

HIERARCHICAL AND CONTRASTIVE REPRESENTATION LEARNING FOR KNOWLEDGE-AWARE RECOMMENDATION

Bingchao Wu^{1,5}, Yangyuxuan Kang², Daoguang Zan^{1,5†}, Bei Guan^{4,5†}, Yongji Wang^{3,4,5}

¹ Collaborative Innovation Center, Institute of Software, Chinese Academy of Sciences, Beijing, China

² Intel Labs China, Beijing, China

³ State Key Laboratory of Computer Science, Institute of Software, Chinese Academy of Sciences, Beijing, China

⁴ Integrative Innovation Center, Institute of Software, Chinese Academy of Sciences, Beijing, China

⁵ University of Chinese Academy of Sciences, Beijing, China

{paulpigwbc@outlook.com, yangyuxuan.kang@intel.com, {daoguang.guanbei}@iscas.ac.cn, ywang@itechs.iscas.ac.cn }

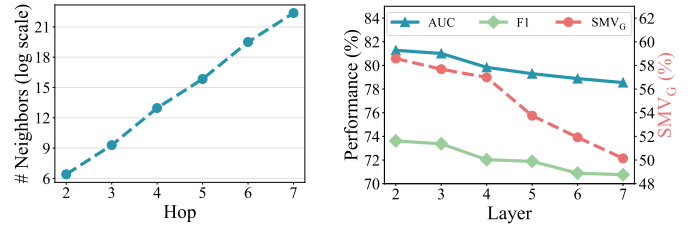
ABSTRACT

Incorporating knowledge graph into recommendation is an effective way to alleviate data sparsity. Most existing knowledge-aware methods usually perform recursive embedding propagation by enumerating graph neighbors. However, the number of nodes' neighbors grows exponentially as the hop number increases, forcing the nodes to be aware of vast neighbors under this recursive propagation for distilling the high-order semantic relatedness. This may induce more harmful noise than useful information into recommendation, leading the learned node representations to be indistinguishable from each other, that is, the well-known over-smoothing issue. To relieve this issue, we propose a Hierarchical and CONTRastive representation learning framework for knowledge-aware recommendation named HiCON. Specifically, for avoiding the exponential expansion of neighbors, we propose a hierarchical message aggregation mechanism to interact separately with low-order neighbors and meta-path-constrained high-order neighbors. Moreover, we also perform cross-order contrastive learning to enforce the representations to be more discriminative. Extensive experiments on three datasets show the remarkable superiority of HiCON over state-of-the-art approaches.

Index Terms— Knowledge-aware recommendation, hierarchical message aggregation, contrastive learning

1. INTRODUCTION

Data sparsity is a significant challenge in the recommendation system. A natural way to alleviate this challenge is incorporating external side information [1]. Knowledge graph (KG), a common and essential source including real-world facts, can



(a) Average number of interacted neighbors. (b) Performance and smoothness metric of KGAT.

Fig. 1: Illustration of the fact that graph neighbors grow exponentially and the over-smoothing issue. Figure (a) counts the average number of interacted neighbors per user at different hops. Since KGAT has to stack as many layers as hops for capturing more complicated semantic relatedness, figure (b) reports recommendation performance (i.e., AUC and F1) and smoothness (i.e., SMV_G) at the different number of layers of KGAT (AUC, F1, SMV_G: the higher, the better).

assist the recommendation system in learning better representations via its rich structural and semantic information [2]. Consequently, it is widely used in existing methods to alleviate the data sparsity issue [3, 4]. Most earlier works [5, 6] utilize knowledge graph embedding (KGE) methods, such as TransR [7] and TransD [8], to pre-train entity embeddings based merely on a single KG triplet (consists of two entities and a connected relation) for enriching item-side semantic representations. However, due to modeling the triplets independently, they may insufficiently capture the high-order relatedness among nodes [4].

Recently, most existing studies develop knowledge-aware recommendation methods based on graph neural networks (GNN) to explore high-order relatedness for accurate recommendation, such as KGAT [3] and KGCN [9]. A typical GNN paradigm usually performs recursive message propagation by

† Corresponding author.

enumerating graph neighbors over the unified graph consisting of the user-item bipartite graph and item-side knowledge graph. There exists a fact that many nodes are aware of exponentially growing neighbors over the unified graph as the number of hops increases (shown as Figure 1a). Thus, to capture the high-order semantic relatedness with these neighbors, this GNN paradigm has to force the nodes to be aware of the information of all exponentially expanding neighbors because of its recursive propagation mechanism. This may bring in more infernal noise than useful information for node representation learning, which greatly harms the recommendation performance [10]. Several existing works attribute this phenomenon to the over-smoothing issue [11]. As shown in Figure 1b, we take KGAT (a representative knowledge-aware recommendation model) as an example to intuitively describe this issue. To reflect the overall smoothness of node representations, we employ SMV_G proposed in DAGNN [11], which is the average pairwise Euclidean distance of two-node representations. We find that the overall smoothness metric SMV_G of node representations learned from KGAT continuously degrades along with the increase of the number of message propagation layers, leading to the degradation of performance on AUC and F1.

In this paper, we propose an effective framework named HiCON for knowledge-aware recommendation. It alleviates the over-smoothing issue in two aspects: 1) selecting and propagating a bundle of valuable neighbors to the central nodes rather than considering all exponentially increasing neighbors to reduce noise disturbance; 2) enhancing the self-discrimination of node representations in the latent space. For the first aspect, we propose a hierarchical message aggregation mechanism consisting of two parts, i.e., low- and high-order message aggregations, to avoid interacting with exponentially expanding neighbors. The low-order aggregation enriches the node representations at the low level by aggregating local graph neighbors over the unified graph. The high-order aggregation aims to learn high-level node representations by first selecting valuable high-order neighbors guided by well-designed meta-paths and then propagating the information of selected neighbors to the central node. For the second aspect, to enhance the discriminativeness of learned node representations, we perform cross-order contrastive learning between low- and high-level semantic representations derived from the message aggregation process. Extensive experiments on three benchmark datasets show that HiCON can greatly improve the recommendation performance and meanwhile alleviate over-smoothing.

Overall, our work makes two major contributions: 1) We propose a hierarchical and contrastive representation learning framework for knowledge-aware recommendation, which alleviates over-smoothing to improve performance from two aspects: i.e., avoiding exponential expansion of interacted neighbors and enhancing the discriminativeness of node representations. 2) Extensive experiments on three datasets val-

idate the effectiveness of HiCON in providing accurate recommendations and the ability to alleviate over-smoothing.

2. RELATED WORK

Knowledge-aware Recommendation. Incorporating the structural and semantic information of knowledge graph (KG) into recommendation is effective to alleviate data sparsity [12]. Existing knowledge-aware recommendation models can be roughly divided into three categories, i.e., embedding-based methods, path-based methods and propagation-based methods. For the embedding-based methods [5, 6, 13], they primarily adopt the knowledge graph embedding (KGE) technique, e.g., TransR [7] and TransD [8], to learn user and item representations by a transition constraint. For the path-based methods [14, 15], they exploit the connectivity patterns among items guided by meta-paths to provide external semantic dependence to make recommendations. For the propagation-based methods [3, 16, 17, 16], they mainly perform recursive message propagation over KG to refine the user and item representations for making accurate recommendations. However, the maximum propagation depth of these models is usually limited to three to avoid the explosive increase of interacted neighbors [18], which hinders the ability of the models to explore high-order relatedness among nodes. Differently, HiCON explicitly derives low- and high-order semantic relatedness by two message aggregation components linked in cascade, helping the model fully exploit structural and semantic information of KG.

Contrastive Learning for Recommendation. Contrastive learning aims to learn better user and item representations by pulling positive pairs closer and pushing negative pairs away [19]. A great number of recommendation tasks, such as social recommendation [20] and knowledge-aware recommendation [17, 4, 21], benefit from this mechanism to improve recommendation performance. For example, for the social recommendation, SEPT [20] exploits social relations between users to generate two complementary views based on tri-training for contrastive learning. As for knowledge-aware recommendation, MCCLK [17] performs a multi-level cross-view contrastive learning based on the local and global views of knowledge graph to learn effective user and item representations. However, existing knowledge-aware models perform contrastive learning based on vanilla GNNs equipped with limited propagation layers, which have defects in capturing deeper semantic information. On the contrary, we propose two message aggregation components linked hierarchically to encode shallow and deep semantic representations, which are treated as positive pairs for contrastive learning.

3. PRELIMINARIES

In this section, we introduce some basic concepts and the knowledge-aware recommendation task in this paper.

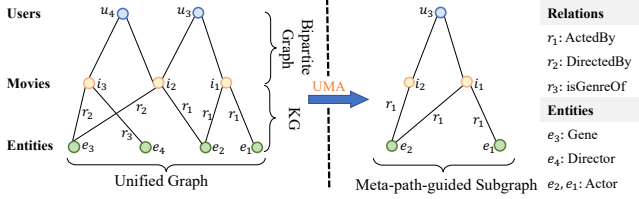


Fig. 2: An example of generating a meta-path-guided subgraph from a unified graph where UMA is the meta-path User-Movie-Actor.

Bipartite Graph. Let $\mathcal{U} = \{u_1, u_2, \dots, u_m\}$ and $\mathcal{I} = \{v_1, v_2, \dots, v_n\}$ be the sets of users and items in the recommendation system, where m and n are the numbers of users and items. We form a bipartite graph \mathcal{G}_r from user-item interactions where its edges mean that users interact with items.

Knowledge Graph. Apart from the bipartite graph \mathcal{G}_r , we also provide a knowledge graph \mathcal{G}_k , and its formal definition is $\{(h, r, t) | h, t \in \mathcal{E}, r \in \mathcal{R}\}$ where h, r and t are the head, relation and tail of triplet (h, r, t) that is the basic unit in knowledge graph. For example, (Brad Pitt, ActorOf, Fight Club) represents the fact that Brad Pitt is an actor in the movie Fight Club. Note that \mathcal{R} contains bidirectional relations between nodes, such as ActorOf and its reverse relation ActedBy, to fully reveal the facts in KG. Besides, there exists an intersection between item set \mathcal{I} and entity set \mathcal{E} , enabling the recommendation model to incorporate knowledge graph.

Unified Graph. Following the collaborative knowledge graph in KGAT [3], we merge the bipartite graph \mathcal{G}_r and knowledge graph \mathcal{G}_k to construct a unified graph $\mathcal{G}_c = \{(h, r, t) | h, t \in \mathcal{E}', r \in \mathcal{R}'\}$, where $\mathcal{E}' = \mathcal{E} \cup \mathcal{U} \cup \mathcal{I}$ and $\mathcal{R}' = \mathcal{R} \cup \{r_{ui}\}$. r_{ui} is the user-item interaction relation.

Meta-path-guided Subgraph. Let $\mathcal{M} = \{m_1, \dots, m_k\}$ be the set of k meta-paths. A meta-path $m \in \mathcal{M}$ is denoted as a path in the form of $m = A_1 \xrightarrow{R_1} A_2 \dots \xrightarrow{R_l} A_{l+1}$ (abbreviated to $A_1 A_2 \dots A_{l+1}$) [22]. It describes a composite relation $R = R_1 \circ R_2 \dots R_l$ between two nodes with type A_1 and A_{l+1} where \circ is the composition operator on relations. For example, Movie-Actor-Movie (MAM) demonstrates the semantic “co-actor relationships between movies”. Then, a meta-path-guided subgraph is defined as the combination of all path instances derived by the given meta-path. For example, given a unified graph and a meta-path User-Movie-Actor (UMA), as shown in Figure 2, we conduct a subgraph including the path instances $u_3-i_2-e_2$, $u_3-i_1-e_2$ and $u_3-i_1-e_1$.

knowledge-aware Recommendation Task. Given a unified graph \mathcal{G}_c , the goal is to recommend items to target users based on the relevance of their representations derived from \mathcal{G}_c .

4. METHOD

4.1. Overall Framework

The overall framework of HiCON is shown in Figure 3. A hierarchical message aggregation module, consisting of two

parts linked in a cascaded manner, learns low- and high-level node representations over the unified graph. The former part aggregates the information of nodes’ local neighbors by a limited number of common propagation layers. The latter part reveals the valuable high-order semantic relations by employing graph convolutional networks over multiple subgraphs generated by the guidance of well-designed meta-paths. Note that both parts help the model avoid propagating a vast number of useless neighbors to the central node. Next, to further alleviate the over-smoothing issue by learning more discriminative node representations, we employ a cross-order contrastive learning module by contrasting low- and high-level representations of the same node.

4.2. Hierarchical Message Aggregation

Most existing GNN-based knowledge-aware methods, such as KGAT [3] and KGCN [9], encode shallow (low-level) and deep (high-level) semantic information by a simple and naive recursive propagation mechanism. However, as the number of propagation layers increases, many central nodes are aware of more unrelated neighbors. This may make the model encounter the over-smoothing issue caused by bringing in a lot of noise information [23]. To alleviate this issue, we decompose the semantic representation learning into two message aggregation components combined in a hierarchical manner.

Low-order Message Aggregation. The unified graph covers not only the user interaction behaviors about items, but also the item-wise semantic relatedness in knowledge graph. To fully derive the shallow semantic information from the unified graph, we propose a dual propagation layer consisting of two embedding propagation layers for the bipartite graph and knowledge graph, respectively.

For the bipartite graph, we utilize a light embedding propagation layer to aggregate the neighbors based on the topological structure of the bipartite graph. Formally, the definition can be summarized as follow:

$$\alpha_{i,j} = \frac{1}{\sqrt{|\mathcal{N}_{u_i}^{(v)}|} \sqrt{|\mathcal{N}_{v_i}^{(u)}|}}, \quad \mathbf{h}_{v_i}^{(k+1)} = \sum_{j \in \mathcal{N}_{v_i}^{(u)}} \alpha_{i,j} \mathbf{e}_{u_j}^{(k)}, \quad (1)$$

where $\mathcal{N}_{u_i}^{(v)}$ is the set of items interacted by user u_i , and $\mathbf{h}_{v_i}^{(k+1)}$ is the user-side representation of item v_i in the $k+1$ -th layer. For knowledge graph, due to there are various types of relations among nodes (e.g., ActedBy and DirectedBy), we propose a relation-aware embedding propagation layer relying on the attention mechanism to exploit the structural and semantic information of KG. Concretely, given an item v_i , and a triplet set $\mathcal{N}_{v_i}^{(e)} = \{(v_i, r, t) | (v_i, r, t) \in \mathcal{G}_k\}$ with the item v_i as the head entity, we perform relation-aware attention to aggregate the neighbors within the set: $\mathbf{g}_{v_i}^{(k+1)} = \sum_{(v_i, r, t) \in \mathcal{N}_{v_i}^{(e)}} \beta_{v_i, t}^{(k)} (\mathbf{e}_r^{(k)} \odot \mathbf{e}_t^{(k)})$ where \odot is the element-wise multiplication. $\mathbf{g}_{v_i}^{(k+1)}$ is the entity-side representation for item v_i in the $k+1$ -th layer, and $\mathbf{e}_t^{(k)}$ and $\mathbf{e}_r^{(k)}$

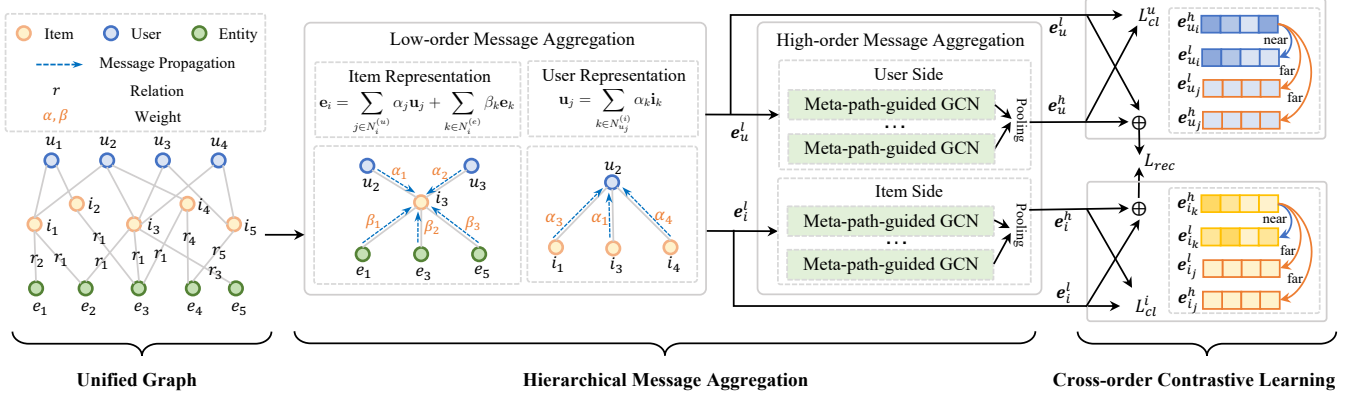


Fig. 3: The framework of our proposed HiCON model.

are the representations for the tail entity t and relation r in the k -th layer. The attention weight $\beta_{v_i,t}^{(k)}$ is normalized by a softmax function based on the importance coefficient $\pi_{v_i,t}^{(k)}$:

$$\beta_{v_i,t}^{(k)} = \frac{\exp(\pi_{v_i,t}^{(k)})}{\sum_{(v_i,r,t) \in \mathcal{N}_{v_i}^{(e)}} \exp(\pi_{v_i,t}^{(k)})}, \quad \pi_{v_i,t}^{(k)} = \mathbf{s}_{v_i}^T \mathbf{s}_t, \quad (2)$$

where $\mathbf{s}_{v_i} = \mathbf{e}_{v_i}^{(k)} \odot \mathbf{e}_r^{(k)} / \|\mathbf{e}_{v_i}^{(k)} \odot \mathbf{e}_r^{(k)}\|$ considers the influence of relation r . $\|\cdot\|$ is the l_2 -norm, and similarly for \mathbf{s}_t .

Since the items are regarded as the bridge connecting the bipartite graph and knowledge graph, we merge both collaborative signals and semantic information of KG to further refine the item representations. Concretely, given the user-side representation $\mathbf{h}_{v_i}^{(k+1)}$ and entity-side representation $\mathbf{g}_{v_i}^{(k+1)}$ of item v_i , we combine these two representations linearly with the sum-pooling operation to acquire the refined item representation $\mathbf{e}_{v_i}^{(k+1)}$ in the $k+1$ layer: $\mathbf{e}_{v_i}^{(k+1)} = \mathbf{h}_{v_i}^{(k+1)} + \mathbf{g}_{v_i}^{(k+1)}$. Finally, to fully exploit the shallow semantic relatedness, we stack two dual propagation layers and sum the output of layers up as the low-level item representations: $\mathbf{e}_{v_i} = \mathbf{e}_{v_i}^{(0)} + \mathbf{e}_{v_i}^{(1)} + \mathbf{e}_{v_i}^{(2)}$, and similarly for the low-level user representations \mathbf{e}_u .

High-order Message Aggregation: In this section, we perform the meta-path-guided message aggregation to derive useful high-order relatedness (shown in Figure 4). Concretely, given a meta-path m and the unified graph \mathcal{G}_c , we collect all path instances derived by the meta-path from \mathcal{G}_c to construct a subgraph \mathcal{G}_s^m (more details on the construction of this subgraph refer to Section 3). To capture the high-order semantic relatedness from \mathcal{G}_s^m , we use a graph convolutional network (GCN) where its single graph convolutional layer employs a nonlinear transformation after the combination of central users and their neighbors:

$$\mathbf{z}_{m,u_i}^{(k+1)} = \sigma \left(\sum_{j \in \mathcal{N}_{u_i}^{(m)}} \frac{1}{\sqrt{|\mathcal{N}_{u_i}^{(m)}|} \sqrt{|\mathcal{N}_j^{(m)}|}} \mathbf{z}_{m,j}^{(k)} \mathbf{W}_m^{(k)} \right), \quad (3)$$

where $\sigma(\cdot)$ is a non-linear activation function, and $\mathbf{W}_m^{(k)}$ is a weight matrix in the k -th layer. $\mathcal{N}_{u_i}^{(m)}$ is the set of neighbors

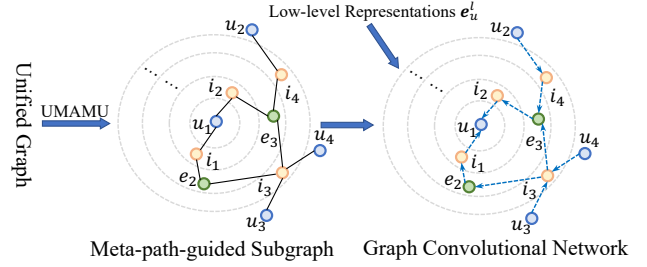


Fig. 4: An example of a meta-path-guided GCN where UMAMU denotes the meta-path User-Movie-Actor-Movie-User. The concentric circles denote different hops of the user.

for the user u_i in \mathcal{G}_s^m . $\mathbf{z}_{m,u_i}^{(k+1)}$ is the representation of user u_i encoded from \mathcal{G}_s^m in the $k+1$ -th layer. We select the output of the last layer of GCN as the final high-level representations $\mathbf{z}_{v_i}^m$ rather than concatenating or summing the output of all layers to reduce the noise disturbance: $\mathbf{z}_{v_i}^m = \mathbf{z}_{m,u_i}^{(l)}$, where l is the layer number in GCN. To effectively exploit diverse high-order relatedness among nodes, we develop multiple meta-paths for constructing corresponding subgraphs, and then utilize independent GCNs to encode multiple high-level user representations. Finally, a mean-pooling operation is applied to combine these representations: $\mathbf{z}_{u_i} = \sum_{m \in \mathcal{M}_u} \mathbf{z}_{u_i}^m$, where \mathcal{M}_u is the set of meta-paths on the user side, and similarly for the high-level item representations \mathbf{z}_v .

4.3. Cross-Order Contrastive Learning

In this section, we propose a cross-order contrastive learning module to learn discriminative node representations by distinguishing positive pairs from negative ones. Concretely, considering the low-level item vectors \mathbf{e}_v and high-order item vectors \mathbf{z}_v as positive pairs, we present an item-side contrastive learning loss to minimize the distance between the positive pairs: $\mathcal{L}_{cl}^I = - \sum_{v_i \in \mathcal{I}} \log \frac{\exp((\mathbf{e}_{v_i} \cdot \mathbf{z}_{v_i})/\tau)}{\sum_{v_j \in \mathcal{C}_{v_i} \cup \{v_i\}} \exp((\mathbf{e}_{v_i} \cdot \mathbf{h}_{z_j})/\tau)}$,

where \mathcal{C}_{v_i} is a negative sample set, consisting of items in the batch except item v_i , and τ is the temperature parameter,

which is set to 0.6. The user-side loss \mathcal{L}_{cl}^U is calculated similarly. After that, we linearly combine the item-side loss \mathcal{L}_{cl}^I and user-side loss \mathcal{L}_{cl}^U to get the final object for cross-order contrastive learning: $\mathcal{L}_{cl} = \mathcal{L}_{cl}^U + \mathcal{L}_{cl}^I$.

4.4. Optimization and Prediction

To optimize HiCON, we adopt a pairwise Bayesian personalized ranking (BPR) loss \mathcal{L}_{bpr} equipped with the aforementioned contrastive loss \mathcal{L}_{cl} as the final loss \mathcal{L}_{HiCON} by the weighted sum operation:

$$\mathcal{L}_{HiCON} = \mathcal{L}_{bpr} + \lambda \mathcal{L}_{cl}, \quad (4)$$

where λ is a hyper-parameter to control the balance between the losses. \mathcal{L}_{bpr} imposes users to have higher predicted scores with interacted items than non-interacted ones, which is defined as: $\mathcal{L}_{bpr} = \sum_{u \in \mathcal{U}} \sum_{i \in \mathcal{N}_u} \sum_{i' \notin \mathcal{N}_u} -\log \sigma(\hat{y}_{u,i} - \hat{y}_{u,i'})$, where \mathcal{N}_u denotes the item neighbors of user u in the bipartite graph. The probability $\hat{y}_{u,i} = [\mathbf{e}_u \parallel \mathbf{z}_u]^T [\mathbf{e}_i \parallel \mathbf{z}_i]$ is calculated by the dot product between the concatenated user and item representations where \parallel denotes the concatenate operation, which is also used to predict user’s preferred items during the inference phase.

5. EXPERIMENT

5.1. Experimental Settings

Datasets. We conduct experiments on three widely used benchmark datasets for recommendation, i.e., LastFM¹, Book-Crossing² and MovieLens-1M³, which are collected from the music, book, and movie domains, respectively. The details of the datasets are included in Appendix A.1 (including statistics and data preprocessing).

Baseline Methods. To verify the effectiveness of HiCON, we select two representative groups of recommendation models as baseline methods: collaborative filtering-based group (BPRMF [24], LightGCN [25]) and knowledge-aware group. In the knowledge-aware group, there are three classic types of methods: an embedding-based method (CKE [5]), a path-based method (PER [14]), and propagation-based methods (RippleNet [10], KGCN [9], KGNN-LS [26], KGAT [3], CKAN [27], KGIN [16], MCCLK [17], KGIC [4]). For the implementation of baselines, please refer to Appendix A.2.

Evaluation Metrics & Parameter Settings. We use two classical metrics *Area Under Curve* (AUC) and F1, which are widely used in click-through rate (CTR) prediction, to evaluate all models. We tailor multiple meta-paths for each relation on both item and user sides. Specifically, the item-side meta-path set \mathcal{M}_v includes Item-User-Item-Entity-Item

Table 1: Performance comparison of HiCON and baselines. The best and second-best results are indicated by bold fonts and underlines, and * indicates statistically significant improvements over the best baseline (t-test with $p < 0.01$).

Model	Book-Crossing		MovieLens-1M		Last.FM	
	AUC	F1	AUC	F1	AUC	F1
BPRMF	0.6583	0.6117	0.8920	0.7921	0.7563	0.7010
LightGCN	0.6134	0.6469	0.8800	0.8091	0.8300	0.7439
CKE	0.6759	0.6235	0.9065	0.8024	0.7471	0.6740
PER	0.6048	0.5726	0.7124	0.6670	0.6414	0.6033
RippleNet	0.7211	0.6472	0.9190	0.8422	0.7762	0.7025
KGCN	0.6841	0.6313	0.9090	0.8366	0.8027	0.7086
KGNN-LS	0.6762	0.6314	0.9140	0.8410	0.8052	0.7224
KGAT	0.7314	0.6544	0.9140	0.8440	0.8293	0.7424
CKAN	0.7420	0.6671	0.9082	0.8410	0.8418	0.7592
KGIN	0.7273	0.6614	0.9190	0.8441	0.8486	0.7602
MCCLK	0.7625	0.6777	<u>0.9351</u>	<u>0.8631</u>	<u>0.8763</u>	<u>0.8008</u>
KGIC	<u>0.7749</u>	<u>0.6812</u>	0.9252	0.8559	0.8592	0.7753
HiCON	0.8045*	0.7169*	0.9410*	0.8718*	0.9045*	0.8186*

(UIEI), Item-Entity-Item-User-Item (IEIUI), and Item-User-Item-User-Item (IUIUI). The user-side set \mathcal{M}_u contains User-Item-User-Item-User (UIUIU) and User-Item-Entity-Item-User (UIEIU). The weight λ for controlling the balance of losses is set to 0.01 on MovieLen-1M, and 0.05 on the remaining datasets.

5.2. Experimental Results

Overall Performance Comparison. Table 1 shows the results of HiCON and baselines on all datasets. We find that HiCON achieves the best performance compared with all baselines. This is because the hierarchical message aggregation effectively derives different-level semantic information from the local semantic relatedness and meta-path-guided high-order relatedness. Besides, the cross-order contrastive learning module can improve the self-discrimination of node representations. We also find that, due to the rich structural and semantic information in KG, most knowledge-aware recommendation models outperform collaborative filtering-based methods without merging KG (BPRMF and LightGCN) by a large margin. In addition, the superior performance of most propagation-based methods over embedding-based (CKE) or path-based ones (PER), which may attribute to recursive propagation mechanisms that mine deeper semantic relatedness.

Ablation Study. To analyze the effectiveness of key components in HiCON, we compare the performance of HiCON with its three variants: 1) w/o Low removes the low-order message aggregation of HiCON; 2) w/o High is the variant of HiCON without the high-order message aggregation; 3) w/o CL wipes out the cross-order contrastive learning module. The results are summarized in Figure 5. First, after removing low- and high-order message aggregations, the performance of HiCON significantly decreases. It shows that both aggregations derive useful external semantic knowledge for learn-

¹<https://grouplens.org/datasets/hetrec-2011/>

²<http://www2.informatik.uni-freiburg.de/cziegler/BX/>

³<https://grouplens.org/datasets/movielens/1m/>

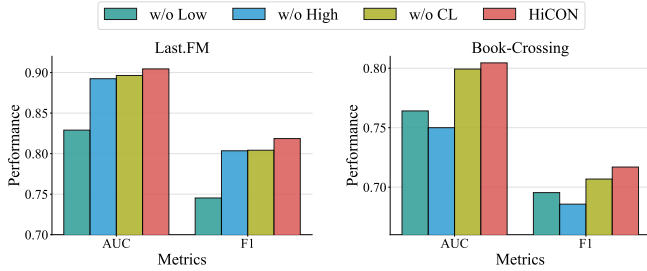


Fig. 5: The performance of HiCON and its variants on the LastFM and Book-Crossing datasets.

ing better user and item representations. Second, hierarchically combining two message aggregations can also improve performance, indicating the low- and high-level representations are mutually enhanced to further provide more precise recommendations. Third, contrastive learning can boost the performance, indicating the discriminative user and item representations are helpful for making accurate recommendation.

In addition, we visualize the feature distribution of user representations in Appendix A.3 to indicate the ability of each module of HiCON to alleviate over-smoothing. The results show that the hierarchical message aggregation and cross-order contrastive learning modules both derive more uniform feature distributions than the variant based on recursive message propagation, indicating both modules help HiCON learn more diverse user preferences to alleviate over-smoothing.

Performance w.r.t Hop (Layer) Number. Here we compare the performance of HiCON with other baselines under different layers (hops). The results are shown in Figure 6. We can observe some interesting findings. First, the performance of some baselines (i.e., KGAT, RippleNet and MCCLK) gradually declines along with the increase of the layer (hop) number. This is because they are aware of exponentially growing neighbors during recursive message propagation, which is harmful to personalized node representation learning. On the contrary, CKAN and KGIC achieve stable performance at different layers (hops). This is because they select a fixed number of neighbors by random sampling. At last, compared with all baselines, HiCON achieves the best performance at all layers (hops), indicating the effectiveness of HiCON in capturing both low and high-order semantic relatedness and alleviating the over-smoothing issue.

Additional Experiments. We perform experiments in Appendix A.4 and A.5 to indicate the effectiveness of HiCON in modeling high-order semantic relatedness and at various interaction sparsity levels, respectively. Besides, we also analyze the impact of hyperparameter γ in Appendix A.6, and find that over-focusing the contrastive loss \mathcal{L}_{cl} may hinder the model optimization process.



Fig. 6: The performance of HiCON and baselines at different layers (hops) on the LastFM dataset.

6. CONCLUSION

In this paper, we propose a hierarchical and contrastive representation learning framework to alleviate the over-smoothing issue for effectively absorbing the structural and semantic information of KG. It relieves this issue by avoiding the exponential expansion of neighbors and enhancing the self-discrimination of node representations. We first propose a hierarchical message aggregation mechanism to explore different-level semantic relatedness via aggregating local neighbors and meta-path-guided high-order neighbors. Besides, we further perform cross-order contrastive learning via maximizing the consistency between low- and high-level views. Extensive experiments on three datasets show the superiority of HiCON over state-of-the-art methods and the ability of alleviating over-smoothing.

7. REFERENCES

- [1] Zhu Sun et al., “Research commentary on recommendations with side information: A survey and research directions,” *ECRA*, vol. 37, 2019.
- [2] Antoine Bordes et al., “Translating embeddings for modeling multi-relational data,” in *NIPS*, 2013.
- [3] Xiang Wang et al., “Kgat: Knowledge graph attention network for recommendation,” in *KDD*, 2019.
- [4] Ding Zou et al., “Improving knowledge-aware recommendation with multi-level interactive contrastive learning,” in *CIKM*, 2022.
- [5] Fuzheng Zhang et al., “Collaborative knowledge base embedding for recommender systems,” in *KDD*, 2016.
- [6] Hongwei Wang et al., “Dkn: Deep knowledge-aware network for news recommendation,” in *WWW*, 2018.
- [7] Yankai Lin et al., “Learning entity and relation embeddings for knowledge graph completion,” in *AAAI*, 2015.
- [8] Guoliang Ji et al., “Knowledge graph embedding via dynamic mapping matrix,” in *ACL*, 2015.
- [9] Hongwei Wang et al., “Knowledge graph convolutional networks for recommender systems,” in *WWW*, 2019.
- [10] Hongwei Wang et al., “Ripplenet: Propagating user preferences on the knowledge graph for recommender systems,” in *CIKM*, 2018.

- [11] Meng Liu et al., “Towards deeper graph neural networks,” in *KDD*, 2020.
- [12] Maksims Volkovs et al., “Dropoutnet: Addressing cold start in recommender systems,” in *NIPS*, 2017.
- [13] Xin Xin et al., “Relational collaborative filtering: Modeling multiple item relations for recommendation,” in *SIGIR*, 2019.
- [14] Xiao Yu et al., “Personalized entity recommendation: A heterogeneous information network approach,” in *WSDM*, 2014.
- [15] Huan Zhao et al., “Meta-graph based recommendation fusion over heterogeneous information networks,” in *KDD*, 2017.
- [16] Xiang Wang et al., “Learning intents behind interactions with knowledge graph for recommendation,” in *WWW*, 2021.
- [17] Ding Zou et al., “Multi-level cross-view contrastive learning for knowledge-aware recommender system,” in *SIGIR*, 2022.
- [18] Yu Wang and Tyler Derr, “Tree decomposed graph neural network,” in *CIKM*, 2021.
- [19] Junliang Yu et al., “Are graph augmentations necessary? simple graph contrastive learning for recommendation,” in *SIGIR*, 2022.
- [20] Junliang Yu et al., “Socially-aware self-supervised tri-training for recommendation,” in *KDD*, 2021.
- [21] Yuhao Yang et al., “Knowledge graph contrastive learning for recommendation,” in *SIGIR*, 2022.
- [22] Yizhou Sun et al., “Pathsim: Meta path-based top-k similarity search in heterogeneous information networks,” *PVLDB*, vol. 4, no. 11, 2011.
- [23] Yongyu Wang, Zhiqiang Zhao, and Zhuo Feng, “Scalable graph topology learning via spectral densification,” in *WSDM*, 2022, pp. 1099–1108.
- [24] Steffen Rendle et al., “Bpr: Bayesian personalized ranking from implicit feedback,” in *UAI*, 2009.
- [25] Xiangnan He et al., “Lightgcn: Simplifying and powering graph convolution network for recommendation,” in *SIGIR*, 2020.
- [26] Hongwei Wang et al., “Knowledge-aware graph neural networks with label smoothness regularization for recommender systems,” in *KDD*, 2019.
- [27] Ze Wang et al., “Ckan: collaborative knowledge-aware attentive network for recommender systems,” in *SIGIR*, 2020.
- [28] Van der Maaten et al., “Visualizing data using t-sne,” *JMLR*, vol. 9, no. 11, 2008.
- [29] Botev Z I et al., “Kernel density estimation via diffusion,” *The annals of Statistics*, vol. 38, no. 5, 2010.

Table 2: Statistics of the LastFM, Book-Crossing, and MovieLens-1M datasets.

		Last.FM	Book-Crossing	MovieLens-1M
Bipartite Interaction	# users	1,872	17,860	6,036
	# items	3,846	14,967	2,445
	# interactions	42,346	139,746	753,772
Knowledge Graph	# entities	9,366	77,903	182,011
	# relations	60	25	12
	# triplets	15,518	151,500	1,241,996

A. APPENDIX

A.1. Datasets

Here we introduce the details and preprocessing of datasets. Last.FM⁴ is collected from Last.fm music platform where musical tracks are regarded as items. Book-Crossing⁵ is a dataset consisting of ratings about books in the Book-Crossing community. MovieLens-1M⁶ is a classical dataset collecting movie ratings from the MovieLens website for movie recommendation. Following previous works [10, 4], we omit the user-item pairs whose ratings are less than the pre-defined threshold (it is set to 4 for MovieLens-1M and 1 for the remaining datasets) and treat the other pairs as positive samples with label 1. The negative instances are randomly sampled from unobserved items with the same number of positive ones for each user. As for the knowledge graph construction, we follow RippleNet [10] to select item-related triplets from Microsoft Satori⁷ for each dataset. The statistics of these datasets are summarized in Table 2.

A.2. Implementation of Baselines

Here we introduce the details of baselines used in this paper, which include two groups of representative methods: i.e., collaborative filtering-based methods (BPRMF, LightGCN) and knowledge-aware methods. In the knowledge-aware methods, we employ three groups of typical methods: i.e., embedding-based methods (CKE), path-based methods (PER), and propagation-based methods (RippleNet, KGCN, KGNN-LS, KGAT, CKAN, KGIN, MCCLK, KGIC):

- *BPRMF* [24] is a representative matrix factorization-based method optimized by the Bayesian personalized ranking (BPR) object.
- *LightGCN* [25] is a simplified graph convolutional network (GCN) by removing the feature transformation and non-linear activation.
- *CKE* [5] is an embedding-based recommendation method that exploits structural, textual and visual information under a Bayesian framework.

⁴<https://grouplens.org/datasets/hetrec-2011/>

⁵<http://www2.informatik.uni-freiburg.de/cziegler/BX/>

⁶<https://grouplens.org/datasets/movielens/1m/>

⁷<https://searchengineland.com/library/bing/satori>

- *PER* [14] is a path-based method that extracts meta-path-aware features to improve the ability to capture the connectivity between users and items.
- *RippleNet* [10] is a classical propagation-based method by propagating the user preference along KG links to discover potential interests.
- *KGCN* [9] aggregates selected neighbors to learn both structural and semantic information of KG based on the graph convolutional network.
- *KGNN-LS* [26] is a propagation-based method that equips GNNs with a label smoothness regularization mechanism to enrich node representations.
- *KGAT* [3] is a propagation-based method that recursively integrates neighbors over the unified graph with an attention strategy.
- *CKAN* [27] independently encodes collaborative information from the bipartite graph and semantic relationships from knowledge graphs relying on distinct message propagation strategies.
- *KGIN* [16] considers fine-grained user intents to profile user-item relationships, and uses relational paths to capture the long-range connectivity.
- *MCCLK* [17] utilizes a multi-level contrastive learning mechanism based on local and global semantic views to enrich user and item representations.
- *KGIC* [4] performs multi-level interactive contrastive learning based on layer-wise augmented views.

A.3. Representation Visualization

In this section, we conduct experiments to analyze the over-smoothing issue by visualizing the user representations learned by HiCON and its two variants:

- HiCON_{w/o Hie} employs the recursive message propagation with the same receptive field of HiCON to replace the hierarchical message aggregation.
- HiCON_{w Hie} merely utilizes the hierarchical message aggregation of HiCON and does not perform cross-order contrastive learning.

Concretely, we first map the user representations into 2-dimensional normalized vectors by t-SNE [28]. Then, we plot feature distributions with Gaussian kernel Density Estimation (KDE) [29] in \mathbb{R}^2 . To present the representations more intuitively, we plot density estimations on angles (i.e., $\arctan 2(y, x)$ for each point (x, y) in the unit hypersphere) [19]. The visualization is shown in Figure 7, from which we observe some findings. First, we observe that the

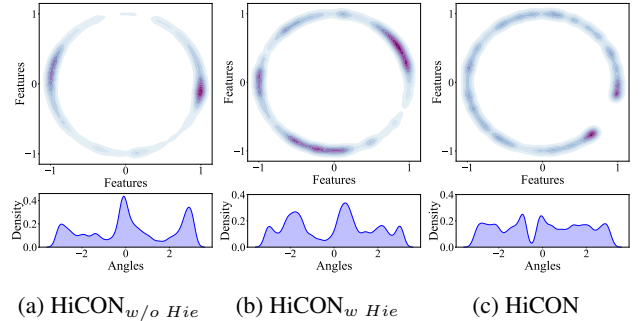


Fig. 7: The visualization of user representations derived by HiCON and its variants on the LastFM dataset. The feature distributions (the darker the color, the more points are distributed into the area [19].) and density estimations (the sharper the curve, the more points are clustered together.) of user representations are plotted in the unit hypersphere.

user representations derived from HiCON_{w/o Hie} are mapped into a narrow area. This may be because the recursive message propagation forces the users to be aware of vast graph neighbors, causing the users to share similar representations in the latent space. Second, we observe a more uniform feature distribution and a more flat density curve when replacing the recursive message propagation with hierarchical message aggregation. This may be because the hierarchical aggregation interacts with a collection of valuable neighbors rather than exponentially growing graph neighbors to reduce noise disturbance. Third, compared with HiCON_{w Hie}, the distribution and the curve of HiCON respectively become more uniform and flat. This may be because performing cross-order contrastive learning enhances the self-discrimination of user representations.

A.4. Performance w.r.t High-order Semantic Relatedness Modeling

In this section, we conduct experiments to validate the effectiveness of HiCON in modeling high-order semantic relatedness. To achieve this goal, we compare the performance of HiCON with two mainstream baselines, i.e., KGAT [3] and MCCLK [17]. To successfully capture the high-order semantic relatedness over the unified graph, we set the number of layers (hops) of all models to six for a fair comparison. More precisely, we conduct the following models:

- KGAT_{high} uses the output of the last layer of KGAT as the final representations for recommendation;
- MCCLK_{high} regards the output of the last layer of MCCLK as the final representations;
- HiCON_{high} merely utilizes the output of the high-order message aggregation for recommendation.

The results of recommendation performance are shown in Figure 8. We observe that HiCON_{high} outperforms the other

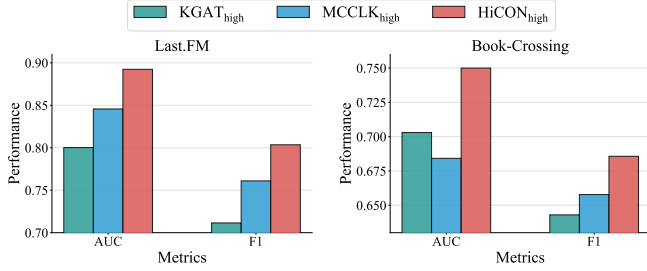


Fig. 8: The performance of KGAT, MCCLK, and HiCON in modeling the high-order semantic relatedness on the LastFM and Book-Crossing datasets.

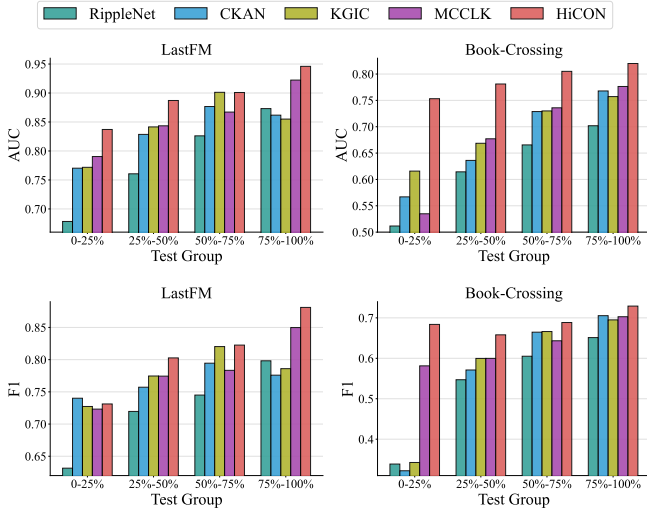


Fig. 9: The performance of HiCON and baselines at different interaction sparsity levels. 25%-50% is a user group sampled from the test set, including the users whose number of interacted items ranges from the 25th to 50th percentile.

two models by a large margin, indicating that our proposed model effectively encodes high-order relatedness to derive high-quality deep semantic representations. This may be because the hierarchical message propagation involved in HiCON selects and propagates a bundle of valuable neighbors to the central nodes rather than considering all exponentially increasing neighbors to reduce noise disturbance.

A.5. Performance w.r.t Interaction Sparsity Levels

In this section, to evaluate the recommendation performance of HiCON at different interaction sparsity levels, we compare HiCON with several representative knowledge-aware baselines, i.e., RippleNet, CKAN, KGIC, and MCCLK. Specially, we split the users included in the test set of LastFM into four groups according to the number of their interacted items. For example, the user group 25%-50% contains the users whose interaction number ranges from 25th to 50th percentile. The

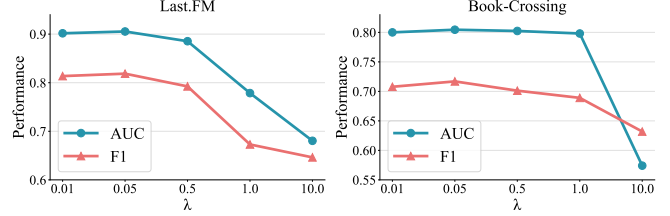


Fig. 10: The impact of the hyper-parameter λ .

group-wise results are shown in Figure 9. We observe the following findings: 1) HiCON outperforms most baselines for the inactive users that interact with a small number of items (e.g., group 25%-50%). This may be because our model can effectively exploit more complex semantic relatedness and learn more discriminative representations to alleviate the data sparsity issue; 2) For the active users, HiCON yields consistent improvement over all baselines in terms of AUC and F1. This may be because the hierarchical message aggregation selects a bundle of valuable neighbors from vast neighbors and then derives semantic relatedness with them, which is helpful for making accurate recommendations.

A.6. Hyper-parameter Analysis

In this section, we analyze the impact of the hyper-parameter λ used for controlling the balance of BPR and contrastive losses. The results are shown in Figure 10. We observe that HiCON achieves relatively stable recommendation performance on both AUC and F1 at first, and then its performance declines significantly when λ is set to a large number (e.g. 10.0). This may be because the contrastive learning with a large weight λ dominates the overall loss, misleading the model optimization process.