SYNHING: SYNTHETIC HETEROGENEOUS INFORMA TION NETWORK GENERATION FOR GRAPH LEARNING AND EXPLANATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph Neural Networks (GNNs) excel in modeling graph structures across diverse domains, such as community analysis and recommendation systems. As the need for GNN interpretability grows, there is an increasing demand for robust baselines and comprehensive graph datasets, especially within the realm of Heterogeneous Information Networks (HIN). To address this, we introduce Syn-HING, a framework for Synthetic Heterogeneous Information Network Generation designed to advance graph learning and explanation. After identifying key motifs in a target HIN, SynHING systematically employs a bottom-up generation process with intra-cluster and inter-cluster merge modules. This process, along with post-pruning techniques, ensures that the synthetic HIN accurately mirrors the structural and statistical properties of the original graph. The effectiveness of SynHING is validated using four datasets - IMDB, Recipe, ACM, and DBLP spanning three distinct application categories, demonstrating both its generality and practicality. Furthermore, SynHING provides ground-truth motifs for evaluating GNN explainer models, establishing a new benchmark for explainable, synthetic HIN generation. This contributes significantly to advancing interpretable machine learning in complex network environments.

028 029 030

031

025

026

027

006

008 009 010

011

013

014

015

016

017

018

019

021

1 INTRODUCTION

032 In recent years, Graph Neural Networks (GNNs) have shown impressive performance in various 033 graph analysis tasks such as community analysis (Shchur & Günnemann, 2019), chemical bond 034 analysis (Stokes et al., 2020), and recommendation systems (Cui et al., 2020). These tasks include node and edge classification, link prediction, and clustering (Chami et al., 2022). Heterogeneous Information Networks (HINs) consist of multiple types of nodes and edges that contain various information, providing a natural representation of real-world data. The development of Heterogeneous 037 Graph Neural Networks (HGNN) has been sparked by HINs. These networks can be classified into different types of models, including meta-path-based models like HAN (Wang et al., 2019) and MAGNN (Fu et al., 2020), transformer-based GNN models (Yun et al., 2020; Hu et al., 2020b), as 040 well as SimpleHGN (Lv et al., 2021) which uses projection layers to map information to a shared 041 feature space and aggregates information using an edge-type attention mechanism. Additionally, 042 TreeXGNN (Hong et al., 2023) combines a tree-based feature extractor with the HGNN model to 043 enhance performance. However, due to the lack of public HIN datasets, most existing HGNNs are 044 trained and evaluated on only a few known public datasets such as IMDB¹, ACM (Wang et al., 2019), and DBLP 2 , leading to potential bias. The scarcity of public HIN datasets compared to other 046 machine learning domains poses significant challenges for HGNNs, potentially causing overfitting 047 and hindering effective generalization (Palowitch et al., 2022).

Due to the scarcity of HIN datasets, it is even more difficult to study trustworthy and interpretable
 models. As trust, transparency, and privacy become essential for machine learning models, it is important to reveal the decision-making process and knowledge hidden behind models. Most existing
 GNN models lack transparency, making them difficult to be trusted and limiting their applicability in

¹https://www.kaggle.com/datasets/karrrimba/movie-metadatacsv
²http://web.cs.ucla.edu/~yzsun/data/

decision-critical scenarios. GNN explainer models were proposed to disclose the black box. Never theless, there is no objective way to evaluate the performances of GNN explainers as there exists no
 suitable dataset in public domains that provides the data, objective, and corresponding ground-truth
 explanations at the same time.

058 Synthetic homogeneous graph datasets have been introduced (Ying et al., 2019), attempting to al-059 leviate the lack of explainable datasets. These datasets are created by randomly attaching specially 060 designed structured network motifs to artificial graphs (Albert & Barabási, 2002). These motifs are 061 then utilized as ground-truth explanations. There have been very few attempts to generalize syn-062 thetic graph generation to HINs. This generalization is non-trivial for two reasons: Firstly, the basic 063 artificially structured motifs, such as houses and grids (Ying et al., 2019), common in the synthesis of 064 homogeneous graphs, are inadequate for HINs. These simplistic motifs often fail to capture the complexity and diversity of real-world heterogeneous graphs, resulting in a poor representation of the 065 intricate structures found in actual HIN environments. Secondly, methods designed for generating 066 homogeneous synthetic graphs cannot be directly applied to heterogeneous information networks 067 (HINs) due to their unique structural constraints. Traditional approaches, which often involve ran-068 domly connecting nodes to form edges, may result in illegitimate connections within HINs. Such 069 connections can violate the semantic rules inherent to HINs, where edges must align with the types of nodes they connect to represent accurate and meaningful relationships. Moreover, these methods 071 randomly add edges, merely increasing node degrees without considering the global structure of 072 the graph. This randomization disrupts the alignment of node degree and structural patterns with 073 real-world HINs. Additionally, as highlighted by Li et al. (2023a), the scarcity of heterogeneous 074 datasets with ground-truth explanations further complicates the development of explanation models 075 for HINs. Researchers often resort to indirect evaluation methods, which may not effectively capture the true performance of these models. 076

077 In light of the above concerns, we introduce Synthetic Heterogeneous Information Network Generation (SynHING), a novel framework that constructs synthetic HINs for graph learning and ex-079 planation by referencing existing real graphs. SynHING methodically generates major motifs es-080 sential for explanations and employs our newly developed Intra-/Inter-cluster Merge method. This 081 method facilitates the merging of multiple subgraph groups to systematically create synthetic HINs of any specified size. Additionally, we introduce the concept of exclusion that allows us to flexibly adjust the complexity of node prediction tasks within these synthetic graphs while maintaining 083 structural similarity to the reference graph. This feature enhances the utility of SynHING in practi-084 cal applications. The synthetic graphs generated by SynHING provide interpretable insights, aiding 085 significantly in the explanation tasks associated with HINs. SynHING has been validated using 086 four datasets: IMDB, Recipe, ACM, and DBLP, which cover three distinct application categories, 087 demonstrating its generality and practicality. Furthermore, theoretical complexity analysis has been 088 conducted to showcase its scalability.³

089 090 091

092

094

2 RELATED WORK

2.1 SYNTHETIC GRAPH GENERATION

Artificially synthesized data has a long history of development (Kingma & Welling, 2013; Bowyer 096 et al., 2011; Dong et al., 2018; Frid-Adar et al., 2018; Karras et al., 2019; Xu et al., 2019; Figueira 097 & Vaz, 2022). With the growth of Graph Neural Networks (GNN), there has been a renewed interest 098 in synthetic graph generation algorithms. Early on, Snijders & Nowicki (1997) introduced a method to generate graph edges based on node clusters. More contemporary approaches, such as those 099 described by Dwivedi et al. (2020), extract subgraphs from real-world graphs to test GNNs' ability 100 to identify specific substructures within Stochastic Block Model (SBM) graphs. The concept of 101 SBM, particularly popular for generating clusters that maintain strong intra-node correlations, has 102 been adapted into various forms, including unsupervised (Tsitsulin et al., 2020) and semi-supervised 103 models (Rozemberczki et al., 2021). Further extending SBM's utility, GraphWorld (Palowitch et al., 104 2022) leverages the Degree-Corrected Stochastic Block Model (DC-SBM) (Abbe, 2017) to create 105 diverse graph datasets using multiple parameters. Yet, these works mainly focus on homogeneous 106

³Open-source code will be released upon acceptance.

graphs, and most synthetic graph generation methods do not provide the ability for interpretation,
 lack ground-truth explanations, and cannot be directly extended to HINs.

110 111

112

2.2 EXPLAINER FOR GRAPH NEURAL NETWORKS

113 Developing trustworthy machine learning models is now a widely acknowledged goal within the 114 community. Explainability methods for Graph Neural Networks (GNNs) have seen considerable development, especially for node and graph classification tasks (Ying et al., 2019; Luo et al., 2020; 115 Yuan et al., 2021; Lin et al., 2022). These methods generally fall into two categories: inherent inter-116 pretable models and post-hoc explainability approaches (Yuan et al., 2020). Inherent interpretable 117 models, such as ProtGNN (Dai & Wang, 2021), integrate explanations directly within the model, 118 using mechanisms like top-K similarity to identify influential subgraphs for their predictions. In 119 contrast, post-hoc explainability methods focus on identifying crucial subgraphs (Luo et al., 2020; 120 Yuan et al., 2021) and key features (Ying et al., 2019), approximating the behavior of black-box 121 models by elucidating the connections between inputs and outputs. While most data validating 122 these explainer models are derived from homogeneous graphs with simplistic motifs (Ying et al., 123 2019), such as houses and grids, these often lack the diversity and complexity of real-world graphs. 124 Some researchers are studying heterogeneous GNN explanations (Li et al., 2023); Lv et al., 2023), 125 but the lack of HINs with explanation ground truths poses challenges. Consequently, generating appropriate heterogeneous graph datasets for validation is emerging as a crucial research area. 126

- 127
- 128 129

2.3 DATASETS WITH GROUND-TRUTH EXPLANATIONS

The conventional evaluation of graph explanatory methods often relies on molecular datasets or 130 synthetic community node classification datasets. These datasets are favored because they offer 131 ground-truth explanations. For instance, MUTAG is a molecular dataset containing graphs labeled 132 according to their mutagenic effect (Debnath et al., 1991). Synthetic datasets, as introduced by Ying 133 et al. (2019), are node classification datasets created by randomly attaching structured network mo-134 tifs to base graphs. These motifs are designed with specific structures such as houses, cycles, and 135 grids. The base graphs are generated either through BA methods or as a balanced binary tree. How-136 ever, these datasets are homogeneous and do not transition well to heterogeneous graph contexts. 137 In contrast, our approach involves extracting relevant motifs directly from actual HINs to serve as 138 the basis for ground-truth explanations and employing a systematic bottom-up method to guarantee 139 inherent semantic rules of HINs when generating HIN datasets.

140 141

142 143

144

152 153

154

3 PROPOSED METHOD: SYNHING

3.1 PRELIMINARIES

145 HINs, also called heterogeneous graphs, consist of multiple node and edge types, which can be 146 defined as a graph $G = (\mathcal{V}, \mathcal{E}, \Phi, \Psi)$, where \mathcal{V} and \mathcal{E} is the set of nodes and edges, respectively. 147 Each node $v \in \mathcal{V}$ has a type $\Phi(v) \in \mathcal{T}_v$, and each edge $e \in \mathcal{E}$ has a type $\Psi(e) \in \mathcal{T}_e$. The node 148 feature matrix is denoted by $F_{\phi} \in \mathbb{R}^{|\mathcal{V}^{\phi}| \times d_{\phi}}$, where \mathcal{V}^{ϕ} is the set of node with node type ϕ , i.e. 149 $\mathcal{V}^{\phi} = \{ v \in \mathcal{V} \mid \Phi(v) = \phi \}$, and d_{ϕ} is the feature dimension of the node type ϕ . Target nodes of 150 the graph G, denoted by \mathcal{V}^{ϕ_0} , are associated with labels collected as $Y \in \mathcal{Y}^{|\mathcal{V}^{\phi_0}|}$, where $\phi_0 \in \mathcal{T}_v$ 151 denotes the target node type.

3.2 OVERVIEW OF SYNHING





- 171
- 172

173 In this study, we proposed SynHING, a novel synthetic HIN generation framework. We aim to 174 generate an arbitrary size synthetic heterogeneous graph G with the explanation ground truths, which 175 closely mimics the property of the given real-world graph G through a bottom-up generation process, 176 which is demonstrated in Fig. 1. Firstly, we generate the major motifs and derive base subgraphs 177 from them with proper node degree distributions. Secondly, we introduce two merge modules to 178 handle partial graphs; Intra-Cluster Merge aims to merge subgraphs within a cluster, and Inter-179 Cluster Merge is used to merge clusters with different labels. Finally, the synthetic HIN will be generated after feature generation and post-pruning. SynHING framework consists of six modules, 181 as shown in Fig. 2: (1) Major Motif Generation, (2) Base Subgraph Generation, (3) Intra-Cluster Merge, (4) Inter-Cluster Merge, (5) Node Feature Generation, and (6) Post-Pruning. The details of 182 the proposed methods will be introduced in the following sections. 183

185 3.3 MAJOR MOTIF GENERATION (MMG)

186 To generate major motifs for ground-truth explanations, we identify meta-paths within the refer-187 enced graph that originate and terminate with target nodes, with all intermediate nodes. The major 188 motif is derived by designating two target nodes as anchors and connecting them through all possible 189 meta-paths within a specified number of hops (Wang et al., 2019; Fu et al., 2020). The maximum 190 number of hops can be set manually or based on the number of layers in the GNN models, reflect-191 ing that GNN computation graphs are represented by *n*-hop subgraphs, where *n* matches the model 192 layers. As depicted in Fig. 2, the MMG module utilizes three one-hop paths to form a major motif 193 as an example.

Further investigation revealed that these derived motifs are common patterns found using the graphlet searching method (Milo et al., 2002) in real-world graphs. We conducted experiments on the IMDB dataset, identifying one of the most common graphlets, *G*20 (Milo et al., 2002), which features multiple minor nodes acting as bridges between two target nodes. A similar pattern is noted in Megnn (Chang et al., 2022). We define these robust graph patterns as the major motif, providing essential ground truths for explanations. In addition, the motif can also be defined by the user and customized for diverse explanation tasks.

201 202

3.4 BASE SUBGRAPH GENERATION (BSG)

Based on the major motifs previously generated, we further develop base subgraphs. In this base
subgraph generation process, we introduce randomness into the major motifs and augment them
with several non-target nodes, designated as minor nodes, attached to target nodes within each motif.
This addition aims to create diverse subgraphs with noise, which are not part of the ground truths
for explanations but mimic the real-world reference graph.

These minor nodes fulfill two primary functions. Firstly, they help match the degree distributions of the target nodes in the subgraphs to those observed in the referenced real-world distribution, denoted as $P^{\phi}(k)$, where k is the number of connections to nodes of type ϕ . Secondly, the minor nodes serve as crucial junction points for subsequent merging processes.

213 Once minor nodes are added, each pair of target nodes within a motif is assigned identical labels, 214 defining this entity as a *base subgraph*. This is denoted as a tuple (S_i, y_i) , where $S_i = (\mathcal{V}_i, \mathcal{E}_i)$ 215 represents the structure of the *i*-th subgraph, and $y_i \in \mathcal{Y}$ is the label of the associated target nodes. The generated subgraphs are then collected into sets \mathcal{K}_y for each label $y \in \mathcal{Y}$.

216 3.5 Merge to Generate HINS217

Conventional methods for constructing graphs often involve adding edges between nodes or subgraphs to create a connected homogeneous graph. However, this approach carries the risk of inadvertently forming illegal connections despite careful node selection to prevent them. It is also challenging to maintain statistical properties, such as node degree distributions relative to different node types, when connecting nodes.

To overcome these challenges, we propose a novel *Merge* operation that combines two nodes into one, which will connect back to the initial neighbors of the two nodes. This method elegantly adheres to existing constraints and preserves the degree distributions of the target nodes within the generated subgraphs, ensuring an accurate representation of the target real-world graph. The general merge function operates as follows: Given a graph $G = (\mathcal{V}, \mathcal{E})$ and two nodes $v_1, v_2 \in \mathcal{V}$. If we merge v_1 and v_2 into v_1 in G, then the process is defined by

$$\mathcal{V}', \mathcal{E}') = \mathsf{Merge}(v_1, v_2; G) \tag{1}$$

$$= (\mathcal{V} \setminus \{v_2\}, \mathcal{E} \cup \{(v_1, u) \mid u \in N(v_2), u \neq v_1\} \setminus \{(v_2, u) \mid u \in N(v_2)\}), \quad (2)$$

where N(v) represents the neighbors of the node v. The Merge operation connects the neighbors of node v_2 to node v_1 while simultaneously removing the merged node v_2 and its original edges to its neighbors. We use Merge($\mathfrak{P}; G$) to denote merging multiple pairs in $\mathfrak{P} \subseteq \mathcal{V} \times \mathcal{V}$ in G (note that the order in \mathfrak{P} does not matter). Similarly, Merge($\mathfrak{P}; G_1 \oplus G_2 \oplus \ldots$) signifies the merging of pairs across multiple graphs G_1, G_2, \ldots , where \oplus denotes the graph disjoint union operator. Next, we generate a complete synthetic HIN from bottom to top with Intra-/Inter-Cluster Merges.

237 238

229 230

3.5.1 INTRA-CLUSTER MERGE (INTRA-CM)

239 Intra-CM is designed to merge the subgraphs with identical labels one by one to form a cluster 240 denoted by $C_y = (\mathcal{V}_y, \mathcal{E}_y)$, which can mirror the "Superstar" phenomenon often observed in com-241 munity networks (Albert & Barabási, 2002; Abbe, 2017). In social networks, early adopters often 242 become central nodes or "Superstars" due to the sequential nature of network growth. As new mem-243 bers join, they tend to connect with already well-established individuals. Early joiners accumulate 244 more connections over time, benefiting from the principle of preferential attachment, where new 245 links are more likely to form with highly connected members. This process leads to a few early 246 adopters amassing a disproportionate number of connections, thereby becoming key influencers or 247 opinion leaders within the network.

248 In brief, for each label y, we denote \mathcal{K}_y as the set of base subgraphs of the same label generated 249 earlier in the BSG step. These base subgraphs are sequentially merged into the cluster C_y . The entire 250 Intra-CM process is conducted $|\mathcal{Y}|$ times to produce all clusters $\{C_y \mid y \in \mathcal{Y}\}$ for all labels. Note 251 that we form each cluster independently. Since different node types need to be handled carefully in 252 HINs, Intra-CM is conducted separately for each node type, denoted by ϕ . Specifically, the initial subgraph S_0 is chosen from K_y to be the original cluster C_y^0 . At each iterations *i*, we select a 253 subgraph $S_i = (\mathcal{V}_i, \mathcal{E}_i)$ from the remaining $\mathcal{K}_y \setminus \{S_0, ..., S_{i-1}\}$ and merge C_y^{i-1} and S_i to generate 254 255 C_{y}^{i} . For each minor node type $\phi \neq \phi_{0}$, we initially determine the number of *Merge* operations to be 256 performed, denoted as n_{intra}^{ϕ} . The number n_{intra}^{ϕ} is sampled from a binomial distribution, i.e. 257

$$n_{\text{intra}}^{\phi} \sim \mathbf{B}(n = |\mathcal{V}_i^{\phi}|, p = p^{\phi}), \tag{3}$$

where p^{ϕ} is the Intra-CM probability for the minor node type ϕ . A higher Intra-CM probability leads to more nodes being merged during this step, resulting in a tighter connection graph within the cluster, which means increased exclusion between clusters. To merge nodes within S_i and C_y , we need to identify the sampling space of the node pairs:

258

$$\mathfrak{M}_{\text{intra}}^{\phi} = \left\{ \left\{ v_y, v_i \right\} \mid v_y \in \mathcal{V}_y^{\phi}, v_i \in \mathcal{V}_i^{\phi} \right\}.$$
(4)

Subsequently, n_{intra}^{ϕ} pairs are then sampled uniformly from $\mathfrak{M}_{intra}^{\phi}$ into $\mathfrak{P}^{\phi} \subseteq \mathfrak{M}_{intra}^{\phi}$ without replacement. The *Merge* operation is performed on the sampled pairs of all minor node types to merge the S_i into the cluster:

268
269
$$C_y^i = \text{Merge}\left(\bigcup_{\phi \in \mathcal{T}_v, \phi \neq \phi_0} \mathfrak{P}^{\phi}; \quad C_y^{i-1} \oplus S_i\right), \tag{5}$$



Figure 3: Intra-Cluster and Inter-Cluster Merges

where \bigcup denotes the union over all minor node types. After the above Intra-CM, we can generate multiple clusters with all labels, shown in Fig. 3a.

282 3.5.2 INTER-CLUSTER MERGE (INTER-CM)

277 278 279

280

281

283

284

285

287

288

292 293 294

295

301

302

303

314

Inter-CM aims to merge clusters with different labels. Unlike intra-clusters, the "Superstar" phenomenon should not occur in inter-clusters because clusters would not appear sequentially to form the whole graph. Hence, we should concurrently merge all clusters rather than sequentially merge subgraphs, as in Intra-CM. In short, the proposed Inter-CM merges node pairs from different clusters to form the complete graph structure \tilde{G} , using clusters { $C_y | y \in \mathcal{Y}$ } produced by Intra-CM.

Analogous to Intra-CM, all merges need to be conducted separately for each node type ϕ . For each node type *phi*, the initial step in Inter-CM involves identifying potential pairs $\mathfrak{M}_{inter}^{\phi}$ to be merged, formally defined as follows:

$$\mathfrak{M}_{\text{inter}}^{\phi} = \left\{ \left\{ v_1, v_2 \right\} \mid v_1 \in \mathcal{V}_{y_1}^{\phi}, \, v_2 \in \mathcal{V}_{y_2}^{\phi}, \, \left\{ y_1, y_2 \right\} \subseteq \mathcal{Y}, y_1 \neq y_2 \right\},\tag{6}$$

where $\mathcal{V}_{y_1}^{\phi}, \mathcal{V}_{y_2}^{\phi}$ are nodes in type ϕ of two different clusters C_{y_1}, C_{y_2} , respectively. The number of pairs n_{inter}^{ϕ} is sampled from a binomial distribution:

$$n_{\text{inter}}^{\phi} \sim B\left(n = \sum_{y \in \mathcal{Y}} |\mathcal{V}_{y}^{\phi}|, k = q^{\phi}\right),\tag{7}$$

where q^{ϕ} is the Inter-Cluster merge probability. A higher value of q^{ϕ} results in more nodes from different clusters being merged, leading to a graph with lower exclusion of clusters. The n_{inter}^{ϕ} pairs are sampled uniformly from $\mathfrak{M}_{\text{inter}}^{\phi}$ to form the set of node pairs that we intend to merge \mathfrak{P}^{ϕ} . The clusters are merged based on \mathfrak{P}^{ϕ} and form a complete graph \tilde{G} :

$$\tilde{G} = \operatorname{Merge}\left(\bigcup_{\phi \in \mathcal{T}'_{v}} \mathfrak{P}^{\phi}; \quad \bigoplus_{y \in \mathcal{Y}} C_{y}\right),$$
(8)

where \bigoplus denotes the graph disjoint union over all labels $y \in \mathcal{Y}$. If the generated graph \tilde{G} is intended to be multi-label, $\mathcal{T}'_v = \mathcal{T}_v$, i.e., all node types, including target node type ϕ_0 , are allowed to be merged. Conversely, in the case of single-label, merging target nodes is not permitted, i.e., $\mathcal{T}'_v = \mathcal{T}_v \setminus \{\phi_0\}$. After the above Inter-CM, we can generate a complete heterogeneous graph structure with multiple labels, as shown in Fig. 3b.

315 3.6 NODE FEATURE GENERATION (NFG)

316 For NFG, we sample node features from within-cluster multivariate normal distributions, following 317 previous studies(Palowitch et al., 2022; Tsitsulin et al., 2022). The node features in the same cluster 318 will be sampled from a shared prior multivariate normal distribution $\mathcal{N}(\mu_y, \alpha)$, where μ_y represents 319 the feature center sampled from another normal distribution $\mu_y \sim \mathcal{N}(0,\beta)$. Here, α/β serves as 320 a hyperparameter, representing the ratio of feature center distance to cluster covariance. It can be 321 interpreted as a signal-to-noise ratio (SNR). For the multi-label nodes, we generate node features based on a joint probability distribution that combines multiple independent probability distribution 322 functions. Afterward, we draw samples from this combined distribution to determine the features of 323 the nodes. In this work, we generated node features for target nodes because minor node features are

often unavailable in the real-world heterogeneous graphs and are commonly preprocessed by either
 constants, node IDs, or propagated features (Lv et al., 2021). Therefore, we follow this common
 setting, using node ID or node type information as node features for minor nodes to approximate
 these real-world datasets.

328

330

3.7 POST-PRUNING (P-P)

P-P is an optional yet critical process applied to synthetic graphs to ensure they adhere to constraints
observed in raw data. For instance, in the IMDB dataset, each movie is linked to no more than
three actors, reflecting inherent limits within the original dataset. During P-P, we first establish the
upper limits of node degrees and scan the HIN node-by-node to remove excess edges until the node
degrees conform to the upper limit constraints. Importantly, we prioritize the retention of edges that
form part of the major motif, thereby preserving the integrity of the explanation ground truths.

337 338

339

3.8 COMPLEXITY OF SYNHING

To explore SynHING's scalability, we conducted a complexity analysis. The process of MMG and BSG are both independent. Therefore, the time complexity of MMG and BSG is O(N). The complexity of Intra-CM is $O(N|\mathcal{V}_i| + N|\mathcal{E}_i|)$ or O(N), as $|\mathcal{V}_i|$ and $|\mathcal{E}_i|$ are the number nodes and edges in the base subgraph, which are constant w.r.t. N. The processes involved in Inter-CM are similar to Intra-CM; the complexity is also O(N). Therefore, the overall time complexity of SynHING is determined by the number of motifs N with a time complexity of O(N), which showcases the scalability of SynHING. More details can be found in the Appendix A.1.

346 347 348

349 350

351

361

362

364

366

4 EXPERIMENTAL SETTINGS

4.1 DATASETS AND HGNNS

To assess the SynHING framework, we generate synthetic graphs using four well-established HIN node classification datasets: IMDB 1, Recipe (Majumder et al., 2019), ACM (Wang et al., 2019), and DBLP 2. Fig. 4a presents the graph schema for these four heterogeneous datasets that illustrate the permissible edge types within each graph.

To identify major motifs, we designate two target nodes as anchors and connect them using all feasible meta-paths, limited by a specified number of hops. Specifically, we utilize two hops for the IMDB, Recipe, and ACM datasets and four hops for the DBLP dataset to align with their respective graph schemas, as depicted in Fig. 8. We further study the minor node degree of the approximate reference graph and real graph. More details can be found in the Appendix A.7.





Figure 4: Graph Schema and Major Motifs of the Four Heterogeneous Graph Datasets

We utilize the transductive learning approach for node classification tasks and randomly select 24%
of the target nodes for training, 6% for validation, and 70% for testing (Wang et al., 2019; Lv et al., 2021; Hong et al., 2023). We used three well-known HGNNs, HGT (Hu et al., 2020b), SimpleHGN (Lv et al., 2021), and TreeXGNN (Hong et al., 2023), as our encoders to validate the synthetic HINs.



Figure 5: Visualization of Synthetic IMDB in Different Intra-/Inter-Cluster Probabilities. Dark blue represents minor nodes, and other colors indicate target nodes on different labels.

4.2 EVALUATION METRICS

We repeated all experiments five times and evaluated performance using average Micro-F1 and Macro-F1 for node prediction (Wang et al., 2019; Lv et al., 2021; Hong et al., 2023) and Fidelity for interpretation evaluation (Yuan et al., 2021; Li et al., 2022). See more details in the Appendix A.3.

5 RESULTS AND DISCUSSION

5.1 Cluster Exclusion Can be Adjusted Flexibly to Affect Model Performance

Intra-CM probability p and Inter-CM probability q determine the degree of cluster exclusion of the synthetic HINs. The higher p and the lower q increases cluster exclusion, which produces a purer and better-organized graph structure. Due to the flexibility of our proposed framework, this degree of exclusion can be flexibly adjusted. To evaluate how such an adjustment alters model performance, we benchmark the performance of the HGT on the synthetic IMDB (Syn-IMDB) with different p, qwith p > q. As illustrated in Fig. 5a, the HGT model performs better in terms of Macro-F1 as the p increases, as Syn-IMDB has more tightly connected clusters or a higher cluster exclusion. Conversely, a decrease in q introduces more noise, leading to a lower cluster exclusion and worse performance. We also visualize the synthetic IMDB graphs with three settings, shown in Fig. 5c, 5d, 5e. The above results demonstrated that we can use p and q to control the exclusion of clusters within the generating synthetic HIN and benchmark the ability of HGNN graph learning. This allows us to control graph generation with high flexibility. Similar trends can also be observed for SimpleHGN and TreeXGNN; details can be found in Fig. 9 in the Appendix.

5.2 MAJOR MOTIFS ACHIEVE LOW FIDELITY

Fig. 5b illustrates the variations in fidelity for HGT across different degrees of cluster exclusion. The
megatrend is similar to HGT's Macro-F1 score, with the purer (higher degree of cluster exclusions)
synthetic HINs fidelity score being better (lower). Notably, the fidelity remains stable regardless of
changes in the Intra-CM probability *p*. While higher Intra-CM probabilities may introduce additional information beyond the major motifs, the fidelity of the model does not undergo significant
changes, shown in Fig. 6a. This finding confirms that the major motif designed indeed serves as
the primary cause for accurate predictions of GNN models. Fig. 6b clearly shows that the fidelity



(a) Intra-Cluster Probability (b) Inter-Cluster Probability (c) SNR of Node Features



Table 1: Ablation studies of SynHING: We replace the critical modules of SynHING with random functions to verify the importance. The Random-motifs turn off the MMG module. The Random-Merge entails randomly merging nodes without conducting Intra-/Inter-CM.

	SynI	MDB	Randon	n-Motifs	Randon	n-Merge
	Macro-F1 (%)	Micro-F1 (%)	Macro-F1 (%)	Micro-F1 (%)	Macro-F1 (%)	Micro-F1 (%)
HAN	82.37 ± 0.45	82.42 ± 0.52	77.21 ± 0.69	77.20 ± 0.81	37.31 ± 5.57	41.47 ± 4.40
HGT	87.86 ± 0.30	87.88 ± 0.31	80.99 ± 0.55	80.97 ± 0.58	68.52 ± 3.74	69.11 ± 3.24
SimpleHGN	87.60 ± 0.49	87.66 ± 0.49	83.04 ± 0.48	83.02 ± 0.48	72.77 ± 6.33	72.88 ± 6.30
TreeXGNN	87.68 ± 0.35	87.70 ± 0.36	81.77 ± 1.14	81.73 ± 1.15	68.92 ± 6.57	69.07 ± 6.56

decreases as the Inter-CM probability q decreases while maintaining a fixed Intra-CM probability p. This decrease can be attributed to the increase in the exclusion of clusters, which causes models to primarily learn from the corresponding major motifs without interference from non-informative nodes and edges. Fig. 6c illustrates the relationship between fidelity and the SNR of node features. As the SNR increases, the amount of information associated with all node features also increases. Consequently, when the graph structure and major motif design remain unchanged, the models incorporate more node feature information, thereby improving fidelity. The above experimental results can confirm that the design of major motifs can indeed represent the most crucial graph patterns and 459 can serve as effective ground-truth explanations.

5.3 ABLATION STUDIES OF SYNHING

463 To evaluate the impact of each module in the SynHING framework, we performed ablation studies by removing critical modules one at a time on the IMDB dataset. More specifically, the Random-464 Motifs study turns off the MMG module and randomly generated motifs. The Random-Merge study 465 entails randomly merging nodes without conducting Intra-Cluster and Inter-Cluster Merges. Table 466 1 shows that compared to the original SynIMDB, the performance of Random-Motifs decreases 467 significantly on all HGNNs, HAN, HGT, SimpleHGN, and TreeXGNN (-5.16%, -6.87%, -4.56%, 468 and -5.91% in Macro-F1, respectively). This drop highlights the importance of the MMG mod-469 ule. Furthermore, the performance of Random-Merge declines even further (-45.06%, -19.34%, 470 -14.83%, and -18.76% in Macro-F1 for HAN, HGT, SimpleHGN, and TreeXGNN, respectively), 471 revealing consistent trends across these HGNNs, demonstrating the effectiveness of the proposed 472 Merge method for generating synthetic HINs.

473

438

439 440

441 442

443

444

452

453

454

455

456

457

458

460 461

462

474 475

5.4 HOW SIMILAR IS THE SYNTHETIC GRAPH TO THE ACTUAL GRAPH?

476 To answer this question, we assess high-level similarity by pretraining on the synthetic graph and 477 finetuning on the responding actual graph, inspired by the fact that without careful selection of pretraining tasks, the transfer of knowledge between diverse semantics can lead to negative transfer (Hu 478 et al., 2020a; Rosenstein et al., 2005). Specifically, we evaluate similarity from two perspectives: 479 (i) If the synthetic and reference graphs are similar, we expect to see a positive transfer. (ii) If we 480 maliciously destroy the structure and features of the synthetic graphs, a negative transfer would 481 occur, which could decrease performance. See more implementation details in Appendix A.8. 482

For the positive transfer experiment, we compare the performance of finetuning the model pre-483 trained on the synthetic graph with and without pertaining. The results, shown in Table 2, demon-484 strated pretraining on synthetic HINs significantly improved performance in the four datasets. In 485 the IMDB dataset, we increased Macro-F1 by up to 3% with HGT and 2% with SimpleHGN. On

IMDB, Recipe, and ACM, the standard deviation was significantly reduced, enhancing the model's
 stability and slightly increasing the performance. The consistent trends across two HGNNs on four
 datasets demonstrate solid positive transfer effectiveness.

Table 2: Performance comparison of HGT and SimpleHGN: With and without pretraining on synthetic HINs and fine-tuning on real-world graphs. We use boldface to highlight performance improvements.

		H	<mark>GT</mark>	Simpl	eHGN
Dataset	Pretrained on	Macro-F1	Micro-F1	Macro-F1	Micro-F1
IMDB		63.00 ± 1.19	67.20 ± 0.57	63.53 ± 1.36	67.36 ± 0.5
INDD	Syn-IMDB	66.10 ± 0.21	68.03 ± 0.53	65.52 ± 0.50	68.45 ± 0.53
Pagina	-	57.26 ± 1.84	56.98 ± 2.02	60.29 ± 1.31	60.15 ± 1.4
Recipe	Syn-Recipe	57.82 \pm 0.46	57.83 ± 0.64	60.40 ± 0.22	60.21 ± 0.2
ACM		91.12 ± 0.76	91.00 ± 0.76	93.42 ± 0.44	93.35 ± 0.4
ACM	Syn-ACM	92.55 ± 0.20	92.54 ± 0.21	94.16 ± 0.43	94.11 ± 0.4
		93.01 ± 0.23	93.49 ± 0.25	94.01 ± 0.24	94.46 ± 0.2
DBLP	Syn-DBLP	93.88 ± 0.25	94.35 ± 0.23	94.27 ± 0.58	94.73 ± 0.5

For the negative transfer experiment, malicious graphs are created by: (i) Node shuffling involves row-wise shuffling the adjacency matrix A corresponding to the graph, breaking the homophily of the synthetic graph. (ii) Feature shuffling involves row-wise shuffling the feature matrix F. Table 3 shows that malicious synthetic graphs obviously cause negative transfer. In most cases, the feature shuffling has a greater impact. It is speculated that the mismatch between the node features and the labels causes more damage to the overall message passing. Noted, when the users focus on explainable ground truths, the generated synthetic graphs do not need to resemble the reference graphs. SynHING can freely generate brand-new HINs through user-defined major motifs.

Table 3: Performance comparison of HGT: Pretraining on node shuffled and feature shuffled synthetic HINs and finetuning on real HINs. We use boldface to highlight the lowest score and an underline to indicate the second-lowest.

	Pretrain on SynHING	Macro-F1	Micro-F1
	w/o Shuffled	66.10 ± 0.21	68.03 ± 0.53
IMDB	Node Shuffled	64.54 ± 0.58	67.44 ± 0.59
	Feature Shuffled	$\overline{\textbf{62.06} \pm \textbf{1.28}}$	$\overline{\textbf{63.96}\pm\textbf{0.79}}$
	w/o Shuffled	57.82 ± 0.46	57.83 ± 0.64
Recipe	Node Shuffled	$\textbf{47.87} \pm \textbf{0.83}$	$\textbf{47.66} \pm \textbf{0.88}$
-	Feature Shuffled	55.46 ± 1.09	55.55 ± 1.11
	w/o Shuffled	92.55 ± 0.20	92.54 ± 0.21
ACM	Node Shuffled	90.45 ± 0.49	90.45 ± 0.48
	Feature Shuffled	$\overline{\textbf{89.02}\pm\textbf{1.54}}$	$\overline{\textbf{89.09}\pm\textbf{1.46}}$
	w/o Shuffled	93.88 ± 0.25	94.35 ± 0.23
DBLP	Node Shuffled	93.56 ± 0.32	94.06 ± 0.30
	Feature Shuffled	$\overline{\textbf{93.25}\pm\textbf{0.29}}$	$\overline{\textbf{93.75}\pm\textbf{0.30}}$

We also applied a statistical method, Comparing Degree Distribution (CDD) (Darabi et al., 2023), measuring the structure similarity between real and synthetic HINs. The results indicate that Syn-HING effectively regulates the generation of synthetic HINs. See more details in Appendix A.7.

6 CONCLUSION

We present SynHING, a novel method for generating synthetic HINs, leveraging the real-world HINs as references, systematically generating the major motifs for explanations, and using Intra-/Inter-Cluster Merges to merge multiple groups of base subgraphs to generate synthetic HINs of any specified sizes. SynHING has been validated using four datasets covering three distinct application categories, demonstrating its generality and practicality. we address the scarcity of heterogeneous graph datasets and overcome the need for such datasets in the domain of GNN explanations. To the best of our knowledge, our work introduces the first framework for generating synthetic heterogeneous graphs with ground-truth explanations. Additionally, we design a comprehensive framework for generating diverse synthetic HINs that can be flexibly adjusted and provide a solid foundation for future research on heterogeneous GNN explanations.

540 REFERENCES

553

- Emmanuel Abbe. Community detection and stochastic block models: recent developments. *The Journal of Machine Learning Research*, 18(1):6446–6531, 2017.
- Réka Albert and Albert-László Barabási. Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47, 2002.
- Kevin W. Bowyer, Nitesh V. Chawla, Lawrence O. Hall, and W. Philip Kegelmeyer. SMOTE: synthetic minority over-sampling technique. *CoRR*, abs/1106.1813, 2011. URL http://arxiv.org/abs/1106.1813.
- Ines Chami, Sami Abu-El-Haija, Bryan Perozzi, Christopher Ré, and Kevin Murphy. Machine learning on graphs: A model and comprehensive taxonomy. *Journal of Machine Learning Research*, 23(89):1–64, 2022.
- Yaomin Chang, Chuan Chen, Weibo Hu, Zibin Zheng, Xiaocong Zhou, and Shouzhi Chen. Megnn: Meta-path extracted graph neural network for heterogeneous graph representation learning. *Knowledge-Based Systems*, 235:107611, 2022. ISSN 0950-7051. doi: https://doi.org/10.1016/j.
 knosys.2021.107611. URL https://www.sciencedirect.com/science/article/ pii/S095070512100873X.
- Zhihua Cui, Xianghua Xu, XUE Fei, Xingjuan Cai, Yang Cao, Wensheng Zhang, and Jinjun Chen.
 Personalized recommendation system based on collaborative filtering for iot scenarios. *IEEE Transactions on Services Computing*, 13(4):685–695, 2020.
- Enyan Dai and Suhang Wang. Towards self-explainable graph neural network. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 302–311, 2021.
- Sajad Darabi, Piotr Bigaj, Dawid Majchrowski, Artur Kasymov, Pawel Morkisz, and Alex Fit Florea. A framework for large scale synthetic graph dataset generation, 2023. URL https:
 //arxiv.org/abs/2210.01944.
- Asim Kumar Debnath, Rosa L. Lopez de Compadre, Gargi Debnath, Alan J. Shusterman, and Corwin Hansch. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity. *Journal of medicinal chemistry*, 34 2:786–97, 1991. URL https://api.semanticscholar.org/ CorpusID:19990980.
- Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson.
 Benchmarking graph neural networks. 2020.
- Alvaro Figueira and Bruno Vaz. Survey on synthetic data generation, evaluation methods and gans.
 Mathematics, 10(15), 2022. ISSN 2227-7390. doi: 10.3390/math10152733. URL https: //www.mdpi.com/2227-7390/10/15/2733.
- Maayan Frid-Adar, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit Greenspan. Synthetic data augmentation using GAN for improved liver lesion classification. *CoRR*, abs/1801.02385, 2018. URL http://arxiv.org/abs/1801.02385.
- Xinyu Fu, Jiani Zhang, Ziqiao Meng, and Irwin King. Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding. In *Proceedings of The Web Conference 2020*, pp. 2331–2341, 2020.
- 591 Ming-Yi Hong, Shih-Yen Chang, Hao-Wei Hsu, Yi-Hsiang Huang, Chih-Yu Wang, and Che Lin.
 592 Treexgnn: can gradient-boosted decision trees help boost heterogeneous graph neural networks?
 593 In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023.

594 Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure 595 Leskovec. Strategies for pre-training graph neural networks. In International Conference 596 on Learning Representations, 2020a. URL https://openreview.net/forum?id= 597 HJ1WWJSFDH. 598 Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. Heterogeneous graph transformer. In Proceedings of the web conference 2020, pp. 2704–2710, 2020b. 600 601 Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative ad-602 versarial networks. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition 603 (CVPR), pp. 4396–4405, 2019. doi: 10.1109/CVPR.2019.00453. 604 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. arXiv preprint 605 arXiv:1312.6114, 2013. 606 607 Tong Li, Jiale Deng, Yanyan Shen, Luyu Qiu, Yongxiang Huang, and Caleb Chen Cao. Towards 608 fine-grained explainability for heterogeneous graph neural network, 2023a. 609 Tong Li, Jiale Deng, Yanyan Shen, Luyu Qiu, Huang Yongxiang, and Caleb Chen Cao. Towards 610 fine-grained explainability for heterogeneous graph neural network. Proceedings of the AAAI 611 Conference on Artificial Intelligence, 37(7):8640-8647, Jun. 2023b. doi: 10.1609/aaai.v37i7. 612 26040. URL https://ojs.aaai.org/index.php/AAAI/article/view/26040. 613 614 Yiqiao Li, Jianlong Zhou, Sunny Verma, and Fang Chen. A survey of explainable graph neural 615 networks: Taxonomy and evaluation metrics. arXiv preprint arXiv:2207.12599, 2022. 616 Wanyu Lin, Hao Lan, Hao Wang, and Baochun Li. Orphicx: A causality-inspired latent variable 617 model for interpreting graph neural networks. In Proceedings of the IEEE/CVF Conference on 618 Computer Vision and Pattern Recognition, pp. 13729–13738, 2022. 619 620 Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang 621 Zhang. Parameterized explainer for graph neural network. Advances in neural information pro-622 cessing systems, 33:19620-19631, 2020. 623 Ge Lv, Chen Jason Zhang, and Lei Chen. Hence-x: Toward heterogeneity-agnostic multi-level 624 explainability for deep graph networks. Proc. VLDB Endow., 16(11):2990-3003, July 2023. 625 ISSN 2150-8097. doi: 10.14778/3611479.3611503. URL https://doi.org/10.14778/ 626 3611479.3611503. 627 628 Qingsong Lv, Ming Ding, Qiang Liu, Yuxiang Chen, Wenzheng Feng, Siming He, Chang Zhou, 629 Jianguo Jiang, Yuxiao Dong, and Jie Tang. Are we really making much progress? revisiting, 630 benchmarking and refining heterogeneous graph neural networks. In Proceedings of the 27th 631 ACM SIGKDD conference on knowledge discovery & data mining, pp. 1150–1160, 2021. 632 Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. Generating person-633 alized recipes from historical user preferences. In Kentaro Inui, Jing Jiang, Vincent Ng, and 634 Xiaojun Wan (eds.), Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-635 guage Processing and the 9th International Joint Conference on Natural Language Processing 636 (EMNLP-IJCNLP), pp. 5976–5982, Hong Kong, China, November 2019. Association for Com-637 putational Linguistics. doi: 10.18653/v1/D19-1613. URL https://aclanthology.org/ 638 D19-1613. 639 Ron Milo, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. Net-640 work motifs: simple building blocks of complex networks. Science, 298(5594):824–827, 2002. 641 642 John Palowitch, Anton Tsitsulin, Brandon Mayer, and Bryan Perozzi. Graphworld: Fake graphs 643 bring real insights for gnns. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge 644 Discovery and Data Mining, pp. 3691–3701, 2022. 645 Michael T. Rosenstein, Zvika Marx, Leslie Pack Kaelbling, and Thomas G. Dietterich. To transfer 646 or not to transfer. In Neural Information Processing Systems, 2005. URL https://api. 647 semanticscholar.org/CorpusID:597779.

- 648 Benedek Rozemberczki, Peter Englert, Amol Kapoor, Martin Blais, and Bryan Perozzi. Pathfinder 649 discovery networks for neural message passing. In Proceedings of the Web Conference 2021, pp. 650 2547-2558, 2021. 651 Oleksandr Shchur and Stephan Günnemann. Overlapping community detection with graph neural 652 networks. arXiv preprint arXiv:1909.12201, 2019. 653 654 Tom AB Snijders and Krzysztof Nowicki. Estimation and prediction for stochastic blockmodels for 655 graphs with latent block structure. Journal of classification, 14(1):75–100, 1997. 656 Jonathan M Stokes, Kevin Yang, Kyle Swanson, Wengong Jin, Andres Cubillos-Ruiz, Nina M 657 Donghia, Craig R MacNair, Shawn French, Lindsey A Carfrae, Zohar Bloom-Ackermann, et al. 658 A deep learning approach to antibiotic discovery. Cell, 180(4):688–702, 2020. 659 660 Anton Tsitsulin, John Palowitch, Bryan Perozzi, and Emmanuel Müller. Graph clustering with graph 661 neural networks. arXiv preprint arXiv:2006.16904, 2020. 662 Anton Tsitsulin, Benedek Rozemberczki, John Palowitch, and Bryan Perozzi. Synthetic graph gen-663 eration to benchmark graph learning. arXiv preprint arXiv:2204.01376, 2022. 664 665 Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, and Philip S Yu. Heterogeneous 666 graph attention network. In The World Wide Web Conference, WWW '19, pp. 2022–2032, New 667 York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450366748. doi: 10. 1145/3308558.3313562. URL https://doi.org/10.1145/3308558.3313562. 668 669 Lei Xu, Maria Skoularidou, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Modeling tabular 670 data using conditional GAN. CoRR, abs/1907.00503, 2019. URL http://arxiv.org/abs/ 671 1907.00503. 672 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. Gnnexplainer: 673 Generating explanations for graph neural networks. Advances in neural information processing 674 systems, 32, 2019. 675 676 Hao Yuan, Haiyang Yu, Shurui Gui, and Shuiwang Ji. Explainability in graph neural networks: A 677 taxonomic survey. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45:5782-678 5799, 2020. URL https://api.semanticscholar.org/CorpusID:229923402. 679 Hao Yuan, Haiyang Yu, Jie Wang, Kang Li, and Shuiwang Ji. On explainability of graph neu-680 ral networks via subgraph explorations. In International Conference on Machine Learning, pp. 681 12241-12252. PMLR, 2021. 682 683 Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J. Kim. Graph transformer networks, 2020. 684 685 686 APPENDIX А 687 688 A.1 SYNHING'S COMPLEXITY AND SCALABILITY 689 690 In this section, we theoretically analyze the complexity of the SynHING framework module by 691 module to demonstrate its scalability. Let N represent the motif number, determining the scale of 692 the generated graph. We demonstrate that the total time complexity for SynHING is O(N). For 693 simplicity, we omit node type in the analysis for both Intra-CM and Inter-CM, as nodes merge only 694 with those of the same type, making the complexity linear to the number of types. The generations 695 of the motif (MMG) and the base subgraph (BSG) can be parallelized, and execution time is linear 696 to the number of items. Therefore, the time complexity of MMG and BSG is O(N). 697 The complexity of Intra-CM is analyzed step-by-step as follows: 698
 - (i) Eq.(3), we determine n_{intra} , the number pairs to be sampled.
 - (ii) Eq.(4), we sample n_{intra} nodes from \mathcal{V}_y and \mathcal{V}_i , and pairing them as $\mathfrak{P}_{\text{intra}}$.
 - (iii) Merge process eq.(5).

700

703						U	0	L		
704							Targ	et Node		
705			#Nodes	#Node Types	#Edges	#Edge Types	Туре	#Features	#Classes	
706		IMDB	19,933	4	80,682	6	Movie	3,489	5	
707		Recipe ACM	53,428 10,967	3 4	1,049,048 551,970	4 8	Recipe Paper	1 1,902	5 3	
708		DBLP	27,303	4	296,492	6	Author	340	4	
709										
710	(iv) V	Ve offset the	"incon	ning" suba	anh S. h	v the maxi	mum II)s of C	(graph d	lisioint union)
711	(10) V				apii \mathcal{S}_i U	y the maxin	inum n	J_{S} of C_{y}	(graph u	ilsjonit union).
712	(v) L	Prop the sele	cted no	des in V_i .						
713	(vi) R	leindex the e	edges in	\mathcal{E}_i based of	on the ma	pping dete	rmined	by \mathfrak{P}_{intra}	ı٠	
714 715 716 717 718	The comp complexity iterations, nodes and	lexity (i), (i y of step (vi) making the edges in the	i), and is $O(\mathcal{E})$ total co base s	(v) are $O(\mathcal{E}_i)$. One its omplexity $O(\mathcal{E}_i)$, where $O(\mathcal{E}_i)$ is a set of the transformation of tra	$ \mathcal{V}_i $). The eration contract $\mathcal{O}(N \mathcal{V}_i)$ hich are	e complexity is omplexity is $+ N \mathcal{E}_i)$ o constant w	ity of s s $O(\mathcal{V}_i $ r $O(N)$.r.t. N.	teps (iv) $ + \mathcal{E}_i $), as $ \mathcal{V}_i $	is $O(\mathcal{V})$ There we and $ \mathcal{E}_i $	$ \mathcal{E}_i + \mathcal{E}_i $). The ill be $(N - \mathcal{Y})$ are the number
719	Following	a similar pr	ocess, f	he compley	city of In	ter-CM is a	nalvze	d:		
720	1 0110 11118	a sinnar pi		ne compre						
721	(i) E	eq.(6) and Ed	q.(7), w	e identify a	$\left(\begin{smallmatrix} \mathcal{Y} \\ 2 \end{smallmatrix} \right) \mathbf{c}$	ombination	ns of cl	usters an	d determ	ine the number
722	0	f pairs that 1	need to	be merged	for each	combinatio	on.			
723	(ii) A	After the pai	r numb	er has beer	n determi	ned, we de	erive th	e node n	umber t	hat needs to be
724	n	nerged for ea	ach clus	ter. We ran	domly se	elect nodes	from e	ach clust	er based	on this number
725	W	ithout repla	cement	. (iii) Merg	e proces	s eq.(8).				
726	(iii) V	Ve offset all	the clus	sters C_y . (the sterma is the second state of the second state	he graph	disjoint un	ion in e	eq.(8)).		
727	(iv) D	Prop one of t	the node	es in each p	bair in \mathfrak{P}_{i}	_{nter} in \mathcal{V}_y for	or each	cluster.		
728 729	(v) R	eindex the e	edges in	\mathcal{E}_y for eac	h cluster	based on the	he map	ping dete	ermined	by \mathfrak{P}_{inter} .
730 731 732 733	The comp Since $\sum_{y} N \mathcal{E}_i = 0$	lexity of (iv) $_{\in \mathcal{Y}} \mathcal{V}_y \leq O(N).$), (v), an $N \mathcal{V}_i ,$	nd (vi) are $\sum_{y \in \mathcal{Y}} \mathcal{E}_y $	$O(\sum_{y \in \mathcal{Y}} \le N z)$	$\mathcal{L}_{i}(\mathcal{V}_{y} + E_{i})$. The c	$(E_y)), C$ complex	$\mathcal{D}(\sum_{y\in\mathcal{Y}} \mathcal{Y})$ with of In	$ \mathcal{V}_y)$, an nter-CM	ad $O(\sum_{\mathcal{Y}} \mathcal{E}_y)$. is $O(N \mathcal{V}_i +$
734 735	Overall, S determine	ynHING ca d by motif n	n gener umber	ate large-so N and com	cale HIN plexity c	s in a reason of $O(N)$.	onable	timefran	ne, with	the graph scale
736	A.2 DA	FASETS ANI	HGN	Ns						
738 739 740	To evaluat node class and DBLF	e the SynHI ification da 2. The IM	NG fra tasets: 1 DB dat	mework, w IMDB 1, R aset is a co	e generat Recipe (Mollection	e synthetic Iajumder e of movie d	graphs t al., 20 lata tha	s based o 019), AC t require	n four w CM (Wan s predict	ell-known HIN ng et al., 2019), ing the various
741 742 743 744	genres ass 2021; Hon and user-re	ociated with ng et al., 202 ecipe interactory ore than 203	n each n 23). The ctions. V	novie and f e Recipe da We exclude	ollowing taset is g d recipes	the commutation of the commutati	on setti om Foo r than t	ng as pro d.com ar hree step	evious pa nd includ s, those y We selec	apers (Lv et al., les food recipes with fewer than
745	target nod	e and identi	fy tech	niques used	1 in recir	es as label	ls. The	en. we ch	noose fiv	e techniques to
746	create the	recipe gran	h. On t	the other h	and, AC	M and DB	LP are	citation	networks	s with different
747	goals. AC	M aims to 1	oredict	naper label	s. while	DBLP foci	ises on	predicti	ng autho	r labels. Fig. 7

 Table 4: Statistics of three heterogeneous graph datasets

Table 4 presents the statistics of the four datasets, including the number of node and edge types, the number of nodes and edges, the number of target node features, and the number of classes.

illustrates the graph schema of the four heterogeneous graph datasets.

751

- 752 A.3 EVALUATION METRICS753
- 754 We use Micro-F1 and Macro-F1 as evaluation metrics for node classification and fidelity for ex-755 planation evaluation. Micro-F1 scoring assesses a model's predictions across all samples, with a rendency to emphasize the majority category. In contrast, Macro-F1 scoring equally weights each



Figure 7: Graph schema of the four heterogeneous graph datasets

category, promoting a balanced evaluation of data across different categories. Therefore, we mainly use Macro-F1 as the major evaluation metric (Wang et al., 2019; Lv et al., 2021; Hong et al., 2023).

Fidelity is a metric commonly used to evaluate the performance of the explanation model (Yuan et al., 2021; Li et al., 2022). It measures how closely related the explanations are to the model's predictions. If the critical information is included in the explanation subgraph, the classification model prediction probability should be close to the original prediction, resulting in low fidelity. We use fidelity as the evaluation metric to support that the major motifs can be excellent explanations of ground truths. The following are the details of the fidelity score:

$$Fidelity = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{L} \sum_{l=1}^{L} \left\| f(G_i)_{y_l} - f(\hat{G}_i)_{y_l} \right\|,\tag{9}$$

where $f(G_i)_{y_l}$ and $f(\hat{G}_i)_{y_l}$ denote the prediction probability of y_l of the original graph G_i and major motifs \hat{G}_i (explanation subgraph), respectively. We denote N as the total number of target node samples and L as the number of node labels.

A.4 BENCHMARK HETEROGENEOUS GRAPH NEURAL NETWORKS 781

We used three different concept HGNN models to validate the synthetic graphs. Model parameters 782 follow paper recommendations. The following briefly introduces the models: (1) HGT (Hu et al., 783 2020b) adopts a transformer-based design for handling different node and edge types without man-784 ually defining the meta-path for the HGNN model. (2) SimpleHGN (Lv et al., 2021) introduces the 785 attention mechanism, projects different node-type features to the shared feature space, and then uses 786 GAT as the HGNN backbone. (3) TreeXGNN (Hong et al., 2023) leverages the decision tree-based 787 model XGBoost to enhance the node feature extraction, assisting the HGNN model in getting more 788 prosperous and meaningful information. 789

As the SynHING framework is modulized and can be highly customized according to the characteristics of the reference datasets, we conduct a series of experiments to examine the influence of different tunable parameters on the synthetic HINs. In order to evaluate the performance of SynHING,
we utilize the transductive learning approach for node classification tasks and randomly select 24%
of the target nodes for training, 6% for validation, and 70% for testing (Wang et al., 2019; Lv et al.,
2021; Hong et al., 2023). We repeated all experiments five times and evaluated performance using
average Micro-F1 and Macro-F1 for node prediction and fidelity for interpretation evaluation.

797 798

799

800 801 802

803 804

761

762 763 764

765

766

773

774 775

A.5 SYNTHETIC HINS WITH GROUND-TRUTH EXPLANATIONS



Figure 8: Major Motifs of the Four Heterogeneous Graph Datasets

We evaluate SynHING using three HGNNs on four synthetic HINs (with Syn- in front) based on
their corresponding real-world graphs, shown in Table 5. HGNNs achieve better performance on
Macro-F1 and Micro-F1 scores for learning and inferencing on synthetic graphs compared to real
graphs. These improvements can be attributed to the designated major motifs in synthetic graphs,
shown in Fig. 8, which provide ground-truth explanations for assessing explainability methods and
result in synthetic graphs containing purer information for graph learning. We mimic the graph

15

811



Table 5: Performance comparison of three HGNNs on real and synthetic HINs

We also explored the impact of adjusting the number of major motifs shown in Fig. 10b, which
 directly affects the number of target nodes and the size of the synthetic graph dataset. It is important
 to note that since we kept the hyperparameter settings of the classification model consistent with the



Figure 10: Macro-F1 and Fidelity of Synthetic IMDB in Different SNR and Number of Motifs

original values, rather than fine-tuning them for each synthetic graph dataset, reducing the dataset size to half caused the model to become overfitted, resulting in a decline in performance.

The fidelity results of HGT, Simple-HGN, and TreeXGNN for varying numbers of major motifs are
shown in Fig. 10c. When adjusting the number of motifs, which corresponds to the size of the graph,
the fidelity performance remains stable.

882 883

884

874 875 876

877

878

A.7 APPROXIMATING REFERENCED GRAPH

Users can customize the synthetic graph for various scenarios using the parameters of SynHING, 885 including the number of major motifs N, the number of clusters $|\mathcal{Y}|$, the Intra-CM probabilities p^{ϕ} , the Inter-CM probabilities q^{ϕ} , and the signal-to-noise ratio (SNR) of features α/β . For example, 887 adjusting Intra-CM probability p^{ϕ} and Inter-CM probability q^{ϕ} results in changes in the exclusion of clusters and difficulty of the synthetic graph. However, these parameters can also be directly de-889 termined by the referenced graph G. Although some statistical properties and network schema have 890 been used for generating graphs, it's further demonstrated that the synthetic graph can approximate 891 the referenced graph more closely by adjusting these parameters. The number of major motifs N892 can be set as half of the number of target nodes in \hat{G} , i.e. $N = \frac{1}{2} |\hat{\mathcal{V}}^{\phi_0}|$, since each motif contains 893 exactly 2 target nodes. The number of clusters can be determined by the number of labels $|\hat{\mathcal{Y}}|$ in \hat{G} . 894 The SNR of features α/β can adjust the difficulty of the task on G, or users can determine the means 895 and variances of clusters of features by maximum likelihood estimation. 896

The Intra-/Inter-CM probabilities p^{ϕ} , q^{ϕ} for minor node type $\phi \neq \phi_0$ control the exclusion of clusters, the degree distributions of source nodes, and their counts in the resulting graph \tilde{G} . For instance, in Fig. 11, we observe the node degree distributions for minor node types in both real-world IMDB and SynIMDB, with p = 0.7 and q = 0.3. In contrast, Fig. 12 compares these distributions with SynIMDB using different probabilities: p = 0.9, q = 0.8, and p = 0.2, q = 0.1. As depicted, improper selection of p and q can lead to notable deviations in the degree distribution of minor node types.



912 913

Figure 11: Degree Distributions of Minor Node Types in IMDB and SynIMDB

We further applied a statistical-based method, Comparing Degree Distribution (CDD) (Darabi et al., 2023), to measure the structure similarity between real and synthetic graphs. The CDD value ranges between 0 and 1, with 1 indicating that the distribution of the two structures is exactly the same.
We applied the settings as Fig. 11 and Fig. 12 for structure similarity analysis. Table 6 indicates that the generated SynIMDB can be controlled by the Intra-CM/Inter-CM ratio that influences the



Figure 12: Comparison of Degree Distribution Deviations in SynIMDB with Varying Intra-/Inter-CM Probabilities

similarity with real IMDB. When p=0.7 and q=0.3, Macro-CDD and Micro-CDD are 0.8545 and 0.8279, respectively, which is the most similar to the real IMDB compared to the other two settings. This result highlights the effectiveness of SynHING in regulating the generation of synthetic HINs.

Table 6: Comparing Degree Distribution (CDD) between IMDB and SynIMDB with Varying Intra-/Inter-CM Probabilities

Intra-CM (p), Inter-CM (q)	Macro-CDD	Micro-CDD
p=0.7, q=0.3	0.8545	0.8279
p=0.9, q=0.8	0.7636	0.7612
p=0.2, q=0.1	0.7597	0.7354

A.8 PRETRAINING AND FINETUNING

In this study, we employ synthetic graphs for pretraining. Models are pretrained based on the recommended settings from their respective original papers, with early stopping applied after 30 epochs without validation set improvement. For finetuning, the weights of the HGNN backbone, excluding the adapter layer that maps the heterogeneous features into shared space, are inherited from the pre-trained model. We note that the weights of the backbone and adapter are trained using different learning rates, as the results are sensitive to the learning rate. For instance, while finetuning from pretrained weights, a lower learning rate for the backbone and a higher learning rate for the adapter generally yield better results, whereas a higher learning rate for the backbone and a lower learning rate for the adapter generally leads to better performance when learning from scratch. Consequently, we conduct a grid search for learning rates in both scenarios, as presented in Tables 2 and 3. For the learning rate of the backbone, we try values of $\{10^{-3}, 10^{-4}\}$. For that of the adapter, we try values of $\{1, 5\} \times \{10^{-2}, 10^{-3}, 10^{-4}\}$.

956 957 958

959 960

961

962

963

964

925

926

927 928 929

930

931

932 933 934

935

944 945

946

947

948

949

950

951

952

953

954

955

A.9 IMPLEMENTED ON HGNN EXPLAINER

For our initial testing, we utilized synthetic ACM and DBLP datasets and inputted them into the xPath framework (Li et al., 2023b). The synthetic IMDB dataset we generated is a multiple-choice dataset. Since xPath does not support this, we skipped it for now. We utilized xPath's default parameters, including the HGNN encoder and explainer. We followed the instructions in xPath, which involved two main steps: (1) Training the HGNN and (2) Generating explanations.

We used HGT as our backbone prediction model. During the training stage, it can effectively converge and achieve solid performance (Macro-F1=99.42%, Micro-F1=99.43%) and (Macro-F1=80.06%, Micro-F1=80.81%) on SynACM and SynDBLP, respectively. During the explanation stage, xPath can successfully generate an explanations subgraph with decent accuracy fidelity and probability fidelity, Facc and Fprob (Yuan et al., 2020): Facc=0.15665, Fprob=0.15297 on SynACM and Facc=0.16935, Fprob=0.08701 on SynDBLP, both presenting quite reasonable scores. The above preliminary results show that our generated synthetic datasets can indeed be used to evaluate HGNN explanation algorithms. This warrants a more complete further exploration in future work.

972 A.10 COMPUTING RESOURCES 973

In our experiments, GNN learning utilized an NVIDIA RTX 3060, with fitting a GNN on a heterogeneous information network (HIN) taking under an hour. Graph generation algorithms were executed
on a CPU, with each graph requiring less than an hour to generate.

976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
1000
111111
1000
1000
1001 1002
1001 1002 1003
1001 1002 1003 1004
1001 1002 1003 1004 1005 1006
1001 1002 1003 1004 1005 1006 1007
1001 1002 1003 1004 1005 1006 1007 1008
1001 1002 1003 1004 1005 1006 1007 1008 1009
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022
1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023