

# Heterogeneous Graph Contrastive Multi-view Learning

Zehong Wang<sup>\*</sup> Qi Li<sup>†</sup> Donghua Yu<sup>†</sup> Xiaolong Han<sup>‡</sup> Xiao-Zhi Gao<sup>§</sup>  
Shigen Shen<sup>¶¶</sup>

## Abstract

Inspired by the success of Contrastive Learning (CL) in computer vision and natural language processing, Graph Contrastive Learning (GCL) has been developed to learn discriminative node representations on graph datasets. However, the development of GCL on Heterogeneous Information Networks (HINs) is still in the infant stage. For example, it is unclear how to augment the HINs without substantially altering the underlying semantics, and how to design the contrastive objective to fully capture the rich semantics. Moreover, early investigations demonstrate that CL suffers from sampling bias, whereas conventional debiasing techniques are empirically shown to be inadequate for GCL. How to mitigate the sampling bias for heterogeneous GCL is another important problem. To address the aforementioned challenges, we propose a novel Heterogeneous Graph Contrastive Multi-view Learning (HGCML) model. In particular, we use metapaths as the augmentation to generate multiple subgraphs as multi-views, and propose a contrastive objective to maximize the mutual information between any pairs of metapath-induced views. To alleviate the sampling bias, we further propose a positive sampling strategy to explicitly select positives for each node via jointly considering semantic and structural information preserved on each metapath view. Extensive experiments demonstrate HGCML consistently outperforms state-of-the-art baselines on five real-world benchmark datasets. To enhance the reproducibility of our work, we make all the code publicly available at <https://github.com/Zehong-Wang/HGCML>.

**Keywords:** Graph contrastive learning, heterogeneous information network, self-supervised learning, graph neural network, multi-view learning.

## 1 Introduction

Considering the capacity for modeling complex systems, Heterogeneous Information Networks (HINs) that preserves

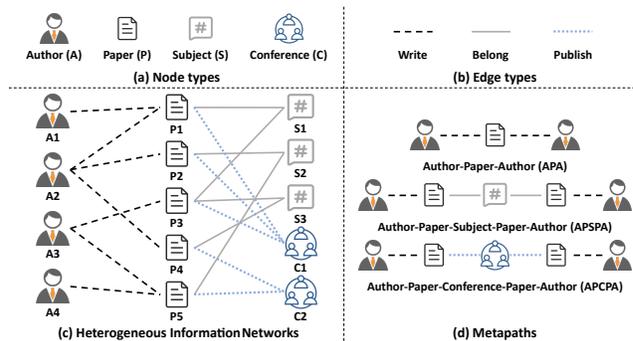


Figure 1: An example of heterogeneous information networks.

rich semantics have become a powerful tool for analyzing real-world graphs. As illustrated in Figure 1, we present a concise example of a heterogeneous bibliography network with four types of nodes and three types of relations. Recently, Graph Neural Networks (GNNs) [1] have emerged as a dominant technique in mining graph structure datasets, and its variant, Heterogeneous Graph Neural Networks (HGNNs) [2, 3, 4, 5, 6, 7], has occupied the mainstream of HIN analysis. In general, HGNNs are trained in an end-to-end manner, which requires abundant, various, and dedicated-designed labels for different downstream tasks. However, in the majority of real-world scenarios, it is highly expensive and/or difficult to collect labels.

Contrastive Learning (CL) [8, 9, 10] that automatically generates supervise signals from data itself is a promising solution for learning representations in a self-supervised manner. By maximizing the confidence (i.e., mutual information) [11] between positive pairs and minimizing the confidence between negative pairs, CL is capable to learn discriminative representations without explicit labels. Inspired by the success of CL in computer vision [9, 10], a wide range of Graph Contrastive Learning (GCL) methods have been proposed. For example, DGI [8] exploits a contrast between graph patches (i.e., nodes) and graph summaries, and GRACE [12] maximizes the mutual information between the same node in two augmented views. Despite some works generalizing the key idea of CL to homogeneous graphs, there are still three fundamental challenges that need to be addressed in explor-

<sup>\*</sup>School of Mathematics, University of Leeds, Leeds, United Kingdom.

<sup>†</sup>Department of Computer Science and Engineering, Shaoxing University, Shaoxing, China

<sup>‡</sup>School of Computer and Software, Nanjing University of Information Science and Technology, Nanjing, China

<sup>§</sup>School of Computing, University of Eastern Finland, Kuopio, Finland.

<sup>¶</sup>School of Information Engineering, Huzhou University, Huzhou, China

<sup>¶¶</sup>To whom correspondence should be addressed(shigens@zjhu.edu.cn)

ing the great potential of CL in heterogeneous graphs:

**(1) How to design distinct views?** Data augmentation that creates corrupted views is shown to be an essential technique to improve the quality of representations [13]. In GCL, prevalent augmentation methods include edge dropping/adding, node dropping/adding, feature shuffling, and so forth. Although these methods excel in homogeneous graphs [12, 14], we believe that they significantly change the latent semantics of HINs. Take a bibliographic network as an example (Figure 1); if the link between Author 3 (A3) and Paper 5 (P5) is dropped, the closest path between Author 3 (A3) and Author 4 (A4) will be changed from 2-hop (A3-P5-A4) to 6-hop (A3-P3-S3-P4-C2-P5-A4). To prevent the knowledge perturbation caused by simple augmentation techniques, we propose to leverage metapaths, the composition of semantic relations, to augment datasets. By applying metapaths, we create multiple different yet complementary subgraphs, referred to as metapath views, without altering the underlying semantics while also capturing the high-order relationships on HINs.

**(2) How to set proper contrastive objectives?** The choice of contrastive objectives (i.e., pretext tasks) determines the discriminativeness of representations in downstream tasks. For HIN, the standard choice of pretext tasks is still unclear. Different works present their own solutions. For example, DMGI [15] proposes to use metapaths to learn a shared consensus vector as node representation, HeCo [16] performs contrast between the aggregation of metapaths (view 1) and network schema (view 2), and HDMI [17] and STENCIL [18] iteratively maximize the mutual information between a single metapath and the aggregation of them. Despite these approaches attempting to incorporate the universal knowledge across all metapaths, we think that they actually assume metapaths are independent, which is different from the complementary nature, failing to capture the consistency between metapaths and thus leading to sub-optimality. To directly model the correlation between metapaths, we propose an intuitive yet unexplored contrastive objective that performs contrast between each pair of metapaths. To be specific, the contrast between two augmented views of a metapath (intra-metapath) aims to learn augmentation-invariant representations, and the contrast between two views generated from two sources (inter-metapath) ensures the alignment across metapaths.

**(3) How to mitigate the sampling bias?** Sampling bias indicates that the negative samples, which are randomly selected from the original datasets, are potential to share the same class with the anchor node (i.e., act as false negatives). Empirical, the sampling bias will lead to a significant performance drop. To prevent the issue, existing works [19, 20] aim to select or synthesize hard negatives to mitigate the impact of false negatives. However, these methods are demonstrated to bring limited benefits or even impose adverse im-

pacts on GCL [21, 22]. To alleviate the issue of false negatives, we propose a positive sampling strategy that collaboratively considers topological and semantic information across metapaths to explicitly decide the positive counterparts for each anchor.

To summarize, we propose a Heterogeneous Graph Contrastive Multi-view Learning (HGCML) model to learn informative node representations on HINs. In particular, we apply metapaths to create multiple views and leverage a GNN model to encode node representations. Then, we employ a novel contrastive objective that aims to maximize the mutual information between any pairs of metapath views (for both intra-metapath and inter-metapath) to explicitly model the complementarity among metapaths, which is neglected in other works. Specifically, we maximize the confidence between two metapaths at node and graph levels to acquire local and global knowledge. To further enhance the expressiveness, we propose a positive sampling strategy that directly picks hard positives for each node based on graph-specific topology and semantics to mitigate the sampling bias inherent in CL. We highlight the contributions as follows:

- We propose a heterogeneous graph contrastive multi-view learning framework, named HGCML, to learn discriminative node representations. The model leverages metapaths in HINs to generate multiple views and employs a novel contrastive objective to model the consistency between any pairs of metapath views at node and graph levels.
- We propose a positive sampling strategy, which selects the most similar nodes as positive counterparts for each anchor by considering semantics and topology across metapath views, to remedy the sampling bias.
- We conduct extensive experiments on five real-world datasets to evaluate the superiority of our model. Experimental results show HGCML outperforms state-of-the-art (SOTA) self-supervised and even supervised baselines.

## 2 Related Work

Following the message passing paradigm, GNNs [23, 24, 25, 26, 4, 6, 27] have received great attention in recent years for learning representations of nodes in HINs. For example, HAN [4] uses attention mechanism to model the correlation between nodes at both metapath-level and semantic-level and MAGNN [6] applies metapath encoders to gain fine-grained knowledge preserved in metapaths. To get rid of the impact of metapaths, RGCN [3] and its variants [5, 7] directly utilize type-specific matrices to model the relationships between different types of nodes in HINs. Despite these models achieving remarkable performance in mining heterogeneous graph datasets, they fail to be performed in a self-supervised

manner.

In another line, GCL that marries the power of GNN and CL has emerged as an important paradigm to learn representations on graphs without annotations. As a pioneering work, DGI [8] treats node embedding and graph summaries as positive pairs and utilizes InfoMAX [11] to optimize the objective. Following this line, MVGRL [28] proposes to use graph diffusion as an augmentation method to generate multiple views and GraphCL [14] further analyzes the role of augmentations in introducing prior knowledge. Inspired by instance discrimination [29], GRACE [12] and GCA [30] propose to leverage the node-level objective in contrasting to preserve node-level discrimination. In addition, BGRL [31] adopts the key idea of BYOL [32] to perform contrast without negative samples via bootstrapping to save memory consumption.

Meanwhile, some studies have generalized the key idea of GCL on HINs. For instance, HDGI [33] extends DGI to heterogeneous graphs and DMGI [15] utilizes a metapath encoder to train consensus vectors as node representations. CKD [34] models the regional and global knowledge between each pair of metapaths, failing to capture node-level properties. CPT-HG [35] applies relation- and subgraph-level pretext tasks to pre-train HGNN on large-scale HINs, and HDMI [17] introduces a triplet loss to further enhance generalization. However, these methods still do not consider the sampling bias inherent in GCL, inevitably leading to sub-optimality. To mitigate the sampling bias, STENCIL [18] and HeCo [16] propose to apply metapath similarity to measure the hardness between nodes to synthesize hard negatives or select semantic positives. However, these models assume metapaths are independent, and treat the aggregation of metapath-induced subgraphs as a single contrastive view, thus failing to model the consistency and complementarity between metapath views.

Different from the aforementioned methods, our model keeps three distinct advancements: (i) applying metapaths as an augmentation approach to generate multi-views for HINs, instead of treating the aggregation of all metapaths as a single view, which ensures keeping fine-grained and complementary properties for each metapath; (ii) performing contrast between any pairs of metapath-induced subgraphs to learn augmentation-invariant representations for a single metapath and to align the consistency between different metapaths; and (iii) explicitly selecting positive samples for each node via considering topology and semantics preserved on metapath views to mitigate sampling bias.

### 3 Preliminary

**DEFINITION 1. Heterogeneous Information Network (HIN)** refers as to a graph consisting of various types of nodes and edges, represented as  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{T}, \mathcal{R}\}$ , where  $\mathcal{V}$  and  $\mathcal{E}$  are the node set and edge set, respectively, and  $\mathcal{T}$

and  $\mathcal{R}$  denote node types and relation types, associated with a node mapping function  $\psi : \mathcal{V} \rightarrow \mathcal{T}$  and an edge mapping function  $\phi : \mathcal{E} \rightarrow \mathcal{R}$ . Note that  $|\mathcal{T}| + |\mathcal{R}| > 2$ .

**DEFINITION 2. Metapath.** Metapath  $\mathcal{P}_m, m \in \mathcal{M}$ , is the composition of relations in HINs, defined as  $\mathcal{P}_m := \mathcal{T}_0 \xrightarrow{\mathcal{R}_0} \mathcal{T}_1 \xrightarrow{\mathcal{R}_1} \dots \xrightarrow{\mathcal{R}_n} \mathcal{T}_{n+1}$ , where  $\mathcal{M}$  is the set of metapaths. For example, we illustrate three metapaths extracted from DBLP in Figure 1 (d), which describe co-author (APA), co-subject (APSPA), and co-conference (APCPA) relationships.

**DEFINITION 3. Metapath-based Neighbors.** Given a metapath  $\mathcal{P}_m$ , metapath-based neighbors  $\mathcal{N}_v^{\mathcal{P}_m}$  is defined as a set of nodes connected to the target node through metapath  $\mathcal{P}_m$ . For example, in Figure 1 (c), the metapath-based neighbors of Author 1 via metapath APA is Author 2.

## 4 Heterogeneous Graph Contrastive Multi-view Learning

In this section, we present a heterogeneous graph contrastive multi-view learning framework to learn representations of nodes in HINs. The overview architecture is illustrated in Figure 2.

**4.1 Data Augmentation** For HINs, collectively applying metapaths to construct multi-views is a natural way to supplement the dataset in opposition to simple augmentation techniques. The created multi-views are actually complementary with each other because metapaths depict various facets of the same HIN. Given a set of metapaths  $\{\mathcal{P}_0, \mathcal{P}_1, \dots, \mathcal{P}_{|\mathcal{M}|}\}$  where  $|\mathcal{M}|$  is the number of metapaths, we extract multiple subgraphs (i.e., metapath views)  $\{\mathcal{G}_0, \mathcal{G}_1, \dots, \mathcal{G}_{|\mathcal{M}|}\}$  from the original graph to sustain the rich semantics preserved in HINs. For the subgraph  $\mathcal{G}_m$  generated through metapath  $\mathcal{P}_m$ , we construct the direct neighborhoods for each node  $v$  as its metapath-based neighbors  $\mathcal{N}_v^{\mathcal{P}_m}$ . Each metapath view is associated with a node feature matrix  $\mathbf{X}_m$  and an adjacent matrix  $\mathbf{A}_m$ . We leverage a GNN encoder  $f(\cdot)$  to learn node representation  $\{\mathbf{H}_0, \dots, \mathbf{H}_m, \dots, \mathbf{H}_{|\mathcal{M}|}\}$  from each metapath-induced view, where  $\mathbf{H}_m = f(\mathbf{X}_m, \mathbf{A}_m)$ . In practice, we leverage additional data augmentations (i.e., feature masking and edge dropping) with specific probabilities  $p_f$  and  $p_e$  to further corrupt metapath views to make the task to be more difficult, which ensures the learned representations to be more discriminative.

**4.2 Contrastive Objectives** To distill rich semantics in HINs, we propose a novel contrastive objective to maximize the correlation between any pair of metapath views. In particular, the contrastive objective is collaboratively performed in intra-metapath (i.e., contrast between two corrupted ver-

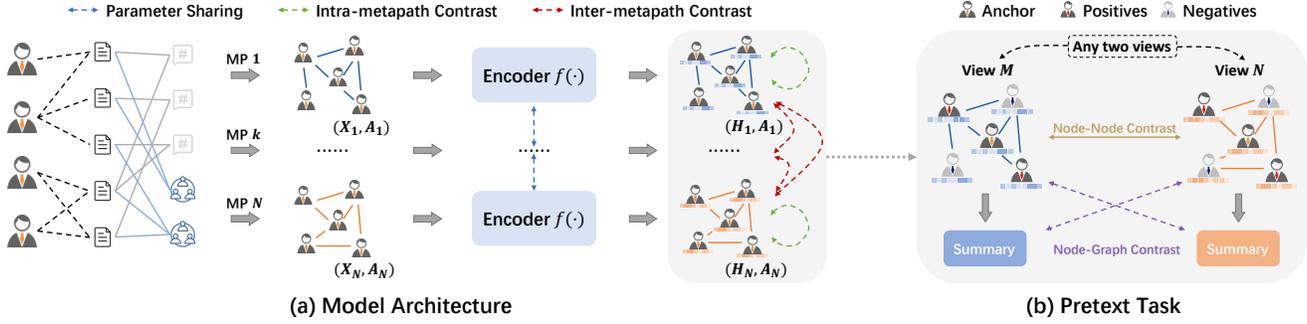


Figure 2: (a) The overview framework of HGCML. Firstly, we generate multi-views via the guide of metapaths (MPs), then leverage a graph neural network (GNN)  $f(\cdot)$  to encode the representations. Following, we employ a novel contrastive objective to capture the consistency between each pair of metapath views. The contrast performed between two corrupted versions of a single view is called intra-metapath contrast and the objective applied between two distinct metapath-induced views is referred to as inter-metapath contrast. (b) The proposed pretext task simultaneously learns from graph patches and graph summaries to acquire local and global knowledge. Note that the task is performed between any two metapath-induced views, where  $M = N$  denotes intra-metapath contrast and  $M \neq N$  indicates inter-metapath contrast. In addition, we perform positive sampling to enhance the expressiveness of node-node contrast.

sions of a metapath view) and inter-metapath (i.e., contrast between two views from different metapaths), demonstrated in Figure 2(a) with green and red dot lines. We argue that the intra-metapath contrast independently learns the augmentation-invariant latent for each metapath view and the inter-metapath contrast is to align the representations gained from various sources to acquire the complementarity inherent in metapaths. Thus, we thoroughly gain the underlying knowledge maintained in individual metapath views and explicitly model the dependencies between pairs of different metapath views. In addition, the pretext task between two views jointly learns from node- and graph-level knowledge to enhance representativeness, as shown in Figure 2(b). Note that in the node-level contrasting, we select hard positives via the proposed sampling strategy to mitigate the sampling bias.

**4.2.1 Node-Node Contrast** Node-node contrast aims to learn discriminative node representations to boost node-level downstream tasks. Specifically, we perform contrast between the anchor and its positive counterparts in two views to maximize (resp. minimize) the confidence between similar (resp. unassociated) nodes:

$$(4.1)$$

$$\mathcal{L}_{local}^{(m,n)}(u, \mathbb{P}_u) = -\log \frac{\sum_{v \in \mathbb{P}_u} \theta(h_u^m, h_v^n)}{\sum_{v \in \mathbb{P}_u} \theta(h_u^m, h_v^n) + \sum_{v \in (\mathcal{V} \setminus \mathbb{P}_u)} \theta(h_u^m, h_v^m) + \sum_{v \in (\mathcal{V} \setminus \mathbb{P}_u)} \theta(h_u^m, h_v^n)},$$

where the values of  $m$  and  $n$  can be the same,  $h_u^m$  is the representation for node  $u$  in view  $m$ ,  $\mathbb{P}_u$  denotes the selected positive samples for  $u$ . We use similarity function  $\theta(h_u^m, h_v^n) = e^{\varphi(\rho(h_u^m), \rho(h_v^n)) / \tau}$  to compute the distance be-

tween node representations where  $\varphi(\cdot, \cdot)$  measures the cosine distance between two vectors,  $\rho(\cdot)$  denotes a non-linear projector head that increases the expressiveness, and  $\tau$  controls the data distribution. This objective function that pulls semantic similar nodes close and pushes dissimilar nodes away contributes to the discrimination of node representations.

**4.2.2 Node-Graph Contrast** Different from node-node contrast that learns local semantics across multi-views, we also perform node-graph contrast as an auxiliary task to facilitate the representation learning by injecting metapath-specific knowledge. We define the node-graph contrast objective as follows:

$$(4.2)$$

$$\mathcal{L}_{global}^{(m,n)}(u) = -\log(\mathcal{D}(h_u^m, s_m)) - \log(1 - \mathcal{D}(h_u^n, s_m)),$$

where the value of  $m$  and  $n$  can be the same, and  $s_m$  is the graph summary of metapath view  $\mathcal{G}_m$  calculated via a  $READOUT(\cdot)$  function (mean pooling in this paper), and  $\mathcal{D}(h, s) = \omega(\rho(h), \rho(s))$  where  $\omega(\cdot, \cdot)$  is a discriminator that consists of a bilinear layer  $BiLinear(\cdot)$  and a sigmoid function  $\sigma(\cdot)$ . By imparting global knowledge brought by metapaths, we ensure the representations of nodes are more informative.

**4.2.3 Overall Objective** The overall objective  $\mathcal{J}$  to be maximized is defined as the aggregation of all pairs of metapaths, formally given by

$$(4.3) \quad \mathcal{J} = \sum_{m \in \mathcal{M}} \sum_{n \in \mathcal{M}} \sum_{u \in \mathcal{V}} \mathcal{L}_{local}^{(m,n)}(u, \mathbb{P}_u) + \mathcal{L}_{global}^{(m,n)}(u),$$

where  $\mathcal{M}$  is the set of metapaths. After optimizing the contrastive objective, we perform late fusion function  $\eta(\cdot)$  (sum

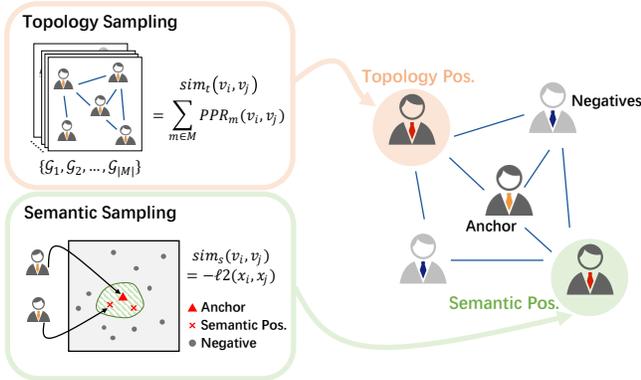


Figure 3: The positive sampling strategy where Personalized PageRank (PPR) is used to measure the topological similarity between nodes, and L2 distance is leveraged to compute the distance between nodes in semantic space to discover semantic associations.

or concatenation) on node representations learned from multiple metapath views to obtain the unified node representations  $h_u$  for downstream tasks as

$$(4.4) \quad h_u = \eta(\{h_u^m, m \in \mathcal{M}\}),$$

where  $h_u^m$  denotes the learned representations for node  $u$  in metapath-induced view  $\mathcal{G}_m$ ,  $\mathcal{M}$  is the metapath set.

**4.3 Positive Sampling Strategy** Sampling bias is an important problem in CL since false negatives will generate adverse signals. However, existing debiasing techniques [20] are theoretically and empirically verified to lead to severer sampling bias for GCL [22], because the message passing mechanism smooths the node representations. To overcome the deficiency, we propose to leverage two different yet reciprocal similarity measurements (i.e., topology and semantics) to define the distance between nodes, as shown in Figure 3, and explicitly select the most similar nodes as positive samples.

**4.3.1 Topology Positive Sampling** To analyze the similarity between nodes based on topological structure, we propose to use the graph diffusion kernel [36] that assesses the global node importance to compute the distance between two arbitrary nodes. In practice, we apply Personalized PageRank (PPR) score  $\mathbf{S}_m$  to measure the node-level relationship for each metapath view  $\mathcal{G}_m$ , which is defined as

$$(4.5) \quad \mathbf{S}_m = \sum_{k=0}^{\infty} \alpha(1-\alpha)^k (\mathbf{A}_m \mathbf{D}_m^{-1})^k,$$

where  $\mathbf{S}_m \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ ,  $\mathbf{A}_m$ , and  $\mathbf{D}_m$  are the diffusion matrix, adjacent matrix and diagonal degree matrix for meta-

path view  $\mathcal{G}_m$ , respectively, and  $\alpha$  denotes teleport probability, whose default is 0.85. Formally, we define the PPR similarity between two nodes  $v_i$  and  $v_j$  under metapath view  $\mathcal{G}_m$  as the  $i$ -th row and  $j$ -th column in the diffusion matrix  $PPR_m(v_i, v_j) = \mathbf{S}_m[i, j]$ . The value in fact describes the stationary probability of starting from  $v_i$  to reach  $v_j$  via an infinite random walk in the metapath view  $m$ . Then, we aggregate the PPR scores computed on all metapath-induced views to determine the topological similarity  $sim_t(v_i, v_j)$  for each node pair as

$$(4.6) \quad sim_t(v_i, v_j) = \sum_{m \in \mathcal{M}} PPR_m(v_i, v_j),$$

and select the top- $k$  similar nodes for each anchor as the topology positives  $\mathbb{P}^t$ .

**4.3.2 Semantic Positive Sampling** Apart from structural information, graph datasets also preserve rich semantics on the node itself. To measure the semantical similarity between nodes, we propose to utilize a simple metric  $sim_s(v_i, v_j)$  to compute the distance between attributes of nodes, which is defined as

$$(4.7) \quad sim_s(v_i, v_j) = -l2(x_i, x_j),$$

where  $x_i$  and  $x_j$  denote attributes on nodes  $v_i$  and  $v_j$ , respectively, and  $l2(\cdot, \cdot)$  measures the L2-distance between two data points. The attributes for each node will not change across metapath views, thus we only need to process once to calculate the distance between pairs of nodes. Finally, we also select the top- $k$  similar nodes for each anchor as the semantic positives  $\mathbb{P}^s$ . At the time, we define the positive samples  $\mathbb{P}_u$  for node  $u$  across metapath views as

$$(4.8) \quad \mathbb{P}_u = \mathbb{P}_u^t \cup \mathbb{P}_u^s.$$

Note that the positive sampling phase is performed in pre-processing, so the module will not significantly increase the computational complexity.

## 5 Experiments

### 5.1 Experimental Setup

**5.1.1 Datasets and Baselines** To demonstrate the superiority of HGCML over SOTA, we conduct extensive experiments on five public benchmark datasets, including ACM, DBLP, IMDB, Aminer, and FreeBase. We evaluate the performance of our model against various baselines from shallow graph representation learning algorithms, including DeepWalk [37], Metapath2vec(MP2vec) [38], HIN2vec [39], HERec [40], to GCL methods (e.g., DGI [8], GRACE [12], DMGI [15], STENCIL [18], HeCo [16]) to supervised GNNs, like GCN [24], GAT [26], HAN [4]. Note that DMGI, STENCIL, and HeCo are dedicated for heterogeneous graphs.

Methods	Data	Node Classification					Node Clustering				
		ACM	DBLP	IMDB	Aminer	FreeBase	ACM	DBLP	IMDB	Aminer	FreeBase
DeepWalk	A	81.78±0.04	88.09±0.07	56.36±0.33	84.93±0.09	69.63±0.05	41.15±0.49	20.13±2.57	5.97±0.23	30.17±2.86	14.56±0.08
MP2vec	A	79.82±0.23	87.67±0.12	50.78±0.18	84.14±0.06	69.66±0.11	37.74±0.09	73.77±0.18	2.71±0.34	26.52±0.36	14.93±1.05
HIN2vec	A	85.23±0.09	91.40±0.08	50.73±0.23	80.77±0.06	67.42±0.14	40.79±0.49	68.83±1.42	3.88±0.29	23.76±0.64	14.26±2.03
HERec	A	67.15±0.85	90.75±0.39	49.12±0.22	80.63±0.10	68.04±0.19	45.39±2.11	70.38±3.29	4.39±1.01	31.05±0.69	15.32±1.05
DGI	X, A	88.44±0.30	90.16±0.60	52.00±0.94	83.24±0.19	68.42±0.31	43.47±2.25	54.44±2.07	4.09±1.88	29.80±1.86	15.16±1.13
GRACE	X, A	87.64±0.31	91.28±0.07	54.80±0.82	83.43±0.22	69.25±0.14	46.50±4.58	67.98±1.32	1.58±1.12	24.12±5.24	16.23±3.37
DMGI	X, A	76.76±1.23	91.60±0.66	51.16±0.63	79.19±0.32	67.69±0.21	52.53±1.73	67.41±0.07	5.45±0.12	28.32±0.44	12.35±0.35
STENCIL	X, A	88.23±0.91	92.56±0.22	57.83±0.62	84.61±0.53	68.26±0.21	56.67±2.51	71.40±1.93	8.25±1.09	29.99±2.69	13.19±1.10
HeCo	X, A	88.97±1.12	92.24±0.48	52.12±0.72	85.22±0.10	69.02±0.07	56.93±1.59	70.03±1.25	7.41±1.26	30.61±3.81	12.07±1.47
HGCMC	X, A	91.02±0.13	93.29±0.12	60.75±0.71	86.63±0.11	71.41±0.04	65.13±1.33	73.28±0.76	<b>9.34±0.86</b>	<b>36.10±2.44</b>	<b>15.46±1.65</b>
HGCMC-P	X, A	<b>91.34±0.17</b>	<b>93.44±0.08</b>	<b>61.02±0.49</b>	<b>87.03±0.06</b>	<b>71.53±0.14</b>	<b>65.75±1.62</b>	<b>74.53±0.48</b>	8.95±1.06	35.62±1.74	<b>16.26±2.56</b>
GCN	X, A, Y	89.87±0.79	92.04±1.03	58.42±1.42	85.42±0.48	69.13±2.51	58.14±0.90	77.71±1.35	8.59±0.84	37.80±1.69	15.77±2.97
GAT	X, A, Y	88.84±0.61	92.51±1.28	57.97±1.64	84.37±0.42	70.42±0.55	62.22±3.67	72.06±1.61	8.04±1.76	36.81±0.66	15.44±1.32
HAN	X, A, Y	89.50±1.21	<u>93.27±0.58</u>	54.78±1.01	85.90±0.43	<u>70.98±1.07</u>	60.98±2.38	<u>78.20±0.83</u>	6.80±2.32	35.37±0.48	<u>16.38±1.73</u>

Table 1: Performance of node classification and clustering on five benchmark datasets in terms of micro-F1 and normalized mutual information (NMI). Boldfaces and underlines denote the best performance among self-supervised and supervised methods, respectively. For our model, we use the suffix -P to indicate the positive sampling version.

Intra-	Inter-	Local	Global	ACM	DBLP	IMDB	Aminer	FreeBase
✓	-	✓	-	89.20	91.94	58.71	84.95	69.61
✓	-	-	✓	83.20	90.95	48.56	83.42	69.38
✓	-	✓	✓	90.52	92.52	60.12	86.17	71.16
-	✓	✓	✓	88.32	92.44	59.80	85.92	70.65
✓	✓	✓	✓	<b>91.02</b>	<b>93.29</b>	<b>60.75</b>	<b>86.63</b>	<b>71.41</b>

Table 2: Ablation study of the proposed HGCMC for pretext tasks on node classification, where the intra- and inter- are abbreviations of intra-metapath and inter-metapath contrasts, and local and global indicate node-node and node-graph contrasts, respectively.

**5.1.2 Evaluation Protocol** We evaluate HGCMC on node classification and node clustering. For node classification, we use Micro-F1 as the metric and follow the linear protocol that utilizes the learned graph encoder as a feature extractor to train a simple linear classifier with 20% random samples as the training set. For node clustering, we apply  $K$ -means to generate clusters and utilize Normalized Mutual Information (NMI) as the metric. To mitigate the impact of initialized centroids, we perform 10 times clustering and report the average results. For all baselines, we run 10 times and present the average scores with standard deviations. For DGI, GRACE, GCN, and GAT, we create homogeneous graphs based on metapaths and report the best results.

**5.1.3 Implementation Details** We leverage a 1-layer GCN as the encoder for each metapath-induced view. The parameters are initialized via Xavier initialization and we apply Adam as the optimizer. We perform grid search to tune the learning rate from  $5e-4$  to  $5e-3$ , the value of temperature from 0.2 to 0.8, the corrupt rate from 0.1 to 0.7, and the number of positives from 0 to 128. Moreover, we set early stop to 20 epochs, node dimension to 64, activation function

Topology Pos.	Semantic Pos.	ACM	DBLP	IMDB	Aminer	FreeBase
-	-	91.02	93.29	60.75	86.63	71.41
✓	-	91.17	93.34	59.92	86.88	71.43
-	✓	90.88	93.35	58.28	86.91	71.45
✓	✓	<b>91.34</b>	<b>93.44</b>	<b>61.02</b>	<b>87.03</b>	<b>71.53</b>

Table 3: Ablation study of the proposed HGCMC for positive sampling strategies on node classification with four variants; ✓ denotes the specific type of positives are selected.

to  $ReLU(\cdot) = \max(\cdot, 0)$ , and use concatenation as the fusion function in ACM and DBLP, and summation in other datasets.

**5.2 Quantitative Results** We report the quantitative result of node classification and node clustering with standard deviations in Table 1. From the table, we observe that GCL methods generally perform better than shallow unsupervised baselines, since the instance discrimination applied on CL captures underlying semantics preserved in HINs but the graph reconstruction adopted in classical methods only considers the topological structure. Our models (HGCMC and HGCMC-P) consistently outperform SOTA self-supervised graph learning methods across all datasets by a large margin on supervised classification and unsupervised clustering tasks, and even achieve competitive results compared to supervised baselines. Beyond that, the performance of HGCMC-P (positive sampling version) is commonly better than its vanilla version that performs contrast between the same node in different views, demonstrating the necessity of introducing correlated nodes as positives to mitigate the sampling bias. Compared with heterogeneous graph contrastive learning methods (i.e., DMGI, STENCIL, and HeCo), our model always acquires higher scores in both classification and clustering. We assume that

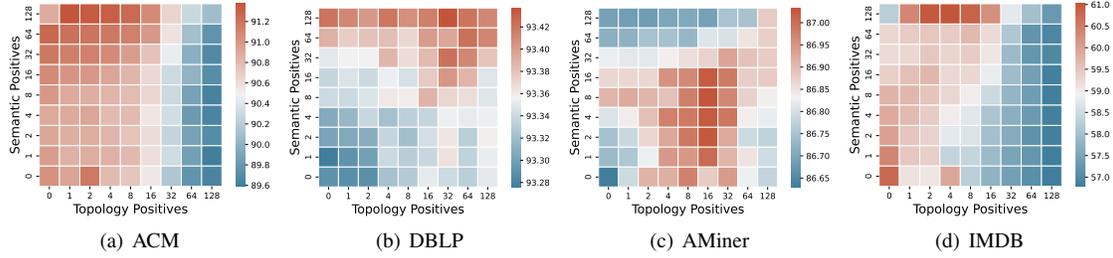


Figure 4: Hyperparameter sensitivity of positive sampling thresholds on node classification. Note that when no extra positives are selected (i.e., the left bottom corner), the model picks the anchor node itself as the positive sample.

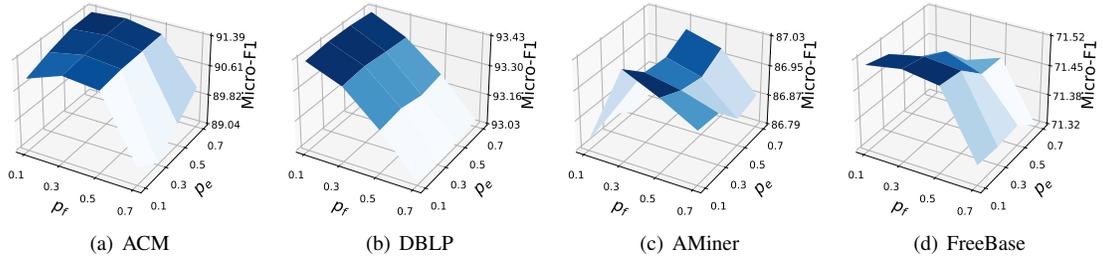


Figure 5: Hyperparameter sensitivity of augmentations (edge dropping  $p_e$  and feature masking  $p_f$ ) on node classification.

it is because (1) the intra-metapath contrast is performed between nodes on two corrupted views induced from the same metapath to learn the discriminative representations and (2) the inter-metapath contrast captures the complementarity between metapaths instead of treating the aggregation of them as a single view under the independent assumption applied in mentioned baselines.

### 5.3 Ablation Study

**5.3.1 Pretext Task** To verify the role of each component in the contrastive objective, we perform ablation studies, as shown in Table 2, to compare the performance of multiple variants on node classification. From the table, we observe that (1) the node-node contrast provides better discrimination ability compared with the node-graph contrast since the fine-grained information (patches) is leveraged in learning representations. When the node- and graph-level objectives are jointly optimized, the performance is significantly improved, showing the necessity of simultaneously modeling local- and global-level dependencies. (2) The intra-metapath contrast is essential in promoting the learning procedure, reflected in the competitive performance obtained in the initialized variant (Intra- & Node) against SOTA self-supervised baselines. (3) The variant with full components persistently achieves the best performance since the intra-metapath contrast captures the latent semantics of each metapath-induced view and the inter-metapath contrast aligns the consistency

between metapaths. If one of them is removed, we cannot thoroughly model the relationship between metapaths, thus encountering model degradation.

**5.3.2 Positive sampling strategy** We also conduct experiments to evaluate the impact of positive sampling, as presented in Table 3. We can find that the selected positives indeed improve the performance by implicitly defining hard negatives. In addition, the significance of these two positive sampling strategies depends on the choice of datasets, i.e., there is no obvious superiority between topology positives and semantic positives. However, when they are jointly leveraged, our model achieves the best scores. The phenomenon demonstrates that the selected positives based on different strategies are distinct yet complementary.

### 5.4 Hyperparameter Analysis

**5.4.1 Positive Sampling Thresholds** In the above section, we analyze the impact of positive sampling strategies in enhancing the quality of representations, here we delve into the positive sampling thresholds to provide a further examination, illustrated in Figure 4. As we can see, the best performance is achieved with a large number of semantic positives and a small number of topology positives (ACM, DBLP, IMDB). When the number of topology positives is too large, the performance generally encounters a drop. We assume the phenomena derive from the inherent property of defined

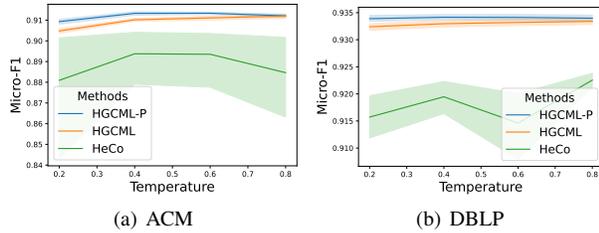


Figure 6: Hyperparameter sensitivity of temperature  $\tau$  in terms of Micro-F1.

similarity functions. To be specific, the semantical similarity is independently measured on attributes of nodes in the representation space, whereas the topological similarity is calculated based on the adjacent matrix, which makes the function naturally biased to nodes with dense connections. Thus, when the number of topology positives is too large, there will contain too many noisy nodes. The Aminer does not follow the observation on the other datasets, whose best performance is achieved when the number of semantic and topology positives are both small. We consider it is because the attributes on Aminer are generated by DeepWalk, a random walk-based algorithm that is biased to hub nodes in the learning procedure.

**5.4.2 Augmentation Probabilities** In this section, we present the impact of two critical data augmentation hyperparameters, i.e., edge dropping  $p_e$  and feature masking  $p_f$ , in Figure 5. We have the following observations. (1) A relatively low corrupt probability 0.3 is desirable to achieve competitive results. (2) When the dropping probability is too large, we face a model degradation because the semantics and/or structures for each metapath-induced view are significantly corrupted, failing to preserve enough augmentation-invariant information. (3) Despite the performance dramatically fluctuating under different probability combinations, the absolute value between the maximum and minimum is generally less than 0.5 except ACM, showing the robustness of HGCML on augmentation probabilities.

**5.4.3 Temperature** The value of temperature  $\tau$  determines the data distribution when measuring the distance between data points in contrasting. As illustrated in Figure 6, we can see that our model is not sensitive to the temperature and have higher scores with lower variance against HeCo, showing its robustness. In addition, we observe that if the value of temperature is smaller, the gap between HGCML-P and HGCML will be larger. It is because the data distribution between positives and negatives will be smoother with the increase in temperature. The observation further proves the effectiveness of the proposed positive sampling strategy, especially with a small temperature.

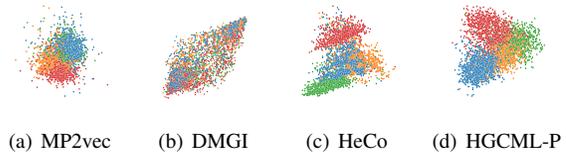


Figure 7: Visualization of node representations on DBLP.

**5.5 Visualization** To profoundly study the expressiveness of HGCML, we visualize the learned node representations of DBLP through  $t$ -SNE. In Figure 7, we visualize node representations obtained from four algorithms, including Meta-path2vec (MP2vec), DMGI, HeCo, and HGCML. As we can see, DMGI presents blurred boundaries between different classes, failing to learn discriminative low-dimensional node representations. For Meta-path2vec and HeCo, despite some types of nodes being categorized clearly, there still exists a large proportion of overlapped data points that cannot be clearly identified. Our model separates nodes into different types, achieving the best performance.

## 6 Conclusion

In this paper, we propose a heterogeneous graph contrastive multi-view learning framework named HGCML. By treating metapaths as data augmentation, we create multi-views without impairing the underlying semantics in HINs. Then, we propose a novel objective that jointly performs intra-metapath and inter-metapath contrasts to model the consistency between metapaths. Specifically, we iteratively utilize graph patches and graph summaries to generate supervision signals to acquire local and global knowledge. To further enhance the quality of representations, we employ a positive sampling strategy that simultaneously considers node attributes and centrality to explicitly select positive samples to mitigate the sampling bias. Experimental results demonstrate the superiority of HGCML across five real-world datasets on node classification and node clustering.

## 7 Acknowledgement

This work was supported in part by the Zhejiang Provincial Natural Science Foundation of China under Grant No. LZ22F020002 and No. LY22F020003, and the National Natural Science Foundation of China under Grant No. 62002226 and No. 62002227.

## References

- [1] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *TNNLS*, vol. 32, no. 1, pp. 4–24, 2020.

- [2] J. Tang, M. Qu, and Q. Mei, "PTE: Predictive Text Embedding through Large-scale Heterogeneous Text Networks," in *KDD*, Aug. 2015, pp. 1165–1174.
- [3] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. v. d. Berg, I. Titov, and M. Welling, "Modeling relational data with graph convolutional networks," in *ESWC*. Springer, 2018, pp. 593–607.
- [4] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous graph attention network," in *WWW*, 2019, pp. 2022–2032.
- [5] C. Zhang, D. Song, C. Huang, A. Swami, and N. V. Chawla, "Heterogeneous graph neural network," in *KDD*, 2019, pp. 793–803.
- [6] X. Fu, J. Zhang, Z. Meng, and I. King, "Magnn: Metapath aggregated graph neural network for heterogeneous graph embedding," in *WWW*, 2020, pp. 2331–2341.
- [7] Z. Hu, Y. Dong, K. Wang, and Y. Sun, "Heterogeneous graph transformer," in *WWW*, 2020, pp. 2704–2710.
- [8] P. Veličković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax," in *ICLR*, 2018.
- [9] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020, pp. 9729–9738.
- [10] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A Simple Framework for Contrastive Learning of Visual Representations," in *ICML*, Nov. 2020, pp. 1597–1607.
- [11] M. I. Belghazi, A. Baratin, S. Rajeswar, S. Ozair, Y. Bengio, A. Courville, and R. D. Hjelm, "Mine: mutual information neural estimation," *arXiv preprint arXiv:1801.04062*, 2018.
- [12] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang, "Deep graph contrastive representation learning," *arXiv preprint arXiv:2006.04131*, 2020.
- [13] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, "What makes for good views for contrastive learning?" *NeurIPS*, vol. 33, pp. 6827–6839, 2020.
- [14] Y. You, T. Chen, Y. Sui, T. Chen, Z. Wang, and Y. Shen, "Graph contrastive learning with augmentations," *NeurIPS*, vol. 33, pp. 5812–5823, 2020.
- [15] C. Park, D. Kim, J. Han, and H. Yu, "Unsupervised attributed multiplex network embedding," in *AAAI*, vol. 34, no. 04, 2020, pp. 5371–5378.
- [16] X. Wang, N. Liu, H. Han, and C. Shi, "Self-supervised heterogeneous graph neural network with co-contrastive learning," in *KDD*, 2021, pp. 1726–1736.
- [17] B. Jing, C. Park, and H. Tong, "Hdmi: High-order deep multiplex infomax," in *WWW*, 2021, pp. 2414–2424.
- [18] Y. Zhu, Y. Xu, H. Cui, C. Yang, Q. Liu, and S. Wu, "Structure-enhanced heterogeneous graph contrastive learning," in *SDM*. SIAM, 2022, pp. 82–90.
- [19] Y. Kalantidis, M. B. Sariyildiz, N. Pion, P. Weinzaepfel, and D. Larlus, "Hard negative mixing for contrastive learning," *NeurIPS*, vol. 33, pp. 21 798–21 809, 2020.
- [20] C.-Y. Chuang, J. Robinson, Y.-C. Lin, A. Torralba, and S. Jegelka, "Debiased contrastive learning," *NIPS*, vol. 33, pp. 8765–8775, 2020.
- [21] Y. Zhu, Y. Xu, Q. Liu, and S. Wu, "An empirical study of graph contrastive learning," in *NeurIPS*, 2021.
- [22] J. Xia, L. Wu, G. Wang, J. Chen, and S. Z. Li, "Progl: Re-thinking hard negative mining in graph contrastive learning," in *ICML*. PMLR, 2022, pp. 24 332–24 346.
- [23] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural Message Passing for Quantum Chemistry," in *ICML*, Jul. 2017, pp. 1263–1272.
- [24] T. N. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," *arXiv:1609.02907 [cs, stat]*, Feb. 2017.
- [25] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *NeurIPS*, vol. 30, 2017.
- [26] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *ICLR*, 2018.
- [27] Z. Wang, Q. Li, D. Yu, and X. Han, "Temporal graph transformer for dynamic network," in *ICANN*. Springer, 2022, pp. 694–705.
- [28] K. Hassani and A. H. Khasahmadi, "Contrastive Multi-View Representation Learning on Graphs," in *ICML*, Nov. 2020, pp. 4116–4126.
- [29] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *CVPR*, 2018, pp. 3733–3742.
- [30] Y. Zhu, Y. Xu, F. Yu, Q. Liu, S. Wu, and L. Wang, "Graph contrastive learning with adaptive augmentation," in *WWW*, 2021, pp. 2069–2080.
- [31] S. Thakoor, C. Tallec, M. G. Azar, M. Azabou, E. L. Dyer, R. Munos, P. Veličković, and M. Valko, "Large-scale representation learning on graphs via bootstrapping," in *ICLR*, 2021.
- [32] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar *et al.*, "Bootstrap your own latent—a new approach to self-supervised learning," *NeurIPS*, vol. 33, pp. 21 271–21 284, 2020.
- [33] Y. Ren, B. Liu, C. Huang, P. Dai, L. Bo, and J. Zhang, "Heterogeneous deep graph infomax," *arXiv preprint arXiv:1911.08538*, 2019.
- [34] C. Wang, S. Zhou, K. Yu, D. Chen, B. Li, Y. Feng, and C. Chen, "Collaborative knowledge distillation for heterogeneous information network embedding," in *WWW*, 2022, pp. 1631–1639.
- [35] X. Jiang, Y. Lu, Y. Fang, and C. Shi, "Contrastive Pre-Training of GNNs on Heterogeneous Graphs," in *CIKM*. New York, NY, USA: Association for Computing Machinery, 2021, pp. 803–812.
- [36] J. Klicpera, S. Weissenberger, and S. Günnemann, "Diffusion improves graph learning," in *NeurIPS*, 2019, pp. 13 366–13 378.
- [37] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *KDD*, 2014, pp. 701–710.
- [38] Y. Dong, N. V. Chawla, and A. Swami, "metapath2vec: Scalable representation learning for heterogeneous networks," in *SIGKDD*, 2017, pp. 135–144.
- [39] T.-y. Fu, W.-C. Lee, and Z. Lei, "Hin2vec: Explore meta-paths in heterogeneous information networks for representation learning," in *CIKM*, 2017, pp. 1797–1806.
- [40] C. Shi, B. Hu, W. X. Zhao, and S. Y. Philip, "Heterogeneous information network embedding for recommendation," *TKDE*, vol. 31, no. 2, pp. 357–370, 2018.