardly achieve a good balance between model performance and computational efficiency. HAN (Wang et al., 2019a) aggregates metapath-guided neighbors, but with information loss caused by discarded intermediate nodes along metapaths. MAGNN (Fu et al., 2020) addresses this issue by encoding metapath instances. But the

computational complexity is significantly increased due to redundant computations introduced. Unlike HAN and MAGNN, which choose metapaths based on human expertise, GTN (Yun et al., 2019) tries to automatically select informative metapaths using learnable weights. However, GTN also suffers from high computation costs and massive memory usage since it multiplies the entire adjacency matrix of the input graph.

These problems with existing models drive us to design a new metapath-based HGNN that can extract comprehensive knowledge from the input graph while maintaining modest computation time and memory footprint. Through analysis, we discover that *metapath contexts*, a new kind of graph structure defined by us, can enable lossless node information aggregation without redundant computations. In this paper, we propose a novel *MEtapath Context Convolution-based Heterogeneous Graph Neural Network* (MECCH) to learn node representations from heterogeneous graphs. Specifically, after feature preprocessing, MECCH first constructs metapath contexts for each node to describe the neighboring structures formed by the metapaths. Next, MECCH efficiently encodes the constructed metapath contexts to integrate information from both metapath-guided neighbors and intermediate nodes along metapaths. In the end, MECCH adaptively combines node representations from different metapaths using an efficient 1-D convolution kernel. In summary, our main contributions are as follows:

1. We provide a unified framework for metapath-based HGNNs and analyze their limitations, facilitating new HGNN design with in-depth model understanding.
2. We introduce metapath contexts, a new kind of graph structure that facilitates lossless node information aggregation while avoiding any redundancy.
3. We propose a novel metapath context convolution-based HGNN (MECCH) that can capture comprehensive information without incurring redundant computations.
4. We conduct extensive experiments on five heterogeneous graph datasets for node classification and link prediction against multiple state-of-the-art HGNNs. MECCH consistently outperforms these baselines while improving computational efficiency.

Paper Structure. Section 2 briefly introduces the historical development of heterogeneous graph representation learning and especially heterogeneous GNNs. Section 3 presents the definitions and graphical illustrations of the important terminologies used in this paper. Section 4 formulates a general framework for metapath-based HGNNs and presents the novel HGNN model we propose for effective and efficient learning on heterogeneous graphs. Section 5 describes the experimental settings and evaluation results on the node classification and link prediction tasks. Section 6 concludes the paper with possible future directions.

2. Related Work

Heterogeneous graph representation learning aims to transform nodes in a heterogeneous graph into low-dimensional vector representations that can preserve the rich semantic information from the original graph structure and node features.

2.1. Shallow Models

Inspired by the homogeneous graph embedding model DeepWalk (Perozzi et al., 2014), metapath2vec (Dong et al., 2017) generates random walks guided by a pre-defined metapath, which are then fed to a skip-gram model to generate node embeddings. HIN2vec (Fu et al., 2017) and HERec (Shi et al., 2019) also leverage metapaths in different ways. Some models capture node proximity at the neighborhood level (Chen et al., 2018) or at the schema level (Zhao et al., 2020). Some other embedding models focus on knowledge graphs, i.e., heterogeneous graphs with rich schemas (Bordes et al., 2013; Yang et al., 2015; Sun et al., 2019). There are also models that conduct heterogeneous graph embedding in non-Euclidean vector spaces (Wang et al., 2019b, 2021a). These shallow models do not utilize node features, resulting in a heavy loss of important information. They are also trained in an unsupervised fashion. Although their generated embeddings can be fed to downstream models for various tasks, there is a performance gap between them and those end-to-end task-dependent models, i.e., GNNs.

2.2. Graph Neural Networks

GNNs have been a recent trend in graph representation learning. The rationale behind GNNs is that each node can be characterized by its features and local neighborhood. Following this idea, GCN (Kipf and Welling, 2017), GAT (Velickovic et al., 2018), and many other powerful GNN architectures (Hamilton et al., 2017; Zhang et al., 2018, 2019b; Yang et al., 2021) have been proposed. GNNs can leverage node features and benefit from various training paradigms of deep learning, achieving state-of-the-art results in many homogeneous graph datasets (Liang et al., 2022; Yang et al., 2023a,b,c) with various applications, such as recommendation (Chen et al., 2023b,c), biochemistry (Meng et al., 2023a,b,c), and natural language (Song and King, 2022; Ma et al., 2023). It is natural for researchers to generalize GNNs to the heterogeneous setting. We summarize recent HGNN designs into relation-based models and metapath-based models.

2.3. Relation-based HGNNs

Relation-based HGNNs aggregate direct neighbors modulated by type-specific weights. By extending GCN with edge-type-specific weights, RGCN (Schlichtkrull et al., 2018) groups and aggregates the neighborhood based on the edge types. HetSANN (Hong et al., 2020) employs type-specific graph attention layers to aggregate local information. HGT (Hu et al., 2020) leverages a Transformer-like encoder (Vaswani et al., 2017) parameterized by node and edge types to aggregate the local neighborhood. Simple-HGN (Lv et al., 2021) extends GAT to the heterogeneous setting via edge-type attention. R-HGNN (Yu et al., 2023) jointly learns node and relation representations via relation-specific and cross-relation message passings. Many follow-up works are based on a similar idea (Zhang et al., 2019a; Yang et al., 2023d; Yu et al., 2020; Yoon et al., 2022; Ahn et al., 2022), which we do not describe in detail. Despite simplicity and efficiency, the only way for relation-based HGNNs to capture information from nodes multiple hops away is to stack

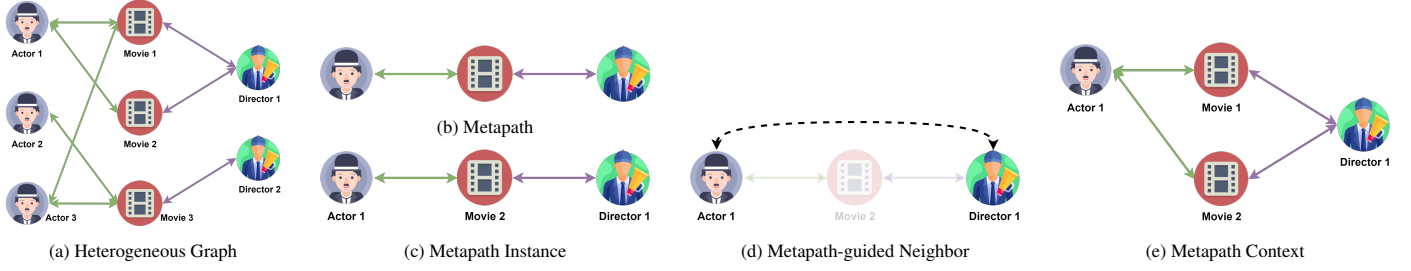


Figure 2: An illustration of the terminologies defined in Section 3.

more GNN layers, bringing along potential performance degradation caused by the over-smoothing problem (Li et al., 2018) and the over-use of transformation operations (Zhou et al., 2021).

2.4. Metapath-based HGNNs

To address the limitations stated above, some researchers designed HGNNs leveraging metapaths to associate nodes far from each other, eliminating the need for stacking multiple layers. By connecting metapath-guided neighbors and discarding intermediate nodes, HAN (Wang et al., 2019a) converts a heterogeneous graph into multiple homogeneous subgraphs. It then applies GAT to each subgraph and combines node representations from different metapaths using attention. MAGNN (Fu et al., 2020) incorporates the intermediate nodes discarded by HAN through leveraging metapath instances. However, this raises an efficiency issue caused by duplicate nodes across an exploding number of metapath instances. Metapaths in HAN and MAGNN are manually selected. To automate this process, GTN (Yun et al., 2019) softly selects relations using learnable weights and applies matrix multiplication to obtain the adjacency matrix of arbitrary metapaths, which is fed to a GCN to get node representations. Useful metapaths can be automatically extracted based on the relation weights. GTN requires a tremendous amount of memory space and computation time due to the large matrix multiplication step. Based on GTN, MEGNN (Chang et al., 2022) rearranges the calculation order of the matrix multiplications to improve the efficiency.

Despite their success in different datasets, existing metapath-based HGNNs cannot have it all for performance and efficiency. HAN suffers from information loss. MAGNN and GTN are limited by high computation costs. MEGNN still faces the scalability issue with large graphs because of large matrix multiplication. Derivative studies based on these models share similar issues (Zhang et al., 2022; Wang et al., 2022). We leave the detailed illustration and discussion of the underlying mechanisms that cause these issues to Figure 3 and Section 4.1.3. Therefore, we want to design a metapath-based HGNN that can capture the most information from the input heterogeneous graph without information loss, and at the same time can be efficient without redundant computations. We achieve these goals by defining and utilizing a special structure called *metapath context*. The advantages of metapath context and how it can help avoid the issues mentioned above will be detailed in Section 4.2.

Table 1: Notations used in this paper.

Notation	Definition
\mathbb{R}^n	n -dimensional Euclidean space
$a, \mathbf{a}, \mathbf{A}$	Scalar, vector, matrix
\mathbf{A}^\top	Matrix/vector transpose
$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	A graph \mathcal{G} with node set \mathcal{V} and edge set \mathcal{E}
v	A node $v \in \mathcal{V}$
P	A metapath
\mathcal{P}_A	Metapaths starting from node type A
\mathcal{N}_v^P	Metapath- P -guided neighbors of node v
\mathbf{x}_v	Raw feature vector of node v
\mathbf{h}_v	Hidden states (embedding) of node v
$\sigma(\cdot)$	Activation function
\odot	Element-wise multiplication
$ \cdot $	The cardinality of a set

3. Preliminary

This section gives formal definitions of some basic terminologies of heterogeneous graphs. Their graphical illustrations are provided in Figure 2. Besides, Table 1 briefly references the mathematical symbols used in this work.

Definition 1 (Heterogeneous Graph). A heterogeneous graph is defined as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ associated with a node type mapping function $\phi : \mathcal{V} \rightarrow \mathcal{A}$ and an edge type mapping function $\psi : \mathcal{E} \rightarrow \mathcal{R}$. \mathcal{A} and \mathcal{R} denote the pre-defined sets of node types and edge types, respectively, with $|\mathcal{A}| + |\mathcal{R}| > 2$.

Definition 2 (Metapath). A metapath P is defined as a path in the form of $A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_K} A_{K+1}$ (abbreviated as $A_1 A_2 \dots A_{K+1}$), describing a composite relation $R = R_1 \circ R_2 \circ \dots \circ R_K$ between node types A_1 and A_{K+1} , where \circ denotes the composition operator on edge types.

Definition 3 (Metapath Instance). Given a metapath $P = A_1 \xrightarrow{R_1} A_2 \xrightarrow{R_2} \dots \xrightarrow{R_K} A_{K+1}$ of a heterogeneous graph G , a metapath instance of P is defined as a path $v_1 \xrightarrow{e_1} v_2 \xrightarrow{e_2} \dots \xrightarrow{e_K} v_{K+1}$ in graph G , such that $\phi(v_i) = A_i$ and $\psi(e_i) = R_i$ for all i .

Definition 4 (Metapath-guided Neighbor). Given a metapath P of a heterogeneous graph \mathcal{G} , the metapath-guided neighbors \mathcal{N}_v^P of a node v is defined as the set of nodes that connect with node v via some metapath instances of P .

Table 2: A summary of metapath-based HGNNs following our framework.

Model	LSR	MGC	MF
HAN	MN	GAT	Attention
MAGNN	MI	GAT+RotatE	Attention
GTN	MN	GCN	Weighted Sum
MECCH	MC	Graph Pooling	1-D Convolution

Definition 5 (Metapath Context). Given a metapath P of a heterogeneous graph \mathcal{G} , the metapath context of a node v is the local subgraph $\mathcal{S}_v^P = (\mathcal{V}_v^P, \mathcal{E}_v^P)$ centered at v , such that edge $(v_1, v_2) \in \mathcal{E}_v^P$ if and only if it belongs to a metapath instance of P starting from v in the original graph \mathcal{G} .

Definition 6 (Heterogeneous Graph Representation Learning). Given a heterogeneous graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, with node feature matrices $\mathbf{X}_{A_i} \in \mathbb{R}^{|\mathcal{V}_{A_i}| \times d_{A_i}}$ for node types $A_i \in \mathcal{A}$, heterogeneous graph representation learning is to learn the d -dimensional node vectors $\mathbf{h}_v \in \mathbb{R}^d$ for all $v \in \mathcal{V}$ with $d \ll |\mathcal{V}|$ that are able to capture rich structural and semantic information involved in \mathcal{G} .

4. Methodology

In this section, we propose MECCH, a novel metapath context convolution-based HGNN. First in Section 4.1, we describe a general framework that can represent most existing metapath-based HGNNs. Then in Section 4.2, 4.3, and 4.4, we present the three novel components that constitute MECCH. In the end, we also introduce the training losses in Section 4.5 and analyze the model complexity in Section 4.6.

4.1. Metapath-based HGNN Framework

We formalize a metapath-based HGNN as the following four components applied in series: (1) *feature preprocessing*, (2) *local subgraph re-construction (LSR)*, (3) *metapath-specific graph convolution (MGC)*, and (4) *metapath fusion (MF)*. Existing metapath-based HGNNs differ in the last three components. We summarize them along with MECCH in terms of the LSR, MGC, and MF components in Table 2, where MN denotes metapath-guided neighbors, MI denotes metapath instances, and MC denotes metapath contexts.

4.1.1. Feature Preprocessing

In heterogeneous graphs, different node types may have unequal dimensions of node features. Therefore, before feeding the input graph into an HGNN, one would apply a node-type-specific linear transformation to project node features into latent vectors with the same dimensionality. For nodes without raw features, we can assign one-hot vectors to them as initial features. Combined with the linear transformation, this is equivalent to an embedding lookup. The projected representation h_v^0 of node v of type $\phi(v)$ can be obtained by:

$$\mathbf{h}_v^0 = \mathbf{W}_{\phi(v)} \mathbf{x}_v + \mathbf{b}_{\phi(v)}, \quad (1)$$

where $\mathbf{x}_v \in \mathbb{R}^{d_{\phi(v)}}$ is the raw feature, $\mathbf{W}_{\phi(v)} \in \mathbb{R}^{d \times d_{\phi(v)}}$ and $\mathbf{b}_{\phi(v)} \in \mathbb{R}^d$ are the learnable weights and bias.

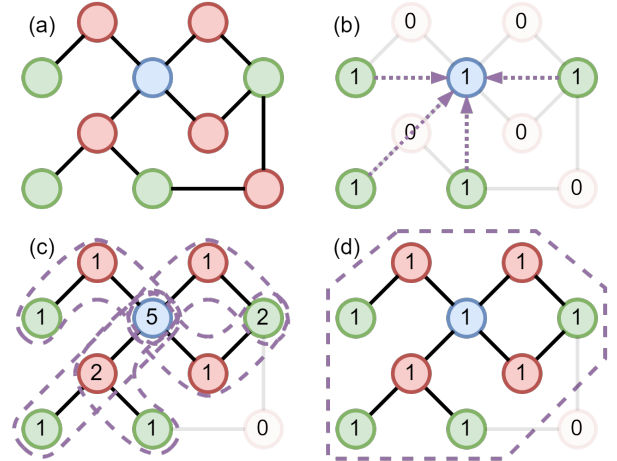


Figure 3: HGNN aggregation for the blue target node using Blue-Red-Green metapath. Numbers indicate how many times the nodes are aggregated. (a) A sample heterogeneous graph. (b) HAN & GTN: metapath-guided neighbors. (c) MAGNN: metapath instances. (d) MECCH (ours): metapath context.

4.1.2. Local Subgraph Re-construction

In metapath-based HGNNs, nodes from the input graph need to be re-connected for message passing through metapath instead of from direct neighbors. Following this, a new local subgraph $\mathcal{S}_v^P = (\mathcal{V}_v^P, \mathcal{E}_v^P)$ would be re-constructed for each node in each related metapath. The definition of \mathcal{S}_v^P differs across HGNNs. For example, \mathcal{S}_v^P in HAN is the metapath- P -guided neighborhood of v , while \mathcal{S}_v^P in MAGNN is the set of metapath- P 's instances incident to v . Each \mathcal{S}_v^P carries unique semantics that will be encoded and captured in the next step.

4.1.3. Metapath-specific Graph Convolution

To compute metapath- P -specific node representations, the model would conduct graph convolution on the re-constructed subgraphs \mathcal{S}_v^P from the last step. For node v at layer l , the metapath- P -specific node representation $\mathbf{h}_{v,P}^l$ is computed by:

$$\mathbf{h}_{v,P}^l = \text{MGC}_P^l(\mathcal{S}_v^P, \{\mathbf{h}_u^{l-1}, \forall u \in \mathcal{V}_v^P\}), \quad (2)$$

where MGC_P^l is the graph convolution function for metapath P at layer l . Different models could have different designs for this function. For example, HAN ignores the intermediate nodes and directly applies GAT to aggregate metapath-guided neighbors. MAGNN incorporates intermediate nodes by encoding and aggregating metapath instances. GTN essentially applies GCN to the metapath-guided neighbors for automatically generated metapaths. A graphical illustration of different MGC_P^l implementations is provided in Figure 3. Basically, the existing metapath-based HGNNs either aggregate deficient nodes (which causes information loss), or aggregate redundant nodes (which leads to high computation costs). Our MECCH introduced later aims to extract adequate information from the graph without any redundancy. After applying this component, every node is now associated with a set of vector representations, each from one metapath.

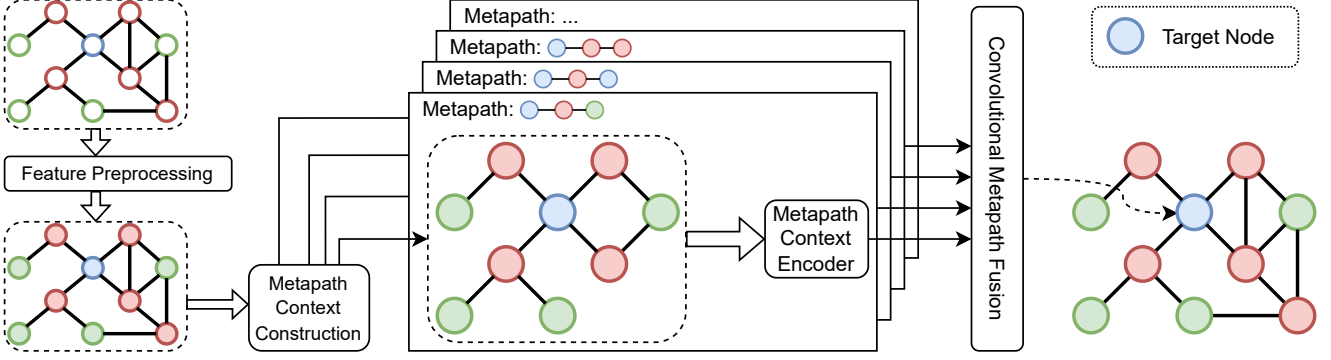


Figure 4: The overall architecture of MECCH. To obtain the representation of the blue target node, after feature preprocessing, MECCH constructs the metapath contexts using metapaths starting from the blue node type. Then, MECCH encodes each metapath context into a vector representation. Finally, MECCH fuses the context representations into one node representation via 1-D convolution.

4.1.4. Metapath Fusion

To combine unilateral metapath-specific representations into one unified node representation \mathbf{h}_v^l , metapath-based HGNNs require a fusion step $\text{MF}_{\phi(v)}^l$:

$$\mathbf{h}_v^l = \text{MF}_{\phi(v)}^l \left(\left\{ \mathbf{h}_{v,P}^l, \forall P \in \mathcal{P}_{\phi(v)} \right\} \right), \quad (3)$$

where $\mathcal{P}_{\phi(v)}$ is the set of metapaths starting from node type $\phi(v)$, either manually selected or automatically generated.

Our proposed MECCH follows this framework with novel designs on the last three components. Specifically, after feature preprocessing, we propose to construct *metapath contexts* for each node in each metapath, then utilize a *metapath context encoder* to embed metapath-specific information, and finally leverage *convolutional metapath fusion* to integrate information from different metapaths together. Figure 4 illustrates the overall architecture of MECCH. The model pseudocode is provided in Algorithm 1.

4.2. Metapath Context Construction

After feature preprocessing, for each node v in each metapath $P \in \mathcal{P}_{\phi(v)}$, MECCH constructs a metapath context, a local subgraph \mathcal{S}_v^P centered at v based on Definition 5. The metapath context \mathcal{S}_v^P describes the entire structure connecting the target node v and the metapath- P -guided neighbors. This notion avoids information loss caused by discarding intermediate nodes in HAN, and reduces redundant computations caused by duplicate nodes across metapath instances in MAGNN. The concept of metapath context is illustrated in Figure 3. As shown, by using metapath contexts, MECCH aggregates each node within the receptive field exactly once, thus can achieve improved efficiency with optimal performance.

For the metapaths considered, MECCH uses all metapaths of a pre-defined length K , which does not require a manual selection process based on human expertise. Any metapath of length less than K must be a sub-sequence of some length- K metapath. The shorter metapath’s context would be subsumed entirely as a subgraph of a longer metapath’s context. Considering any metapath shorter than K would lead to extracting redundant information. Therefore, unlike GTN, which uses all

metapaths with length $\leq K$, MECCH only considers length- K metapaths. In this way, MECCH requires no prior knowledge to select metapaths and at the same time avoids a significant amount of redundant computations.

4.3. Metapath Context Encoder

After constructing the metapath context \mathcal{S}_v^P , MECCH employs a metapath context encoder to learn structural and semantic information embedded in \mathcal{S}_v^P . Since metapath contexts are essentially graphs, it is natural to leverage graph pooling functions to encode \mathcal{S}_v^P . Therefore, the metapath context encoder MGC_P^l of MECCH is formulated as:

$$\mathbf{h}_{v,P}^l = \text{READOUT}_P^l \left(\mathcal{S}_v^P, \left\{ \mathbf{h}_u^{l-1}, \forall u \in \mathcal{S}_v^P \right\} \right), \quad (4)$$

where READOUT_P^l could be a permutation-invariant function like summation, or a more sophisticated GNN-based graph encoder to consider the graph structure. Each $\mathbf{h}_{v,P}^l$ obtained can be interpreted as the information of node v exhibited by metapath P . Our experiments show that an attention-based encoder cannot well capture the semantics embedded in the metapath context. Eventually, we employ the mean graph pooling strategy for this component, which is highly efficient without the need to tune any parameters.

4.4. Convolutional Metapath Fusion

After encoding metapath contexts, MECCH applies a fusion layer to combine representations from different metapaths $\left\{ \mathbf{h}_{v,P}^l, \forall P \in \mathcal{P}_{\phi(v)} \right\}$ into one unified node representation. A straightforward solution is to take the element-wise mean. But the mean operation does not consider that metapaths might have different importance to node representations. HAN and MAGNN compute metapath weights by applying attention to metapath summary vectors calculated by averaging projected representations of all nodes, which restricts the model to full-batch training only. To bypass this limitation, MECCH adopts an efficient 1-D convolution kernel to adaptively learn the impact of each metapath on the downstream task:

$$\mathbf{h}_{v,\mathcal{P}_{\phi(v)}}^l = \sum_{P \in \mathcal{P}_{\phi(v)}} \alpha_P^l \odot \mathbf{h}_{v,P}^l, \quad (5)$$

Algorithm 1 MECCH forward propagation.

Input: Heterogeneous graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with node type mapping $\phi : \mathcal{V} \rightarrow \mathcal{A}$ and edge type mapping $\psi : \mathcal{E} \rightarrow \mathcal{R}$; Node features $\{\mathbf{x}_v, \forall v \in \mathcal{V}\}$; Metapath sets $\{\mathcal{P}_A, \forall A \in \mathcal{A}\}$; Number of layers L

Output: Node representations

```
1:  $\mathbf{h}_v^0 = \mathbf{W}_{\phi(v)}\mathbf{x}_v + \mathbf{b}_{\phi(v)}, \forall v \in \mathcal{V}$  ▷ Feature preprocessing
2: for all  $A \in \mathcal{A}$  do
3:   for all  $P \in \mathcal{P}_A$  do
4:     Construct  $\mathcal{S}_v^P = (\mathcal{V}_v^P, \mathcal{E}_v^P), \forall v \in \mathcal{V}_A$  ▷ 4.2
5:   end for
6: end for
7: for  $l = 1 \dots L$  do
8:   for all  $A \in \mathcal{A}$  do
9:     for all  $P \in \mathcal{P}_A$  do
10:     $\mathbf{h}_{v,P}^l = \text{READOUT}_P^l(\mathcal{S}_v^P, \{\mathbf{h}_u^{l-1}, \forall u \in \mathcal{S}_v^P\})$ , ▷ 4.3
11:     $\forall v \in \mathcal{V}_A$ 
12:   end for
13:    $\mathbf{h}_{v,\mathcal{P}_A}^l = \sum_{P \in \mathcal{P}_A} \alpha_P^l \odot \mathbf{h}_{v,P}^l, \forall v \in \mathcal{V}_A$  ▷ 4.4
14:    $\mathbf{h}_v^l = \sigma(\mathbf{W}_A^l \mathbf{h}_{v,\mathcal{P}_A}^l + \mathbf{b}_A^l), \forall v \in \mathcal{V}_A$  ▷ 4.4
15: end for
16: return  $\{\mathbf{h}_v^L, \forall v \in \mathcal{V}\}$ 
```

where α_P^l is a learnable vector of length equal to the hidden dimension d . Then MECCH further applies a projection to obtain the final node vector of the desired output dimension:

$$\mathbf{h}_v^l = \sigma(\mathbf{W}_{\phi(v)}^l \mathbf{h}_{v,\mathcal{P}_{\phi(v)}}^l + \mathbf{b}_{\phi(v)}^l), \quad (6)$$

where $\mathbf{W}_{\phi(v)}^l$ and $\mathbf{b}_{\phi(v)}^l$ are the weight matrix and the bias term, σ is an activation function (ReLU by default).

4.5. Training Losses

For the node classification task, the loss function would be the cross-entropy between ground truths and model outputs over a small set of labeled nodes:

$$\mathcal{L}_{nc} = -\frac{1}{|\mathcal{V}_{label}|} \sum_{v \in \mathcal{V}_{label}} \sum_{c=1}^C \mathbf{y}_v[c] \cdot \log \hat{\mathbf{y}}_v[c], \quad (7)$$

where \mathcal{V}_{label} is the set of labeled nodes for training, C is the number of classes, \mathbf{y}_v is the one-hot ground truth of node v , and $\hat{\mathbf{y}}_v = \mathbf{h}_v^L$ is the model's predicted probability vector.

For the link prediction task, we leverage the negative sampling strategy (Mikolov et al., 2013b) to train the model along with a DistMult (Yang et al., 2015) link decoder using a binary cross-entropy loss function:

$$\mathcal{L}_{lp} = -\frac{1}{|\Omega|} \sum_{(u,v) \in \Omega} \left(\log \sigma(\mathbf{h}_u^T \mathbf{W}_{lp} \mathbf{h}_v) + \mathbb{E}_{v^- \sim p_n} \left[\log \sigma(-\mathbf{h}_u^T \mathbf{W}_{lp} \mathbf{h}_{v^-}) \right] \right), \quad (8)$$

where $\mathbf{W}_{lp} = \text{diag}(\mathbf{w}_{lp})$ is a learnable diagonal matrix, σ is the sigmoid function, Ω is the set of positive (observed) edges, and p_n is the distribution of v^- for sampling negative edges (u, v^-) .

4.6. Complexity Analysis

The time complexity of MECCH consists of three parts, corresponding to the feature preprocessing, the metapath context encoder, and the convolutional metapath fusion, respectively. The metapath context construction does not contribute to the complexity because it can be pre-computed. To simplify the complexity analysis, we consider the following scenario: there are N nodes in the input graph, every node is associated with M metapaths of length K , and every node has t type- A neighbors for each node type A . The feature preprocessing is done by Eq. (1), which requires $O(N_A \cdot d \cdot d_A)$ for each node type A . If we assume every node type's input features have the same dimensionality d_m . Then this part requires $O(N \cdot d \cdot d_m)$ in total. The metapath context encoder is defined by Eq. (4), where \mathcal{S}_v^P contains at most $1 + t + t^2 + \dots + t^K = \frac{t^{K+1}-1}{t-1}$ nodes. Hence this part requires $O(N \cdot M \cdot \frac{t^{K+1}-1}{t-1} \cdot d) = O(N \cdot M \cdot t^K \cdot d)$. The convolutional metapath fusion consists of Eq. (5) and Eq. (6). Eq. (5) requires $O(N \cdot M \cdot d)$, while Eq. (6) requires $O(N \cdot d \cdot d_{out})$.

In practice, the input dimensionality d_{in} and the output dimensionality d_{out} are typically very small compared to $M \cdot t^K$. Thus, the overall time complexity of a 1-layer MECCH is $O(N \cdot M \cdot t^K \cdot d)$. In comparison, MAGNN requires $O(N \cdot M \cdot K \cdot t^K \cdot d)$, GTN requires $O(N^3 \cdot K)$, and HAN requires $O(N \cdot M \cdot t^K \cdot d)$. Therefore, we have $\text{MECCH} \approx \text{HAN} < \text{MAGNN} < \text{GTN}$ for computational complexity.

5. Experiments

In this section, we demonstrate the effectiveness and efficiency of our proposed MECCH by conducting experiments in the node classification and link prediction tasks on five heterogeneous graph datasets. The experiments are designed to answer the following research questions.

- RQ1. How does MECCH perform in classifying nodes?
- RQ2. How does MECCH perform in predicting links?
- RQ3. How efficient is MECCH?
- RQ4. How effective is each proposed component?

5.1. Experimental Setups

5.1.1. Datasets

We select datasets of varying sizes from different domains to simulate real-life applications. For *node classification*, we adopt the IMDB, ACM, and DBLP datasets from (Yun et al., 2019) to evaluate the performance of MECCH and baselines. For *link prediction*, we choose the LastFM dataset from (Fu et al., 2020) and the PubMed dataset from (Yang et al., 2022) to compare the models. Since previous researchers usually preprocess the raw data under different pipelines, scores reported in their original papers cannot be compared directly. To ensure fairness, we reran experiments of each model on the same pre-processed datasets. We try our best to follow the default settings (e.g., training/validation/testing splits) of these datasets. Statistics of the five adopted datasets are summarized in Table 3.

Table 3: Statistics of datasets.

Dataset	# Nodes (Types)	# Edges (Types)	# Feats	Target	# Training	# Validation	# Testing
IMDB	12,722 (3)	37,288 (4)	1,256	Movie	300 (10%)	300 (10%)	2,339 (80%)
ACM	8,994 (3)	25,922 (4)	1,902	Paper	600 (20%)	300 (10%)	2,125 (70%)
DBLP	18,405 (3)	67,946 (4)	334	Author	800 (20%)	400 (10%)	2,857 (70%)
LastFM	20,612 (3)	201,908 (5)	-	User-Artist	64,984 (70%)	9,283 (10%)	18,567 (20%)
PubMed	63,109 (4)	368,245 (16)	200	Disease-Disease	29,845 (70%)	4,264 (10%)	8,528 (20%)

- **IMDB** is a movie information network for *node classification*. The preprocessed subset contains three types of nodes: movies, actors, and directors. A bag-of-words representation of plot keywords describes each node in the dataset. The movie nodes are labeled by their genres.
- **ACM** is an academic network for *node classification*, with three types of nodes: papers, authors, and subjects. Nodes in ACM are associated with bag-of-words representations of keywords. The paper nodes are labeled based on the conference in which they were published.
- **DBLP** is also an academic network for *node classification*. It contains three types of nodes: papers, authors, and conferences. Nodes in the dataset are associated with bag-of-words representations of keywords. The author nodes are labeled based on their research fields.
- **LastFm** is a music listening information network for *link prediction*. The preprocessed subset contains three types of nodes: users, artists, and tags. Nodes in this dataset are not associated with any features. The prediction targets are the user-artist edges.
- **PubMed** is a biomedical knowledge graph for *link prediction*, with four types of nodes: genes, diseases, chemicals, and species. Nodes in PubMed are associated with aggregated word2vec embeddings. The prediction targets are the disease-disease edges.
- **HAN** (Wang et al., 2019a) is a *metapath-based HGNN*. It learns metapath-specific node embeddings by aggregating metapath-guided neighbors and leverages the attention mechanism to combine metapaths.
- **MAGNN** (Fu et al., 2020) is a *metapath-based HGNN* that enhances HAN by incorporating the intermediate nodes along metapaths with metapath instance encoders.
- **GTN** (Yun et al., 2019) is a *metapath-based HGNN* that can automatically discover valuable metapaths via multiplying weighted sums of relation subgraphs.

5.1.3. Hyperparameters

For fairness, we configure the node embedding dimension of all HGNNs to 64. For models with multi-head attention, we set the number of heads as 8 and ensure it would not multiply the embedding dimension, as suggested in (Hu et al., 2020). For metapaths used by HAN and MAGNN, we adopt the ones suggested in their papers. For MAGNN, we use the relational rotation metapath instance encoder. For GTN, we use the adaptive learning rate as suggested in the authors’ implementation. For MECCH, we choose the metapath length K within $\{1, 2, 3, 4, 5\}$. For each model-dataset pair, we experiment with the number of HGNN layers within $\{1, 2, 3, 4, 5\}$, the dropout rate within $\{0, 0.5\}$, the learning rate within $\{1, 2, 5\} \times \{10^{-3}, 10^{-2}\}$, the weight decay within $\{0, 10^{-3}\}$. The hyperparameter combinations that can fit in the GPU memory and yields the best results are selected (except for GTN, which runs on CPU for ACM, DBLP, and LastFM). We use the same data splits for training, validation, and testing across models. We train all HGNNs using the Adam optimizer for 500 epochs with early stopping applied after a patience of 50 epochs.

5.1.4. Platform Specifications

All of our experiments are conducted on the same platform with the following system specifications. **CPU**: 40x Intel Xeon Platinum 8268 cores. **GPU**: 1x NVIDIA TITAN V card. **Memory**: 775 GiB. **Operating System**: Ubuntu 18.04.6 LTS. **Libraries**: DGL 0.7.2, PyTorch 1.10.1, and CUDA 11.3.

5.2. Node Classification (RQ1)

We evaluate MECCH against other baselines for node classification on IMDB, ACM, and DBLP. As shown in Table 3, only a small set of labeled nodes are available for training and validation, which is a typical semi-supervised node classification setting for GNNs. The averages and standard deviations of

5.1.2. Baselines

We compare MECCH against state-of-the-art HGNNs, including three *relation-based HGNNs* (RGCN, HGT, and Simple-HGN), three *metapath-based HGNNs* (HAN, MAGNN, and GTN).

- **RGCN** (Schlichtkrull et al., 2018) is a *relation-based HGNN* that groups and aggregates neighbors based on their edge types with edge-type-specific weights.
- **HGT** (Hu et al., 2020) is a *relation-based HGNN*. It leverages a Transformer-like self-attention architecture parameterized by node types and edge types to capture the importance weights of the heterogeneous neighborhood.
- **Simple-HGN** (Lv et al., 2021) is a *relation-based HGNN*. It extends GAT by adding edge-type attention, residual connection, and output normalization.

Table 4: Experiment results (%) of the node classification task on the IMDB, ACM, and DBLP datasets. \uparrow (\downarrow) indicates the higher (lower), the better. For each column, the best result is in **bold** font, and the second-best result is underlined.

Model	IMDB		ACM		DBLP	
	Macro-F1 \uparrow	Micro-F1 \uparrow	Macro-F1 \uparrow	Micro-F1 \uparrow	Macro-F1 \uparrow	Micro-F1 \uparrow
RGCN	56.85 \pm 0.41	58.59 \pm 0.42	<u>92.17\pm0.28</u>	<u>92.07\pm0.27</u>	92.46 \pm 1.40	93.41 \pm 1.17
HGT	57.32 \pm 1.04	59.02 \pm 1.33	<u>90.55\pm0.63</u>	<u>90.52\pm0.63</u>	93.96 \pm 0.24	94.76 \pm 0.24
Simple-HGN	<u>58.70\pm0.63</u>	60.27 \pm 0.85	89.36 \pm 3.12	89.25 \pm 3.17	92.60 \pm 0.53	93.40 \pm 0.52
HAN	<u>58.54\pm1.67</u>	59.60 \pm 1.80	91.67 \pm 0.18	91.55 \pm 0.18	92.29 \pm 0.35	93.16 \pm 0.34
MAGNN	57.87 \pm 1.29	59.78 \pm 1.65	90.63 \pm 0.27	90.58 \pm 0.26	<u>94.14\pm0.23</u>	<u>94.86\pm0.23</u>
GTN	58.66 \pm 1.30	<u>60.28\pm1.74</u>	91.94 \pm 0.57	91.83 \pm 0.57	<u>93.25\pm0.25</u>	<u>94.12\pm0.21</u>
MECCH	62.59\pm1.96	64.62\pm2.38	92.74\pm0.40	92.67\pm0.41	94.34\pm0.29	95.08\pm0.25

Table 5: Experiment results (%) of the link prediction task on the LastFM and PubMed datasets.

Model	LastFM	PubMed
RGCN	82.85 \pm 0.61	68.22 \pm 1.26
HGT	82.36 \pm 0.29	73.24 \pm 1.46
Simple-HGN	<u>83.02\pm0.91</u>	<u>73.74\pm0.50</u>
HAN	82.26 \pm 0.42	71.50 \pm 0.84
MAGNN	79.26 \pm 0.89	65.74 \pm 2.35
GTN	76.94 \pm 0.09	-
MECCH	84.27\pm0.05	76.69\pm0.24

Macro-F1 and Micro-F1 scores from 5 runs of each model on each dataset are reported in Table 4. From the table, we have the following observations:

1. MECCH consistently outperforms previous state-of-the-art HGNNs in node classification across different datasets. On average, the performance gain of our model over the best baseline is around 2.59%.
2. Generally speaking, metapath-based HGNNs produce higher and more stable results than relation-based HGNNs. But the margin is tiny. Nevertheless, our MECCH proves that properly leveraging metapath contexts can bring noticeable advantages.
3. It is interesting that MAGNN and GTN can hardly benefit from their high computation costs. One possible reason might be that redundant aggregations overly smooth meaningful signals from important nodes. For MAGNN, another possible culprit is the parameter-intensive attention mechanism, which is removed in MECCH to avoid overfitting and enhance generalization ability.

5.3. Link Prediction (RQ2)

On LastFM and PubMed, link prediction is formulated as a binary classification problem to distinguish whether a given edge exists or not in the original graph. During training, the negative edge (u, v^-) is generated by replacing the destination node v of a positive edge (u, v) by uniformly sampling v^- from $\mathcal{V}_{\phi(v)}$. For validation and testing, negative edges are pre-generated in a 1:1 ratio to the positive ones and are kept the same across

Table 6: Efficiency comparisons of metapath-based HGNNs on ACM and DBLP. TPE: time per epoch (s); #E: number of epochs; RAM: CPU memory usage (GiB); VRAM: GPU memory usage (GiB).

Dataset	Model	TPE \downarrow	#E \downarrow	RAM \downarrow	VRAM \downarrow
ACM	HAN	<u>0.062</u>	<u>21.6</u>	0.855	<u>2.649</u>
	MAGNN	0.151	45.0	4.967	5.205
	GTN	24.769	24.4	8.342	-
	MECCH	0.049	4.0	<u>0.905</u>	1.727
DBLP	HAN	0.068	41.6	3.671	2.937
	MAGNN	(7.022)	(22.2)	(4.265)	(1.569)
	GTN	253.875	<u>21.4</u>	55.734	-
	MECCH	<u>0.181</u>	4.2	<u>4.142</u>	<u>3.645</u>

different models for a fair comparison. We also follow (Lv et al., 2021) to sample hard negative edges for validation and testing. The averages and standard deviations of ROC-AUC scores in 5 runs are reported in Table 5, from which we have the following observations:

1. MECCH performs consistently better than other models, with 1.51% and 4.00% improvements over the strongest baselines on LastFM and PubMed, respectively.
2. The small ROC-AUC variances show that MECCH is much more stable than other HGNNs.
3. One interesting observation is that relation-based HGNNs usually achieve better scores than metapath-based HGNNs in the link prediction task. Nevertheless, our MECCH still outperforms all of them.
4. On PubMed, MAGNN has poor performance, and GTN cannot run due to significant memory consumption. They could both suffer from the diluted neighborhood messages caused by aggregating duplicate information. Moreover, MAGNN may also be negatively affected by the over-parameterized attention modules.

5.4. Time and Memory Usage (RQ3)

Besides comparing HGNNs based on prediction accuracy, we also evaluate models in terms of computational efficiency. Table 6 records the time spent per epoch, the number of epochs

Table 7: Ablation study for node classification (IMDB, ACM, and DBLP) and link prediction (LastFM and PubMed).

Model	IMDB		ACM		DBLP		LastFM	PubMed
	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	ROC-AUC	ROC-AUC
MECCH-KHop	53.53	56.28	91.32	91.26	93.74	94.52	83.29	72.74
MECCH-HAN	60.03	62.25	91.55	91.48	92.67	93.58	83.34	73.42
MECCH-MAGNN	59.42	61.44	90.59	90.51	94.05	94.71	80.97	70.92
MECCH-ACE	62.70	64.22	90.74	90.68	92.42	93.30	82.44	75.25
MECCH-MMF	61.90	64.12	92.36	92.32	94.14	94.89	84.22	75.76
MECCH-nAMF	61.58	63.23	91.86	91.76	94.01	94.67	83.10	75.38
MECCH-gAMF	62.17	64.25	92.38	92.30	94.39	94.86	84.15	76.26
MECCH	62.59	64.62	92.74	92.67	94.34	95.08	84.27	76.69

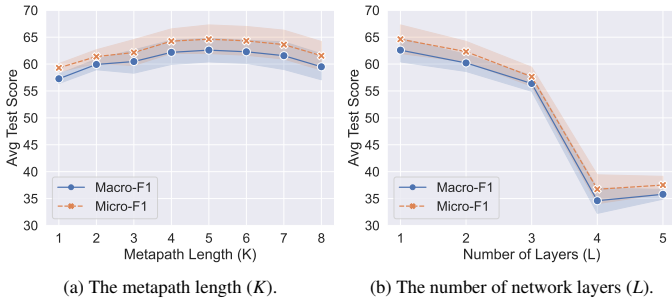


Figure 5: MECCH performance on the IMDB dataset with varying metapath lengths (K) and model depths (L).

used (excluding the patience epochs), and CPU/GPU memory consumption of MECCH and other metapath-based HGNNs. All models except MAGNN on DBLP run in full-batch.

As shown, **MECCH is much more efficient than MAGNN and GTN** in training time and memory usage, with a better or comparable resource usage compared to HAN. This is accomplished when MECCH incorporates more information than HAN (i.e., the intermediate nodes along metapaths). MAGNN on DBLP can only run in mini-batches due to its high GPU memory usage. Hence its reported values cannot be compared directly. On both ACM and DBLP, GTN cannot fit in the limited GPU memory. Instead, it runs on CPU with a significantly longer running time. These results show that MECCH can achieve both superior model performance and decent computational efficiency, further proving the advantages of leveraging metapath contexts over metapath-guided neighbors and metapath instances.

5.5. Ablation Study (RQ4)

To analyze the effectiveness of the three proposed components, we conduct an ablation study by comparing our MECCH with the following model variants: (1) **MECCH-KHop** is a variant without leveraging metapath contexts, equivalent to aggregating all the nodes within a radius of K regardless of node or edge types; (2) **MECCH-HAN** is a variant replacing metapath contexts with metapath-guided neighbors used in HAN; (3) **MECCH-MAGNN** is a variant replacing metapath contexts with metapath instances used in MAGNN; (4) **MECCH-**

ACE is a variant using an attention-based module as the metapath context encoder; (5) **MECCH-MMF** is a variant using the element-wise mean as the metapath fusion strategy; (6) **MECCH-nAMF** is a variant that fuses metapaths by applying the attention mechanism per node; (7) **MECCH-gAMF** is a variant that fuses metapaths by using the attention scores computed globally based on the summaries of all nodes' representation. The experiment results are reported in Table 7, from which we have the following observations.

1. By comparing MECCH against MECCH-KHop, MECCH-HAN, and MECCH-MAGNN, we find that MECCH can achieve a considerable performance improvement, especially on the IMDB dataset. This helps to validate the superiority of leveraging metapath contexts, which is the core idea of MECCH.
2. By comparing MECCH against MECCH-ACE, it is interesting that a simple mean graph pooling is more effective for encoding metapath contexts than an attention-based module in most cases. Besides model performance, mean graph pooling as the metapath context encoder is much more efficient than other options.
3. The difference between the results of MECCH and MECCH-MMF suggests that fusing metapaths adaptively with a 1-D convolution kernel is generally a better choice than fusing with the element-wise mean. This tells us that metapaths have different importance levels concerning the downstream task, and it is beneficial to model them.
4. The results of MECCH-nAMF and MECCH-gAMF suggest that the attention mechanism is not suitable for metapath fusion (MF). The node-wise attention MF in MECCH-nAMF results in different metapath weights among nodes, which is counter-intuitive and potentially jeopardizes the performance. Despite having decent performance, the global attention MF in MECCH-gAMF constrains the model to only train and make inferences in full batch, because all nodes' representations are required to calculate the summary vector of each metapath. Thus, 1-D convolution is a better choice for metapath fusion, considering its simplicity, effectiveness, and efficiency.
5. Using the same neighborhood subgraph definition, MECCH-MAGNN in Table 7 outperforms or matches the original

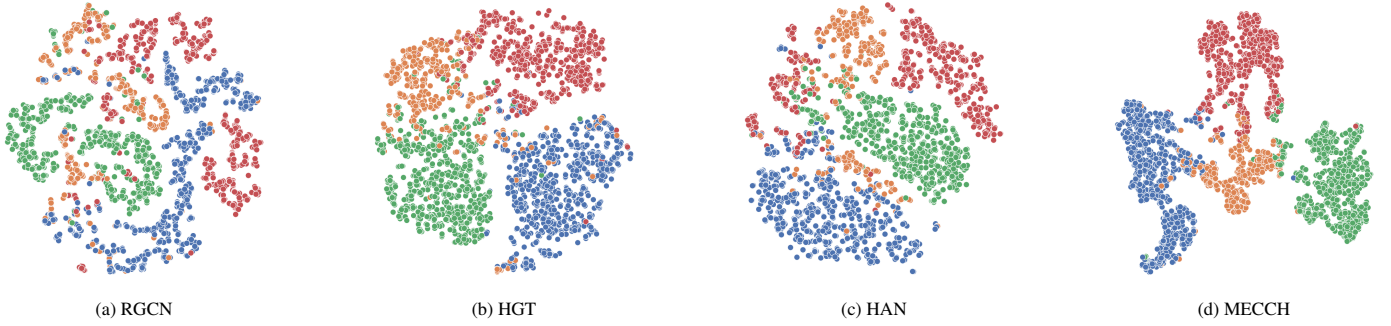


Figure 6: Embedding visualization on DBLP. Each point is an author where color indicates the class label.

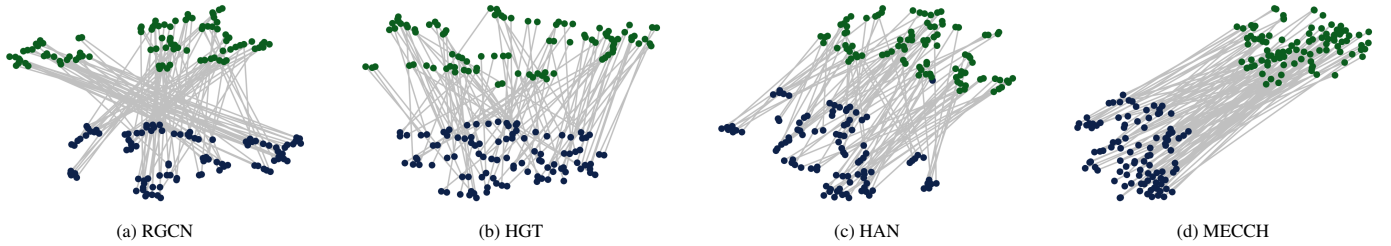


Figure 7: Embedding visualization on LastFM. Blue and green points are users and artists, respectively.

MAGNN model in Table 4 and Table 5. This also indicates that the attention-free design of MECCH components improves the model performance.

5.6. Hyperparameter Study

To investigate how MECCH is affected by the hyperparameters, we conduct pilot tests on the IMDB dataset with different choices for the metapath length K and the number of MECCH layers L . Results are shown in Figure 5, where Figure 5 (a) shows results when L is fixed at 1, and Figure 5 (b) shows results under $K = 5$. Except for the concerning hyperparameter, all other setups are kept unchanged.

From Figure 5 (a), we observe significant increase in performance when leveraging larger metapath contexts ($K > 1$). The best results on IMDB is reached when K is around 5. Further increasing the metapath length K has diminishing returns. Obviously, the optimal choice of K is closely related to the structural characteristics of the dataset and the downstream task. Referring to Figure 5 (b), MECCH is sensitive to the choice of L . This makes sense because the metapath length K is fixed at 5. Adding one more MECCH layer introduces a much larger receptive field that could dramatically impact the final performance. This also suggests that MECCH is capable of learning the complex semantics embedded in the graph with very few layers, avoiding the potential performance degradation issues when the model gets deep.

5.7. Visualization

To provide an intuitive assessment of the model’s embedding quality, we visualize the node representations generated by different HGNNs on the DBLP and LastFM datasets. On DBLP, we take the generated node vectors of the authors in the testing

set just before applying the last GNN layer/output projection. On LastFM, we randomly sample 100 positive user-artist pairs from the testing set and take the node vectors just before applying the DistMult decoder. Then they are projected into a two-dimensional vector space using t-SNE. Here we plot the t-SNE node embedding visualizations of R-GCN, HGT, HAN, and MECCH in Figure 6 for DBLP, where each color corresponds to one author class (research area), and in Figure 7 for LastFM, where blue points are users and green points are artists.

From the two figures, one can clearly observe the advantage of MECCH over other HGNNs. Figure 6 shows that our proposed MECCH can generate more informative node representations in the node classification task. On DBLP, MECCH makes nodes with the same class labels close to each other and nodes with different class labels far away, resulting in a clear separation between classes. Other HGNNs would produce dispersed intra-class distribution and obscure inter-class boundaries. Figure 7 shows that MECCH can well separate user and artist nodes and preserve the aligned correlation of the positive user-artist pairs. While other HGNNs can roughly divide user nodes and artist nodes into two different groups, they cannot maintain good alignment between users and artists.

6. Conclusion

In this paper, we formulate a unified framework for metapath-based HGNNs to facilitate model understanding and new model design. We analyze existing models and observe that they cannot balance model performance and computational efficiency. To address these issues, we propose a metapath context convolution based HGNN (MECCH) with three novel components: (1) metapath context construction, (2) metapath context encoder, and (3) convolutional metapath fusion. The three components

are designed to prevent information loss and redundant computations simultaneously.

In experiments, MECCH consistently outperforms state-of-the-art baselines on five widely adopted heterogeneous graph datasets for node classification and link prediction, with improved computational efficiency. Ablation studies validate the effectiveness of MECCH’s novel components in boosting model performance. Embedding visualizations qualitatively demonstrate the superiority of MECCH.

For future work, we plan to improve MECCH so that it can be applied to web-scale datasets, such as AMiner (Tang et al., 2008) and MAG (Sinha et al., 2015), where existing metapath-based HGNNs can hardly run due to substantial computational costs. We also plan to incorporate various training paradigms by using MECCH as an encoder to further boost the prediction accuracy, such as self-supervised learning (Ren et al., 2019; Park et al., 2020), contrastive learning (Wang et al., 2021b; Qian et al., 2022; Chen et al., 2023a), and federated learning (Fu and King, 2023).

Acknowledgement

The work described in this paper was partially supported by the National Key Research and Development Program of China, 2018AAA0100204 and the RGC General Research Fund (GRF), 2151185 (CUHK 14222922).

References

Ahn, H., Yang, Y., Gan, Q., Moon, T., Wipf, D., 2022. Descent steps of a relation-aware energy produce heterogeneous graph neural networks, in: *NeurIPS*.

Bordes, A., Usunier, N., Garcia-Durán, A., Weston, J., Yakhnenko, O., 2013. Translating embeddings for modeling multi-relational data, in: *NIPS*, pp. 2787–2795.

Chang, Y., Chen, C., Hu, W., Zheng, Z., Zhou, X., Chen, S., 2022. Megnn: Meta-path extracted graph neural network for heterogeneous graph representation learning. *Knowl. Based Syst.* 235, 107611.

Chen, H., Yin, H., Wang, W., Wang, H., Nguyen, Q.V.H., Li, X., 2018. PME: projected metric embedding on heterogeneous networks for link prediction, in: *KDD*, ACM. pp. 1177–1186.

Chen, M., Huang, C., Xia, L., Wei, W., Xu, Y., Luo, R., 2023a. Heterogeneous graph contrastive learning for recommendation, in: *WSDM*, ACM. pp. 544–552.

Chen, Y., Truong, Q.T., Shen, X., Wang, M., Li, J., Chan, J., King, I., 2023b. Topological representation learning for e-commerce shopping behaviors, in: *Proceedings of the 19th International Workshop on Mining and Learning with Graphs (MLG)*.

Chen, Y., Zhang, Y., Yang, M., Song, Z., Ma, C., King, I., 2023c. WSFE: wasserstein sub-graph feature encoder for effective user segmentation in collaborative filtering, in: *SIGIR*, ACM. pp. 2521–2525.

Dong, Y., Chawla, N.V., Swami, A., 2017. metapath2vec: Scalable representation learning for heterogeneous networks, in: *KDD*, ACM. pp. 135–144.

Fu, T., Lee, W., Lei, Z., 2017. Hin2vec: Explore meta-paths in heterogeneous information networks for representation learning, in: *CIKM*, ACM. pp. 1797–1806.

Fu, X., King, I., 2023. Fedhgn: A federated framework for heterogeneous graph neural networks, in: *IJCAI*, ijcai.org. pp. 3705–3713.

Fu, X., Zhang, J., Meng, Z., King, I., 2020. MAGNN: metapath aggregated graph neural network for heterogeneous graph embedding, in: *WWW*, ACM / IW3C2. pp. 2331–2341.

Hamilton, W.L., Ying, Z., Leskovec, J., 2017. Inductive representation learning on large graphs, in: *NIPS*, pp. 1024–1034.

Hong, H., Guo, H., Lin, Y., Yang, X., Li, Z., Ye, J., 2020. An attention-based graph neural network for heterogeneous structural learning, in: *AAAI*, AAAI Press. pp. 4132–4139.

Hu, Z., Dong, Y., Wang, K., Sun, Y., 2020. Heterogeneous graph transformer, in: *WWW*, ACM / IW3C2. pp. 2704–2710.

Kipf, T.N., Welling, M., 2017. Semi-supervised classification with graph convolutional networks, in: *ICLR (Poster)*, OpenReview.net.

Li, Q., Han, Z., Wu, X., 2018. Deeper insights into graph convolutional networks for semi-supervised learning, in: *AAAI*, AAAI Press. pp. 3538–3545.

Liang, L., Xu, Z., Song, Z., King, I., Ye, J., 2022. Resnorm: Tackling long-tailed degree distribution issue in graph neural networks via normalization. *CoRR* abs/2206.08181.

Lv, Q., Ding, M., Liu, Q., Chen, Y., Feng, W., He, S., Zhou, C., Jiang, J., Dong, Y., Tang, J., 2021. Are we really making much progress?: Revisiting, benchmarking and refining heterogeneous graph neural networks, in: *KDD*, ACM. pp. 1150–1160.

Ma, Y., Song, Z., Hu, X., Li, J., Zhang, Y., King, I., 2023. Graph component contrastive learning for concept relatedness estimation, in: *AAAI*, AAAI Press. pp. 13362–13370.

Meng, Z., Li, Y., Zhao, P., Yu, Y., King, I., 2023a. Meta-learning with motif-based task augmentation for few-shot molecular property prediction, in: *SDM*, SIAM. pp. 811–819.

Meng, Z., Zhao, P., Yu, Y., King, I., 2023b. Doubly stochastic graph-based non-autoregressive reaction prediction, in: *IJCAI*, ijcai.org. pp. 4064–4072.

Meng, Z., Zhao, P., Yu, Y., King, I., 2023c. A unified view of deep learning for reaction and retrosynthesis prediction: Current status and future challenges, in: *IJCAI*, ijcai.org. pp. 6723–6731.

Mikolov, T., Chen, K., Corrado, G., Dean, J., 2013a. Efficient estimation of word representations in vector space, in: *ICLR (Workshop Poster)*.

Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J., 2013b. Distributed representations of words and phrases and their compositionality, in: *NIPS*, pp. 3111–3119.

Park, C., Kim, D., Han, J., Yu, H., 2020. Unsupervised attributed multiplex network embedding, in: *AAAI*, AAAI Press. pp. 5371–5378.

Perozzi, B., Al-Rfou, R., Skiena, S., 2014. Deepwalk: online learning of social representations, in: *KDD*, ACM. pp. 701–710.

Qian, Y., Zhang, Y., Wen, Q., Ye, Y., Zhang, C., 2022. Rep2vec: Repository embedding via heterogeneous graph adversarial contrastive learning, in: *KDD*, ACM. pp. 1390–1400.

Ren, Y., Liu, B., Huang, C., Dai, P., Bo, L., Zhang, J., 2019. Heterogeneous deep graph infomax. *CoRR* abs/1911.08538.

Schlichtkrull, M.S., Kipf, T.N., Bloem, P., van den Berg, R., Titov, I., Welling, M., 2018. Modeling relational data with graph convolutional networks, in: *ESWC*, Springer. pp. 593–607.

Shi, C., Hu, B., Zhao, W.X., Yu, P.S., 2019. Heterogeneous information network embedding for recommendation. *IEEE Trans. Knowl. Data Eng.* 31, 357–370.

Sinha, A., Shen, Z., Song, Y., Ma, H., Eide, D., Hsu, B.P., Wang, K., 2015. An overview of microsoft academic service (MAS) and applications, in: *WWW (Companion Volume)*, ACM. pp. 243–246.

Song, Z., King, I., 2022. Hierarchical heterogeneous graph attention network for syntax-aware summarization, in: *AAAI*, AAAI Press. pp. 11340–11348.

Sun, Z., Deng, Z., Nie, J., Tang, J., 2019. Rotate: Knowledge graph embedding by relational rotation in complex space, in: *ICLR (Poster)*, OpenReview.net.

Tang, J., Zhang, J., Yao, L., Li, J., Zhang, L., Su, Z., 2008. Arnetminer: extraction and mining of academic social networks, in: *KDD*, ACM. pp. 990–998.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need, in: *NIPS*, pp. 5998–6008.

Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y., 2018. Graph attention networks, in: *ICLR (Poster)*, OpenReview.net.

Wang, C., Zhou, S., Yu, K., Chen, D., Li, B., Feng, Y., Chen, C., 2022. Collaborative knowledge distillation for heterogeneous information network embedding, in: *WWW*, ACM. pp. 1631–1639.

Wang, L., Gao, C., Huang, C., Liu, R., Ma, W., Vosoughi, S., 2021a. Embedding heterogeneous networks into hyperbolic space without meta-path, in: *AAAI*, AAAI Press. pp. 10147–10155.

Wang, X., Ji, H., Shi, C., Wang, B., Ye, Y., Cui, P., Yu, P.S., 2019a. Heterogeneous graph attention network, in: *WWW*, ACM. pp. 2022–2032.

Wang, X., Liu, N., Han, H., Shi, C., 2021b. Self-supervised heterogeneous graph neural network with co-contrastive learning, in: *KDD*, ACM. pp.

1726–1736.

- Wang, X., Zhang, Y., Shi, C., 2019b. Hyperbolic heterogeneous information network embedding, in: AAAI, AAAI Press. pp. 5337–5344.
- Yang, B., Yih, W., He, X., Gao, J., Deng, L., 2015. Embedding entities and relations for learning and inference in knowledge bases, in: ICLR (Poster).
- Yang, C., Xiao, Y., Zhang, Y., Sun, Y., Han, J., 2022. Heterogeneous network representation learning: A unified framework with survey and benchmark. *IEEE Trans. Knowl. Data Eng.* 34, 4854–4873.
- Yang, M., Zhou, M., Pan, L., King, I., 2023a. *khgcn*: Tree-likeness modeling via continuous and discrete curvature learning, in: KDD, ACM. pp. 2965–2977.
- Yang, M., Zhou, M., Xiong, H., King, I., 2023b. Hyperbolic temporal network embedding. *IEEE Trans. Knowl. Data Eng.* 35, 11489–11502.
- Yang, M., Zhou, M., Ying, R., Chen, Y., King, I., 2023c. Hyperbolic representation learning: Revisiting and advancing, in: ICML, PMLR. pp. 39639–39659.
- Yang, Y., Guan, Z., Li, J., Zhao, W., Cui, J., Wang, Q., 2023d. Interpretable and efficient heterogeneous graph convolutional network. *IEEE Trans. Knowl. Data Eng.* 35, 1637–1650.
- Yang, Y., Liu, T., Wang, Y., Zhou, J., Gan, Q., Wei, Z., Zhang, Z., Huang, Z., Wipf, D., 2021. Graph neural networks inspired by classical iterative algorithms, in: ICML, PMLR. pp. 11773–11783.
- Yoon, M., Palowitch, J., Zelle, D., Hu, Z., Salakhutdinov, R., Perozzi, B., 2022. Zero-shot transfer learning within a heterogeneous graph via knowledge transfer networks, in: NeurIPS.
- Yu, L., Shen, J., Li, J., Lerer, A., 2020. Scalable graph neural networks for heterogeneous graphs. *CoRR abs/2011.09679*.
- Yu, L., Sun, L., Du, B., Liu, C., Lv, W., Xiong, H., 2023. Heterogeneous graph representation learning with relation awareness. *IEEE Trans. Knowl. Data Eng.* 35, 5935–5947.
- Yun, S., Jeong, M., Kim, R., Kang, J., Kim, H.J., 2019. Graph transformer networks, in: NeurIPS, pp. 11960–11970.
- Zhang, C., Song, D., Huang, C., Swami, A., Chawla, N.V., 2019a. Heterogeneous graph neural network, in: KDD, ACM. pp. 793–803.
- Zhang, J., Shi, X., Xie, J., Ma, H., King, I., Yeung, D., 2018. GaAN: Gated attention networks for learning on large and spatiotemporal graphs, in: UAI, AUAI Press. pp. 339–349.
- Zhang, J., Shi, X., Zhao, S., King, I., 2019b. STAR-GCN: stacked and reconstructed graph convolutional networks for recommender systems, in: IJCAI, ijcai.org. pp. 4264–4270.
- Zhang, M., Wang, X., Zhu, M., Shi, C., Zhang, Z., Zhou, J., 2022. Robust heterogeneous graph neural networks against adversarial attacks, in: AAAI, AAAI Press. pp. 4363–4370.
- Zhao, J., Wang, X., Shi, C., Liu, Z., Ye, Y., 2020. Network schema preserving heterogeneous information network embedding, in: IJCAI, ijcai.org. pp. 1366–1372.
- Zhou, K., Dong, Y., Wang, K., Lee, W.S., Hooi, B., Xu, H., Feng, J., 2021. Understanding and resolving performance degradation in deep graph convolutional networks, in: CIKM, ACM. pp. 2728–2737.