

# Contrastive Multi-Level Graph Neural Networks for Session-based Recommendation

Fuyun Wang, Xingyu Gao\*, *Member, IEEE*, Zhenyu Chen, *Member, IEEE*, Lei Lyu

**Abstract**—Session-based recommendation (SBR) aims to predict the next item at a certain time point based on anonymous user behavior sequences. Existing methods typically model session representation based on simple item transition information. However, since session-based data consists of limited users' short-term interactions, modeling session representation by capturing fixed item transition information from a single dimension suffers from data sparsity. In this paper, we propose a novel contrastive multi-level graph neural networks (CM-GNN) to better exploit complex and high-order item transition information. Specifically, CM-GNN applies local-level graph convolutional network (L-GCN) and global-level network (G-GCN) on the current session and all the sessions respectively, to effectively capture pairwise relations over all the sessions by aggregation strategy. Meanwhile, CM-GNN applies hyper-level graph convolutional network (H-GCN) to capture high-order information among all the item transitions. CM-GNN further introduces an attention-based fusion module to learn pairwise relation-based session representation by fusing the item representations generated by L-GCN and G-GCN. CM-GNN averages the item representations obtained by H-GCN to obtain high-order relation-based session representation. Moreover, to convert the high-order item transition information into the pairwise relation-based session representation, CM-GNN maximizes the mutual information between the representations derived from the fusion module and the average pool layer by contrastive learning paradigm. We conduct extensive experiments on multiple widely used benchmark datasets to validate the efficacy of the proposed method. The encouraging results demonstrate that our proposed method outperforms the state-of-the-art SBR techniques.

**Index Terms**—Session-based Recommendation, Contrastive Learning, Graph Neural Networks

## I. INTRODUCTION

AS an important way to alleviate information overload, recommendation systems [1], [2], [3], [4], [5] play an important role in many e-commerce platforms, online entertainment platforms, and web search applications. Typically, on many websites, a user engages in brief interactions in a scenario where the system is not logged on. These interactions between the anonymous user and the system are organized into a session. The task of session-based recommendation (SBR) is to predict the user's next interaction based on the user's historical interaction session. Compared with conventional sequential recommendations methods, SBR methods model user preferences by capturing item transitions in the sessions

and have attracted widespread attention in recent years [6], [7], [8].

Given the characteristics of session-based data, many effective recommendation paradigms are proposed in early studies (e.g., item co-occurrence relations-based recommendation methods [9] and Markov chain-based recommendation methods [10], [11]). However, these approaches either ignore users' sequential behavior patterns or fail to capture complex item transitions in sessions. With the rapid growth of Recurrent Neural Networks (RNNs), many RNNs-based recommendation methods have achieved encouraging results in the fields of SBR [12], [13], [14], [6], [15], [16], [17]. Despite the promising progress, RNNs-based approaches overemphasize the sequential patterns of item transitions and thus suffer from overfitting problems. Besides, since RNNs-based approaches only model the relations between adjacent items and fail to capture implicit connections among all items, these methods still have shortcomings in their ability to mine the features of session-based data. Recently, self-attention-based approaches (e.g. SASRec [18], [19], [20], [21], [22], [23]) have been proposed to capture long-term dependencies in the sequence. But these kinds of approaches have learned too much item transition information which is not related to the current session and thus increases computational complexity.

In recent years, graph neural networks (GNNs) have been widely used in SBR [24], [25], [26], [27], [28]. SR-GNN [25] models complex transitions among items by adopting a gated graph neural network. But SR-GNN only calculates the relative importance of each item to the last item when learning session representation, which fails to capture the specific item transition pattern within the session. Compared with SR-GNN, FGNN [26] fully considers the inherent order of the item transition patterns and has achieved better performance. Although these methods alleviate the problem of previous work failing to capture implicit connections among all items to some extent, they suffer from two inherent algorithmic drawbacks. On the one hand, existing methods only use the current session to make recommendations, ignoring item transition patterns in other related sessions. On the other hand, existing methods model each session as a directed subgraph and regard item transitions as pairwise relations, failing to exploit the high-order information of item transitions.

To tackle the above issues, we propose a novel contrastive multi-level graph neural network (CM-GNN) to model the complex item transition patterns over all the sessions. Specifically, we design a local-level graph convolutional network (L-GCN) and a global-level graph convolutional network (G-GCN) to capture pairwise relations of the current session

F. Wang and L. Lyu are with School of Information Science and Engineering, Shandong Normal University, Jinan 250358, China.

\*Corresponding author: Xingyu Gao is with Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China.

Zhenyu Chen is with Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China.

and all the sessions, respectively. To be specific, G-GCN can capture item transition patterns related to the current session from other sessions, which is helpful to learn more accurate item representations. To learn pairwise relation-based session representation, an attention-based fusion module is designed to fuse the item representations outputted by L-GCN and G-GCN. Furthermore, CM-GNN designs a hyper-level graph convolutional network (H-GCN) for capturing high-order relations. By averaging the item representations generated by H-GCN, we obtain high-order relation-based session representation. Since the above two levels of session representation can be seen as two different views which describe the same session, CM-GNN converts the high-order item transition information into the pairwise relation-based session representation by maximizing the mutual information between the different session representations through contrastive learning. The main contributions of our work can be concluded as follows:

- We propose a novel contrastive multi-level graph neural network (CM-GNN) for session-based recommendation (SBR), which can model complex and high-order item transition patterns over all the sessions.
- We design local-level graph convolutional network (L-GCN), global-level graph convolutional network (G-GCN) and hyper-level graph convolutional network (H-GCN) for capturing pairwise relations and high-order relations.
- We obtain a pairwise relation-based session representation through an attention-based fusion module and a high-order relation-based session representation through an average pool layer, respectively, and incorporate high-order item transition information with pairwise item transition information through contrastive learning.
- We conduct comprehensive experiments on three benchmark datasets. The encouraging results show that our proposed CM-GNN outperforms the state-of-the-art methods.

The remainder of this paper is organized as follows. Section II reviews the related work. Section III describes the proposed CM-GNN in detail. Section IV presents extensive experiments. Finally, Section V concludes the paper.

## II. RELATED WORK

### A. Session-based Recommendation

The task of SBR is to predict users' next interaction according to their historical behavior sequences. Conventional SBR methods include item co-occurrence relation-based methods [9] and Markov chain-based methods [10] [11], [29]. Item co-occurrence relation-based methods are recommended based on the similarity between items and these methods fail to capture sequential patterns within sessions. Markov chain-based approaches use latest users' click to model users' short-term preferences, ignoring useful information used to model users' long-term preferences.

With the rapid development of neural networks in the field of natural language processing, deep learning-based methods have also achieved advanced performance in SBR. Hidasi et

al. [12] first apply Recurrent Neural Networks (RNNs) to sequential recommendation and propose a GRU4Rec model. Afterward, Li et al. [6] introduce attention mechanisms into GRU4Rec and capture more representative item transition information. Liu et al. [7] propose a short-term memory network named STAMP that uses attention mechanisms to replace RNNs encoding, which captures the users' general interests from the long-term memory of session context and takes into account the user's current interest from the short-term memory of the last click. SASRec [19] proposed by Wang et al. is another classical method in sequential recommendation, which models the users' historical behaviors through self-attention mechanism and fully captures the transition information of the item. Since Convolutional Neural Networks (CNNs) do not depend on the previous state and are sufficiently parallel in calculation, Wang et al. [30] utilize memory networks to predict users' intent based on collaboration sessions. The disadvantages of RNNs-based models and CNNs-based models are that they only consider the relations of adjacent items and fail to capture the complex relations among non-adjacent items [31], [32].

In recent years, Graph Neural Networks (GNNs) have achieved promising results in capturing complex relations between nodes [33], [34], [35], [24], [36]. In the field of SBR, Wu et al. [25] first propose the SR-GNN which used the graph neural network for sequence modeling. SR-GNN regards the problem of sequence modeling as a problem of graph modeling and learns the hidden vector representation of each item through a gating graph neural network. With the successful application of SR-GNN, some variants of GNNs have also achieved encouraging performance. GC-SAN [27] uses gating graph neural network (GGNN) to learn transition information between items and to learn long-term interests of users based on multiple self-attention layers. By studying inherent order of the item transition pattern in the sequence, FGNN [26] compensates for the SR-GNN's failure to take into account the transition pattern of specific items within the session, further enhancing the recommendation performance.

Although the above approaches have achieved promising performance in SBR, they still fail to obtain more useful information. First of all, the above methods construct each session into a directed subgraph that regards item transitions as pairwise relations. However, in the real-world scenario, there are often high-order relations among items, which can not be captured by the above methods. In addition, since the above methods are only for single-session graphs and do not take cross-session information into account, we often obtain the inaccurate session representations. Wang et al. [37] propose HyperRec which is the earliest work used by hypergraph convolutional network for sequential recommendation. HyperRec models the short-term preferences of users by constructing the sequence as a hypergraph [38], [39]. Afterward, Xia et al. [40] propose S<sup>2</sup>-DHCN which is a dual-channel hypergraph convolutional network for SBR. S<sup>2</sup>-DHCN models users' short-term interests based on hypergraph channels and constructs self-supervision signals based on line graph channels to enhance hypergraph modeling. In order to capture relevant item transition information from other sessions, CA-TCN [41]

and GCE-GNN [42] first learn the item representation based on the session graph and then calculate the local-level and global-level session representations respectively. Unlike CATCN, GCE-GNN constructs a global-level graph based on all the sessions and combines local-level session representations representing the user's short-term interests with global-level session representations representing the user's long-term interests to obtain the final session representation. GCE-GNN is the state-of-the-art method and has achieved the best performance in SBR. However, GCE-GNN only takes into account pairwise relations and fails to model the high-order relations. In addition, existing studies fail to combine pairwise relations with high-order relations, which is made up by our work.

### B. Contrastive Learning Recommendation

Recently, contrastive learning [43], [44], [45], [46], [47], [48] has been widely applied to mine its own supervision information from large-scale unsupervised data by using pretext tasks. It is first applied to the field of computer vision [49], [50], [51], [52], [53], and then make important progress in the field of audio processing [54], and natural language processing [55], [56], [57], [58]. An important branch of contrastive learning is mutual information maximization [46], [59], which maximizes the mutual information between these views by using views from the same input as positive samples and views from different inputs as negative samples.  $S^3$ -Rec [60] is a classical application of contrastive learning in sequential recommendation, which maximizes the mutual information between attribute, item, and sequence views by masking different levels of granularity of the contextual information. In the field of graph representation learning, contrastive learning has also made encouraging progress. SGL [61] aims to explore contrastive learning on the user-item graph. By changing the structure of the graph to generate multiple views of a node, SGL maximizes the agreement between different views of the same node and the views of other nodes, enriching the feature representation of the node. In the field of SBR,  $S^2$ -DHCN models session-based data as a hypergraph, and then constructs self-supervision signals based on the hypergraph-induced line graph, enhancing the modeling of the hypergraph by maximizing the mutual information between session representations learned by different graphs. Recently, contrastive learning has also been widely used in the field of social recommendation [62], and GroupIM [63] and MHCN [64] have gained more comprehensive user representation by constructing self-supervision signals from users' social relationships. Unlike the above approaches, since the session typically consists of brief interactions by anonymous users, this work focuses on exploring contrastive learning on item transitions in SBR to improve user intent learning.

## III. METHODOLOGY

In this section, we describe the algorithm details of contrastive multi-level graph neural network (CM-GNN). In subsection III-A, we describe the problem definitions for session-based recommendation (SBR) and introduce the notion

definitions used throughout the paper. In subsection III-B, we present an overview of CM-GNN. In subsection III-C, we describe the local-level graph convolutional network (L-GCN), global-level graph convolutional network (G-GCN), and hyper-level graph convolutional network (H-GCN). In subsection III-D, we describe how to learn the pairwise relation-based session representation (pairwise relation-based session representation is generated by items that contain pairwise item transition information) and the high-order relation-based session representation (high-order relation-based session representation is generated by items that contain high-order item transition information) based on the item representations generated by different GCNs. In subsection III-E, we describe the prediction process of our model. In subsection III-F, we show the optimization process of CM-GNN. In the subsection III-G, we describe how to convert the high-order item transition information into the pairwise relation-based session representation through contrastive learning paradigm.

### A. Problem Definition

Let  $V = \{v_1, v_2, \dots, v_n\}$  denote the set of all items, where  $n$  is the number of all unique items. Each anonymous session  $s$  can be represented by list  $s = \{v_{s,1}, v_{s,2}, v_{s,3}, \dots, v_{s,m}\}$ , where  $m$  is the length of  $s$  and  $v_{s,k} \in V (1 \leq k \leq m)$  represents an interacted item of an anonymous user within the session  $s$ . We describe each item  $v_i \in V$  with embedding vector  $\mathbf{x}_i^{(t)} \in \mathbb{R}^{d^{(t)}}$ , where  $d^{(t)}$  is the dimension of item in the  $t$ -th layer of a GCN. Given a specific session  $s$ , the goal of session-based recommendation (SBR) is to predict the next item  $v_{s,m+1}$ . Formally, our model aims to output the click probability of all items  $\hat{y} = \{\hat{y}_1, \hat{y}_2, \hat{y}_3, \dots, \hat{y}_n\}$ , where  $\hat{y}_i \in \hat{y}$  indicates the likelihood score of clicking item  $v_i$ . We rank the prediction scores in descending order, and recommend the candidate items ranked in the top-K. In this paper, we adopt bold capital letters to denote matrices, bold lowercase letters to denote vectors, and italic lowercase letters to represent scalars.

### B. Overview

We illustrate the framework of CM-GNN in Fig. 1. Firstly, we build all the sessions, the current session, and all the sessions as a global-level graph (G-G), local-level graph (L-G), and hyper-level graph (H-G) respectively. Then for each type of graph, we apply global-level graph convolutional network (G-GCN), local-level graph convolutional network (L-GCN), and hyper-level graph convolutional network (H-GCN) to capture global-level pairwise relations, local-level pairwise relations, and high-order relations respectively. Afterward, we use an attention-based fusion module to learn the pairwise relation-based session representation by fusing the global-level item representations and the local-level item representations. We also obtain the high-order relation-based session representation by averaging the hyper-level item representations. Moreover, we integrate contrastive learning into CM-GNN and incorporate high-order item transition information into the pairwise relation-based session representation by maximizing the mutual information across different levels of session representation. Finally, we obtain the next-click one for a given session by computing the scores of each candidate item.

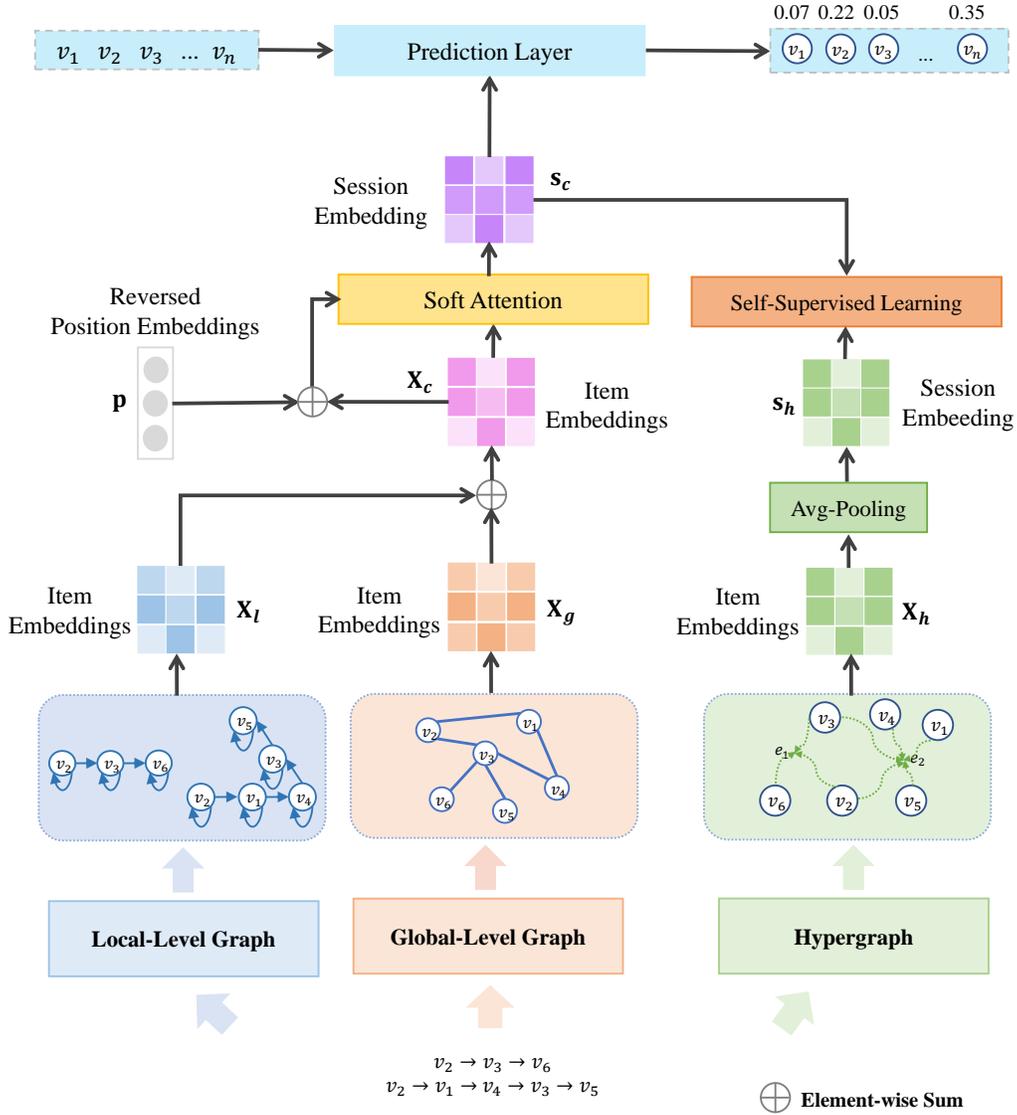


Fig. 1: The overall framework of our proposed CM-GNN.

### C. Different Levels of Graph Convolutional Network

In this section, we introduce the construction process of different levels of graph and describe the corresponding graph convolutional network on them.

1) *L-GCN*: To capture the pairwise relations among items in the current session, we construct a L-G for each session  $s$ . Let  $G_l = (\mathcal{V}_l, \mathcal{E}_l)$  denote L-G, where  $\mathcal{V}_l \subseteq V$  and  $\mathcal{E}_l = \{e_{ij}^l \mid (v_i, v_j) \mid v_i, v_j \in V\}$  represent the node set and the edge set, respectively. Each edge  $(v_i, v_j) \in \mathcal{E}_l$  implies that a user clicks item  $v_j$  after  $v_i$  in the session  $s$  and we capture pairwise item transition patterns of current session based on these edges. To get more item transition patterns, we add self-loop to each item. Moreover, since L-G is a directed graph, the types of directed edge are distinguished into four types based on the item transition order, i.e.,  $r_{in}$ ,  $r_{out}$ ,  $r_{in-out}$  and  $r_{self}$ . Given a specific edge  $(v_i, v_j)$ ,  $r_{in}$  indicates there is only one transition from  $v_i$  to  $v_j$  while  $r_{out}$  indicates there is only one transition

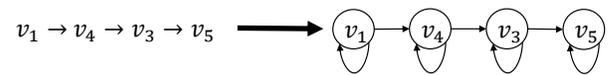


Fig. 2: The construction process of the local-level graph.

from  $v_j$  to  $v_i$ ,  $r_{in-out}$  indicates there are both transitions from  $v_i$  to  $v_j$  and from  $v_j$  to  $v_i$  and  $r_{self}$  indicates an item has a self-transition. The construction process of L-G is illustrated in Fig. 2.

Since not all neighbors of item  $v_i$  have the same importance to item representation learning, we introduce an attention mechanism to measure the contribution of neighbors by learning weights of edges between different nodes. Formally, the attention coefficients are calculated as follows:

$$a_{ij} = \text{LeakyRelu}(\mathbf{w}_{r_{ij}}^T (\mathbf{x}_{v_i} \odot \mathbf{x}_{v_j})), \quad (1)$$

where  $a_{ij}$  measures the importance of node  $v_j$  to node  $v_i$  and  $\text{LeakyRelu}(\cdot)$  is chosen as activation function.  $r_{ij}$  is the relation between node  $v_i$  and node  $v_j$ .  $\odot$  is the element-wise multiplication operation and  $\mathbf{w}_* \in \mathbb{R}^d$  are weight vectors.

There may be more than one neighbor for a same node and the contributions of these neighbors are different when learning node representation. As a consequence, if the representation of the current node is learned based on an equal contribution of its neighbors, it would result in inaccurate learning of node representation. To avoid this problem, we first normalize the attention coefficients through the softmax function:

$$\alpha_{ij} = \frac{\exp(\text{LeakyRelu}(\mathbf{w}_{r_{ij}}^T (\mathbf{x}_{v_i} \odot \mathbf{x}_{v_j})))}{\sum_{v_k \in \mathcal{N}_{v_i}^l} \exp(\text{LeakyRelu}(\mathbf{w}_{r_{ik}}^T (\mathbf{x}_{v_i} \odot \mathbf{x}_{v_k})))}, \quad (2)$$

where  $\alpha_{ij}$  denotes the attention coefficients after normalization,  $\mathcal{N}_{v_i}^l$  represents the collection of nodes which are adjacent to node  $v_i$  in L-G.

Then we obtain output features for each node by calculating the linear combination of the attention coefficient and its corresponding item representation.

$$\mathbf{x}_{v_i}^l = \sum_{v_j \in \mathcal{N}_{v_i}^s} \alpha_{ij} \mathbf{x}_{v_j}. \quad (3)$$

We use an attention mechanism to aggregate the features of the item itself and its neighbors in the current session and obtain the local-level item representation.

We denote the pairwise local-level item representations as  $\mathbf{X}_l$ , where each row in  $\mathbf{X}_l$  represents the local-level item representation after L-GCN learning, i.e.  $\mathbf{x}_{v_i}^l$ .

2) *G-GCN*: Since a session is generally a short sequence, modeling the current session only based on L-G often suffers from the problem of data sparsity and results in difficulty generating a comprehensive and accurate item representation. To overcome the above problems, we construct G-G over all the sessions to capture the item transition patterns within all sessions. Based on the fact that the same item appears in both the current session and the other sessions, we follow the concept of  $\varepsilon$ -neighbor proposed by GCE-GNN [42] (i.e.,  $\varepsilon$ -neighbor set) to construct G-G. Let  $G_g = (\mathcal{V}_g, \mathcal{E}_g)$  denote G-G, where  $\mathcal{V}_g$  is the node set which contains all the items in  $V$ ,  $\mathcal{E}_g = \{e_{ij}^g \mid (v_i, v_j) \mid v_i \in V, v_j \in \mathcal{N}_{\varepsilon(v_i)}\}$  indicates the edge set, where  $\mathcal{N}_{\varepsilon(v_i)}$  represents the  $\varepsilon$ -neighbor set of item  $v_i$ . Each edge  $e_{ij}$  corresponds to a pairwise relation between item  $v_i$  and item  $v_j$  existing in the  $\varepsilon$ -neighbor of item  $v_i$  within all sessions. Different from L-G, to distinguish the importance of different neighbors to item  $v_i$ , each edge is assigned a weight based on the occurrence frequency of corresponding item transition patterns within all sessions. G-G is an undirected graph and the construction process is illustrated in Fig. 3.

Based on G-G, we develop G-GCN to capture the pairwise relations derived from all the sessions. Specifically, we first identify the multiple sessions where item  $v_i$  is involved and then capture the relevant item transition patterns. Since not

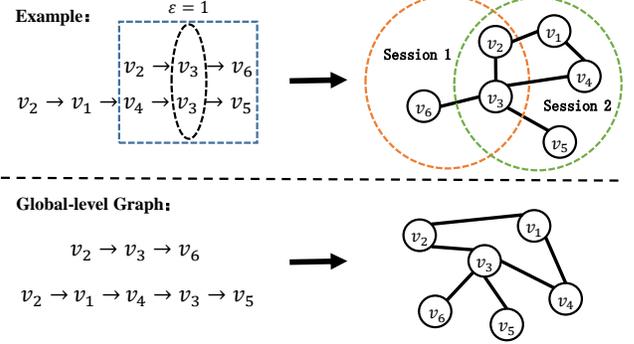


Fig. 3: The construction process of the global-level graph.

all of the items in  $\mathcal{N}_{\varepsilon(v_i)}$  are relevant to the user preference of the current session, we introduce a session-aware attention mechanism to learn the importance of each neighbor to  $v_i$ .

First, to obtain relevance between items in  $\mathcal{N}_{\varepsilon(v_i)}$  and the current session, we calculate representation of the current session by averaging the item representations:

$$\mathbf{s} = \frac{1}{m} \sum_{v_i \in s} \mathbf{x}_{v_i}, \quad (4)$$

where  $s$  denotes the current session and  $m$  is the length of  $s$ .

Then, we calculate importance weight of different neighbors to the current session:

$$b_{ij} = \mathbf{q}_1 \text{LeakyRelu}(\mathbf{W}_1[(\mathbf{s} \odot \mathbf{x}_{v_j}) \parallel w_{ij}]), \quad (5)$$

where  $b_{ij}$  denotes the weight coefficient of different neighbors of node  $v_i$  in G-G to the current session,  $w_{ij}$  denotes the weight of edge  $(v_i, v_j)$  in G-G,  $\mathbf{q}_1$  and  $\mathbf{W}_1$  are two trainable parameters.  $\odot$  and  $\parallel$  represent the element-wise product and concatenation operation, respectively.

To make the importance weight comparable across all neighbors connected with  $v_i$ , we normalize the importance weight of different neighbors by softmax function:

$$b_{ij} = \frac{\exp(b_{ij})}{\sum_{v_k \in \mathcal{N}_{v_i}^g} \exp(b_{ik})}, \quad (6)$$

Obviously, the neighbor nodes with higher importance weight should be given more attention.

Then, we calculate the neighbor representation of node  $v_i$  as follows:

$$\mathbf{x}_{\mathcal{N}_{v_i}^g} = \sum_{v_j \in \mathcal{N}_{v_i}^g} b_{ij} \mathbf{x}_{v_j}, \quad (7)$$

Finally, the global-level item representations are obtained by aggregating their own representations and their neighbors' representations:

$$\mathbf{x}_v^g = \text{Relu}(\mathbf{W}_2[\mathbf{x}_v \parallel \mathbf{x}_{\mathcal{N}_{v_i}^g}]). \quad (8)$$

where  $\mathbf{W}_2$  denotes the trainable parameter and here  $\text{Relu}(\cdot)$  is chosen as the activation function.

We denote the pairwise global-level item representations as  $\mathbf{X}_g$ , where each row in  $\mathbf{X}_g$  represents the local-level item representation after G-GCN learning, i.e.  $\mathbf{x}_{v_i}^g$ .

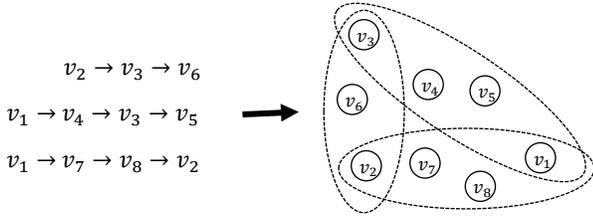


Fig. 4: The construction process of the hyper-level graph.

3) *H-GCN*: In this section, we first present the definition of H-G. Then we describe the construction process of H-G in detail. Finally, H-GCN is proposed to learn the high-order relations over all the sessions.

Let  $G_h = (\mathcal{V}_h, \mathcal{E}_h)$  denote H-G, where  $\mathcal{V}_h$  and  $\mathcal{E}_h$  denote the set of nodes and the set of hyperedges, respectively. Each hyperedge  $e_{ij}^h \in \mathcal{E}_h$  contains all the items in each session, and all hyperedges are connected by the shared items in the sessions. The construction process of the H-G is shown in Fig. 4.

After the H-G is constructed, we develop the H-GCN to capture high-order item transition patterns. Inspired by S<sup>2</sup>-DHCN [40], we regard the process of H-GCN as a two-stage refinement of the 'node-hyperedge-node' feature transformation. We follow S<sup>2</sup>-DHCN to assign each hyperedge the same weight with the value of 1 and define H-GCN with row normalization as:

$$\mathbf{X}_h^{t+1} = \mathbf{D}^{-1} \mathbf{H} \mathbf{W} \mathbf{B}^{-1} \mathbf{H}^T \mathbf{X}_h^{(t)}, \quad (9)$$

where  $\mathbf{H} \in \mathbb{R}^{(N \times M)}$  indicates the matrix representation of H-G,  $\mathbf{W}$  is a diagonal matrix which consists of the weight of each hyperedge,  $\mathbf{D}$  and  $\mathbf{B}$  are two diagonal matrices which consist of the degree of each vertex and the degree of each hyperedge, respectively.

To get the high-order hyper-level item representations, we extend H-GCN from one layer to multiple layers and average the items representations generated from each layer. Formally, we define this process as follows:

$$\mathbf{X}_h = \frac{1}{L+1} \sum_{t=0}^L \mathbf{X}_h^{(t)}. \quad (10)$$

where  $L$  denotes the number of the convolutional layer and  $t$  denotes the current H-GCN layer.

We denote each row in  $\mathbf{X}_h$  as the hyper-level item representation after H-GCN learning, i.e.,  $\mathbf{x}_{v_i}^h$ .

#### D. Session Representation Learning Layer

Aided by the above three GCNs, we can obtain three different levels of item representations. On this basis, we describe how to learn corresponding session representations. First, we learn the pairwise relation-based item representation by aggregating the representation of items obtained by L-GCN and G-GCN. Specifically, we calculate the above item representation by sum pooling:

$$\mathbf{x}'_{v_i} = \mathbf{x}_{v_i}^g + \mathbf{x}_{v_i}^l, \quad (11)$$

We denote the pairwise relation-based item representations as  $\mathbf{X}_c$ , where each row in  $\mathbf{X}_c$  represents item representation after the above sum pooling operation, i.e.  $\mathbf{x}'_{v_i}$ .

After that, to obtain a more comprehensive session representation, we introduce an attention-based fusion module to adaptively calculate the contribution of different items to the next prediction. First, to incorporate the sequential information into the session representation, we add a learnable position embeddings  $\mathbf{p} \in \mathbb{R}^d$  into the item representations:

$$\mathbf{z}_{v_i} = \text{Tanh}(\mathbf{W}_3 [\mathbf{x}'_{v_i} \parallel \mathbf{p}_{t-i+1}] + \mathbf{b}_3), \quad (12)$$

where  $\mathbf{W}_3$  and  $\mathbf{b}_3$  are trainable parameters and here  $\text{Tanh}(\cdot)$  is chosen as the activation function.

We get the static representation  $\mathbf{s}'$  of the session by averaging item representations in it:

$$\mathbf{s}' = \frac{1}{t} \sum_{i=1}^t \mathbf{x}'_{v_i}, \quad (13)$$

Following SR-GNN [25], we calculate the attention coefficients as follows:

$$\beta_i = \mathbf{q}_2^T \sigma(\mathbf{W}_4 \mathbf{z}_{v_i} + \mathbf{W}_5 \mathbf{s}' + \mathbf{b}_4), \quad (14)$$

where  $\sigma(\cdot)$  is the sigmoid function, and  $\mathbf{W}_4, \mathbf{W}_5, \mathbf{q}_2, \mathbf{b}_4$  are trainable parameters.

Finally, we obtain the pairwise relation-based session representation  $\mathbf{s}_c$  by calculating the linear combination of the attention coefficient and its corresponding item representation:

$$\mathbf{s}_c = \sum_{i=1}^t \beta_i \mathbf{x}'_{v_i}. \quad (15)$$

After obtaining the pairwise relation-based session representation, we obtain the high-order relation-based session representation  $\mathbf{s}_h$  by averaging the hyper-level item representations:

$$\mathbf{s}_h = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_{v_i}^h. \quad (16)$$

#### E. Prediction Layer

In this section, we describe the prediction layer to compute the score of each candidate item based on the learned session representation. Given a session  $s$ , we calculate the scores  $\hat{\mathbf{y}}_{v_i}$  for all the candidate items by multiplying the pairwise relation-based session representation with all item representations:

$$\hat{\mathbf{y}}_{v_i} = \text{Softmax}(\mathbf{s}_c^T \mathbf{x}_{v_i}). \quad (17)$$

here  $\text{Softmax}(\cdot)$  is chosen as our activation function.

#### F. Model Optimization

In this section, we present the optimization process of our model. Specifically, we define the learning objective as the cross entropy loss function, which has been extensively used in recommendation systems:

$$\mathcal{L}_r = - \sum_{i=1}^n \mathbf{y}_{v_i} \log(\hat{\mathbf{y}}_{v_i}) + (1 - \mathbf{y}_{v_i}) \log(1 - \hat{\mathbf{y}}_{v_i}). \quad (18)$$

where  $\mathbf{y}$  denotes the one-hot encoding vector of the ground truth item.

### G. Enhancing CM-GNN with Contrastive Learning

In this section, we show how to mine self-supervision signals based on contrastive learning to enhance the training process of the model. In the above subsections, we capture the pairwise item transition patterns over all the sessions through L-GCN and G-GCN, and obtain the pairwise relation-based session representation based on an attention-based fusion module. We also use H-GCN to capture the high-order item transition patterns over all the sessions and obtain the high-order relation-based session representation by averaging the hyper-level item representations. Both the above two levels of session representation only capture a specific item transition patterns in the session, which suffers more from the problem of data sparsity and result in suboptimal recommendation performance. To get a more accurate and comprehensive session representation, we convert the high-order item transition information into the pairwise relation-based session representation. Specifically, we regard the two levels of session representation as two views characterizing different aspects of the same session and innovatively integrate contrastive learning into the training of CM-GNN. CM-GNN improves the recommendation performance by maximizing the mutual information between two different levels of session representation.

Technically, we adopt InfoNCE [45], which can maximize the mutual information between the different levels of session representation, as our learning objective:

$$\mathcal{L}_s = -\log\sigma(f_D(\mathbf{s}_{v_i}^c, \mathbf{s}_{v_i}^h)) - \log\sigma(1 - f_D(\tilde{\mathbf{s}}_{v_i}^c, \mathbf{s}_{v_i}^h)), \quad (19)$$

where  $f_D(\cdot) : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  is the discriminator function that takes two vectors as the input and then scores the agreement between them. We simply implement the discriminator as the dot product between two representations. Since the pairwise relation-based session representation and the high-order relation-based session representation both model the same session, they can be the ground truth of each other. We corrupt  $\mathbf{s}_{v_i}^c$  by both row-wise and column-wise shuffling to create negative examples  $\tilde{\mathbf{s}}_{v_i}^c$ .

Finally, we unify the recommendation task and the contrastive task and optimize it by means of joint learning. Formally, we define the joint learning objective as follows:

$$\mathcal{L} = \mathcal{L}_r + \beta\mathcal{L}_s. \quad (20)$$

where  $\beta$  is a hyper-parameter used to control the magnitude of the contrastive task.

## IV. EXPERIMENTS

In this section, we conduct extensive experiments to validate our model and answer the following four key research questions:

**RQ 1** Does CM-GNN achieve state-of-the-art performance?

**RQ 2** Does each component in CM-GNN contribute?

**RQ 3** How does the number of layers in different GCNs influence the performance of CM-GNN?

**RQ 4** What is the effectiveness of contrastive learning in CM-GNN?

**RQ 5** How does hyper-parameter  $\beta$  in contrastive learning influence the performance of CM-GNN?

TABLE I: Dataset Statistics

| Dataset           | Tmall  | Diginetica | Nowplaying |
|-------------------|--------|------------|------------|
| training sessions | 351268 | 719470     | 825304     |
| test sessions     | 25898  | 60858      | 89824      |
| # of items        | 30968  | 40728      | 43097      |
| average lengths   | 6.69   | 5.12       | 7.42       |

### A. Datasets

We evaluate our model on three real-world benchmark datasets, i.e., Diginetica, Tmall, and Nowplaying. The Diginetica<sup>1</sup> dataset is obtained from the CIKM Cup 2016, which consists of typical transaction data. Tmall<sup>2</sup> dataset comes from the IJCAI-15 competition, which consists of anonymized users' shopping logs on the Tmall online shopping platform. Nowplaying<sup>3</sup> is obtained from [65], which consists of the music listening behavior of users. For a fair comparison, we follow the experiment environment in [40] to preprocess the three benchmark datasets. Specifically, we filter out all the sessions containing only one item and items appearing less than five times for all the datasets. We regard the latest data (such as the sessions of last week) as test data and the remaining data as the training set. Moreover, we augment and label both the training dataset and the test dataset by employing a sequence splitting approach for all the datasets, then we generate multiple labeled sequences with the corresponding labels, i.e.,

$$([v_{s,1}, v_{s,2}], ([v_{s,1}, v_{s,2}], v_{s,3}), \dots, ([v_{s,1}, v_{s,2}, v_{s,m-1}], v_{s,m}).$$

We summarize the statistics of datasets after preprocessing in Table I.

### B. Baseline Methods

We compare CM-GNN with eleven strong and commonly used methods, which are presented as follows:

- **POP**: POP recommends top-K items based on their popularity.
- **Item-KNN**[9]: Item-KNN computes the cosine similarity between items of the current session and other sessions and recommends the most similar top-K items.
- **FPMC**[10]: FPMC captures users' short-term preferences by combining the matrix factorization and the first-order Markov chain.
- **GRU4REC**[12]: GRU4REC employs Gated Recurrent Unit (GRU) to model user behavior sequences.
- **NARM**[6]: NARM is an RNN-based method, which adopts the attention mechanism to model sequential patterns.
- **STAMP**[7]: STAMP utilizes attention layers to replace all RNN encoders in the previous work and resorts to the last item as the user's short-term interest in the current session to make the prediction.
- **CSR**M[30]: CSR utilizes a parallel memory module to capture the sequential patterns of the latest m sessions.

<sup>1</sup><https://competitions.codalab.org/competitions/11161>

<sup>2</sup><https://tianchi.aliyun.com/dataset/dataDetail?dataId=42>

<sup>3</sup><http://dbis-nowplaying.uibk.ac.at/nowplaying>

- **SR-GNN**[25]: SR-GNN applies a Gated Graph Neural Network (GGNN) to learn item embeddings and obtains the session embeddings by applying a soft-attention mechanism.
- **FGNN**[26]: FGNN extends NARM by considering the inherent order of the item transition patterns.
- **GCE-GNN**[42]: GCE-GNN applies different levels of graph convolutional networks (GCNs) on both session-level graph and global-level graph respectively and captures comprehensive pairwise item transition patterns.
- **S<sup>2</sup>-DHCN**[40]: S<sup>2</sup>-DHCN constructs a hypergraph to capture the high-order relations and constructs a linegraph for contrastive learning to enhance network training.

### C. Evaluation Metrics

Following the previous work [25], [40], [42], we adopt two relevancy-based metrics P@K(Precision) and MRR@K(Mean Reciprocal Rank) to evaluate our approach. In the top-K items, P@K indicates the proportion of correctly recommended items. MRR@K represents the average of reciprocal ranks of the correctly recommended items in the top-k ranking list. In the paper, we adopt P@K and MRR@K with K = 10, 20 to evaluate all compared methods.

### D. Parameter Setup

For CM-GNN, we split the random 10% subset of the training set into the validation set and keep the hyper-parameters of each model consistent for a fair comparison. We set the dimension of latent vectors to 100 for all the datasets. Besides, we set mini-batch to 256 and we use Gaussian distribution  $\mathcal{N}(0, 0.1^2)$  to initialize all parameters in our model. We use the Adam optimizer to optimize these parameters, where the initial learning rate is set to 0.001 and decay by 0.1 after every 3 epochs. The L2 penalty is set to  $10^{-5}$  and the number of neighbors and the maximum distance of adjacent items  $\varepsilon$  are set to 12 and 3, respectively. Here, we follow their best parameter setups reported in the original papers and present their best results for all the baseline methods.

### E. Recommendation Performance (RQ 1)

To evaluate the performance of CM-GNN, we report the comparison results on three widely used benchmark datasets in Table 2, where we highlight the best result of each column in boldface. As can be observed, our proposed CM-GNN outperforms all the baseline methods on all datasets.

Traditional methods include POP, FPMC, and Item-KNN. POP only recommends top-N frequent items according to popularity, showing the worst performance. FPMC utilizes first-order Markov chain and matrix factorization to recommend the next item by taking into account the simple contextual information of the session, which achieves better performance than POP. Among the traditional methods, Item-KNN achieves the best results, which recommends the top-K items by computing the cosine similarity between items of the current session and other sessions. However, Item-KNN fails to capture the sequential patterns within the session and still does not achieve the desired performance.

Compared with the conventional methods, neural network-based methods consider the sequential dependency in the sessions and achieve decent performance for SBR. Methods based on deep recurrent neural structure are the earliest deep learning methods for modeling user preferences. Specifically, GRU4REC is the first work to adapt Recurrent Neural Networks (RNNs), which employs Gated Recurrent Unit (GRU) to model user behavior sequences. Compared with GRU4Rec, NARM and STAMP introduce the attention mechanism into RNN and improve the performance by a large margin. The superior performance of NARM and STAMP indicates that the attention mechanism has played an important role in sequence modeling. However, NARM and STAMP are modeled only based on the explicit interactions of user behavior, which neglects the implicit connection of the items in sessions. By incorporating collaborative information from others into the current session, CSRM enriches the session representation and achieves decent performance. However, CSRM often encodes irrelevant information from other sessions into the current session and suffers from noise impact. To overcome the above issues, subsequent studies begin to model session-based data as graphs and employ Graph Neural Networks (GNNs) to capture implicit connections among items in sessions. SR-GNN is the first work to employ GNN for sequential patterns capturing in sessions, which captures pairwise relations among items based on the Graph Gating Neural Network (GGNN). Compared with SR-GNN, FGNN takes into account the specific item transition patterns within the session for better performance. However, SR-GNN and FGNN only capture the sequential patterns in the current session and do not take advantage of item transition information in other sessions. GCE-GNN captures pairwise item transition information in all the sessions by constructing session-based data into a session-level graph and global-level graph, respectively. However, GCE-GNN only captures pairwise relations among items in sessions which neglects the complex and high-order item transition information. In recent work, S<sup>2</sup>-DHCN has achieved decent recommended performance by modeling high-order relations among items in sessions. However, S<sup>2</sup>-DHCN fails to explicitly capture pairwise item transition information among items. Besides, since sessions are generally short sequences formed by the interactions of anonymous users, modeling session-based data directly as a hypergraph can suffer from data sparsity problem.

Compared with the above baselines, our proposed CM-GNN shows overwhelming superiority on all evaluation metrics over all the datasets. On average, CM-GNN outperforms the best result by 3.7% on Tmall, 0.8% on Deginetica, 3.3% on Nowplaying. CM-GNN shows significant recommendation advantages over most methods, which are benefit from the consideration of global item transition information. By capturing the item transition patterns over all the sessions, CM-GNN generates more comprehensive and accurate session representations. Compared with GCE-GNN, CM-GNN still achieves decent performance improvements on the Tmall and Nowplaying datasets, which attributes to the ability to model high-order relations based on hyper-level graph. In contrast, CM-GNN's performance improvement on the Diginetica dataset was not significant, one possible explanation is that the average session

TABLE II: The overall performances based on three datasets. The best is bold and the second is underline.

| Datasets<br>Methods  | Tmall        |              |              |              | Diginetica   |              |              |              | Nowplaying   |              |             |             |
|----------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|-------------|
|                      | P@10         | P@20         | MRR@10       | MRR@20       | P@10         | P@20         | MRR@10       | MRR@20       | P@10         | P@20         | MRR@10      | MRR@20      |
| POP                  | 1.67         | 2.00         | 0.88         | 0.90         | 0.76         | 1.18         | 0.26         | 0.28         | 1.86         | 2.28         | 0.83        | 0.86        |
| Item-KNN             | 6.65         | 9.15         | 3.11         | 3.31         | 25.07        | 35.75        | 10.77        | 11.57        | 10.96        | 15.94        | 4.55        | 4.91        |
| FPMC                 | 13.10        | 16.06        | 7.12         | 7.32         | 15.43        | 22.14        | 6.20         | 6.66         | 5.28         | 7.36         | 2.68        | 2.82        |
| GRU4Rec              | 9.47         | 10.93        | 5.78         | 5.89         | 17.93        | 30.79        | 7.73         | 8.22         | 6.74         | 7.92         | 4.40        | 4.48        |
| NARM                 | 19.17        | 23.30        | 10.42        | 10.70        | 35.44        | 48.32        | 15.13        | 16.00        | 13.6         | 18.59        | 6.62        | 6.93        |
| STAMP                | 22.63        | 26.47        | 13.12        | 13.36        | 33.98        | 46.62        | 14.26        | 15.13        | 13.22        | 17.66        | 6.57        | 6.88        |
| CSRM                 | 24.54        | 29.46        | 13.62        | 13.96        | 36.59        | 50.55        | 15.41        | 16.38        | 13.20        | 18.14        | 6.08        | 6.42        |
| SR-GNN               | 23.41        | 27.57        | 13.45        | 13.72        | 38.42        | 51.26        | 16.89        | 17.78        | 14.17        | 18.87        | 7.15        | 7.47        |
| FGNN                 | 20.67        | 25.24        | 10.07        | 10.39        | 37.72        | 50.58        | 15.95        | 16.84        | 13.89        | 18.78        | 6.8         | 7.15        |
| GCE-GNN              | <u>28.01</u> | <u>33.42</u> | <u>15.08</u> | <u>15.42</u> | <u>41.16</u> | <u>54.22</u> | <u>18.15</u> | <u>19.04</u> | 16.94        | 22.37        | <u>8.03</u> | <u>8.40</u> |
| S <sup>2</sup> -DHCN | 26.22        | 31.42        | 14.60        | 15.05        | 40.21        | 53.66        | 17.59        | 18.51        | <u>17.35</u> | <b>23.50</b> | 7.87        | 8.18        |
| <b>CM-GNN</b>        | <b>29.72</b> | <b>34.69</b> | <b>16.06</b> | <b>16.41</b> | <b>41.56</b> | <b>54.67</b> | <b>18.28</b> | <b>19.19</b> | <b>17.50</b> | <u>23.30</u> | <b>8.05</b> | <b>8.44</b> |

TABLE III: The ablation study.

| Dataset<br>Measures | Tmall        |              |              |              | Diginetica   |              |              |              | Nowplaying   |             |              |             |
|---------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------------|--------------|-------------|
|                     | P@10         | MRR@10       | P@20         | MRR@20       | P@10         | MRR@10       | P@20         | MRR@20       | P@10         | MRR@10      | P@20         | MRR@20      |
| CG-L                | 12.10        | 6.64         | 14.56        | 6.81         | 35.50        | 14.42        | 48.71        | 15.34        | 15.42        | 7.47        | 20.69        | 7.83        |
| CG-G                | 25.94        | 13.86        | 30.41        | 14.17        | 41.53        | 18.14        | 54.53        | 19.05        | 17.49        | 7.36        | 22.89        | 7.76        |
| CG-H                | 28.85        | 15.40        | 33.85        | 15.76        | 41.55        | 18.21        | 54.53        | 19.11        | 17.35        | 8.12        | 22.79        | 8.49        |
| CG-OL               | 28.55        | 14.78        | 32.74        | 15.02        | 37.25        | 15.77        | 51.34        | 16.71        | 16.47        | 7.09        | 22.10        | 7.35        |
| CG-OG               | 18.29        | 9.82         | 20.54        | 10.71        | 30.80        | 12.02        | 45.58        | 13.42        | 14.47        | 6.46        | 18.39        | 6.64        |
| CG-OH               | 26.03        | 12.44        | 31.27        | 13.93        | 37.88        | 16.01        | 52.18        | 16.84        | 15.63        | 6.72        | 21.34        | 7.02        |
| CG-A                | 24.49        | 13.94        | 28.98        | 14.25        | 36.34        | 15.90        | 48.48        | 16.73        | 13.32        | 5.71        | 18.83        | 6.09        |
| CG-S                | 28.54        | 15.14        | 33.70        | 15.50        | 41.39        | 18.20        | 54.55        | 19.11        | 17.36        | <b>8.13</b> | 22.93        | <b>8.52</b> |
| <b>CM-GNN</b>       | <b>29.72</b> | <b>16.07</b> | <b>34.69</b> | <b>16.41</b> | <b>41.56</b> | <b>18.28</b> | <b>54.69</b> | <b>19.19</b> | <b>17.83</b> | 8.05        | <b>23.30</b> | 8.44        |

length of Diginetica is relatively smaller than the other two datasets, which makes the modeling of high-order relations does not capture more effective session features.

#### F. Ablation Study

The overwhelming superiority of CM-GNN presented in the last section can be seen as the result of joint effect of different components. To investigate the contributions of each component in CM-GNN, we conduct extensive experiments on variants of CM-GNN for all the datasets.

1) *Impact of each component in CM-GNN (RQ2)*: In this subsection, we conduct extensive experiments on three datasets to evaluate the effectiveness of local-level graph convolutional network (L-GCN) module, global-level graph convolutional network (G-GCN) module, hyper-level graph convolutional network (H-GCN) module, attention-based fusion module and contrastive learning module. Technically, we develop eight variants of CM-GNN: CG-L, CG-G, CG-H, CG-OL, CG-OG, CG-OH, CG-A, and CG-S. Specifically, CG-L represents the version without L-GCN module, CG-G represents the version without G-GCN module, CG-H represents the version without H-GCN module, CG-OL represents the version that we remove G-GCN module and H-GCN module and only remain L-GCN module, CG-OG represents the version that we remove L-GCN module and H-GCN module and only remain G-GCN module, CG-OH represents the version that we remove L-GCN module and G-GCN module and only remain H-GCN module, CG-A represents the version that we remove the attention-based fusion module and replace it with averaging item representations as the representation of each session and CG-S represents the version that we replace the InfoNCE loss

function with MSE loss function. We report the results of CM-GNN and these five variants above in Table III.

As can be observed from Table III, CM-GNN's performance on both the Tmall and Diginetica datasets consistently outperforms the five variants, which demonstrates that each component contributes to the CM-GNN recommendation process. Overall, L-GCN module is one of the most critical parts of CM-GNN, which serves for the capture of the item transition patterns within the current session. When removing L-GCN module, we can observe a remarkable performance drop on all the four metrics on both the Tmall dataset and Diginetica dataset. Correspondingly, it can be observed that all the metrics on the Nowplaying dataset have a significant drop when the attention-based fusion module is replaced with the average pooling module, which demonstrates that the attention-based fusion module avoids bringing up a bias of unrelated items when propagating information from all items within all the sessions. For the Tmall and Diginetica datasets, G-GCN module and the attention-based fusion module are the modules that play the second key role in CM-GNN, while for the Nowplaying dataset, the second most important modules of CM-GNN are L-GCN module and G-GCN module. We can conclude that the capture of global-level item transition patterns helps the model make more accurate predictions and learning different item importance across sessions is better than directly averaging representations of contained items for learning session representations in SBR. Moreover, Table III shows that all the metrics on the Tmall dataset and the Diginetica dataset have a significant drop when either the H-GCN module is removed or the InfoNCE loss function is replaced with MSE loss, which indicates that the application of H-GCN module and contrastive learning help the model learn more

TABLE IV: The comparison performances of different layers in GCNs module.

| Datasets       | Tmall        |              | Diginetica   |              | Nowplaying   |             |
|----------------|--------------|--------------|--------------|--------------|--------------|-------------|
|                | P@20         | MRR@20       | P@20         | MRR@20       | P@20         | MRR@20      |
| G-1-H-1        | 28.92        | 14.08        | 54.30        | 19.14        | 23.17        | 8.30        |
| <b>G-2-H-1</b> | <b>34.56</b> | <b>16.41</b> | <b>54.66</b> | <b>19.19</b> | <b>23.24</b> | 8.39        |
| G-2-H-2        | 34.53        | 16.33        | 54.61        | 19.17        | 23.12        | 8.41        |
| G-2-H-3        | 34.02        | 16.09        | 54.56        | 19.15        | 23.13        | <b>8.44</b> |

TABLE V: The effectiveness of contrastive learning in CM-GNN.

| Datasets      | Tmall        |              | Diginetica   |              | Nowplaying   |             |
|---------------|--------------|--------------|--------------|--------------|--------------|-------------|
|               | P@20         | MRR@20       | P@20         | MRR@20       | P@20         | MRR@20      |
| CG-POOL       | 31.21        | 14.47        | 50.72        | 16.22        | 21.71        | 7.87        |
| <b>CM-GNN</b> | <b>34.69</b> | <b>16.41</b> | <b>54.67</b> | <b>19.19</b> | <b>23.30</b> | <b>8.44</b> |

high-order information. Relatively speaking, although H-GCN module and the contrastive learning module also contribute to the P@K metric on the Nowplaying dataset, they do not exhibit competitive performance on the MRR@K metric. One possible explanation is that due to the relatively small average session length of Nowplaying, the pairwise item transition patterns and the high-order item transition patterns captured by CM-GNN may be very similar, which results in the model collapse when training based on contrastive learning.

2) *Impact of Model Depth (RQ3)*: We evaluate the influence of the number of layers in different GCNs on CM-GNN recommendation performance. Specifically, we perform all the experiments on our RTX 2080 Ti GPU, and to avoid suffering from the out-of-memory problem, we fix the number of layers of L-GCN to 1. Moreover, we range the number of layers of G-GCN in {1,2} and we range the number of layers of H-GCN in {1,2,3}. Technically, We develop four variant versions of CM-GNN: G-1-H-1, G-2-H-1, G-2-H-2, and G-2-H-3 to explore the influence of the number of layers of different levels in GCNs on performance. In G-1-H-1, we set the number of layers of G-GCN and H-GCN to 1 and 1, respectively. In G-2-H-1, we set the number of layers of G-GCN and H-GCN to 2 and 1, respectively. In G-2-H-2, we set the number of layers of G-GCN and H-GCN to 2 and 2, respectively. In G-2-H-3, we set the number of layers of G-GCN and H-GCN to 2 and 3, respectively. We show the results of four variants of CM-GNN in Table IV.

From the table IV, CM-GNN with G-GCN building by two layers performs better one layer, which indicates that high-level exploring might capture more effective information from global-level graph. Furthermore, we fix the number of layers of the G-GCN to 2 and range the number of layers of H-GCN in {1, 2, 3} for comparison. It can be observed that H-GCN with 1 layer performs the best on both the Tmall dataset and Diginetica dataset. As the number of layers of the H-GCN increases, the overall performance tends to drop, which indicates that higher-level exploring might also introduce noise. However, The H-GCN with more layers achieves better performance on the Nowplaying dataset, one possible explanation is that due to the relatively small average session length of Nowplaying, multi-layer H-GCN can capture more item transition patterns.

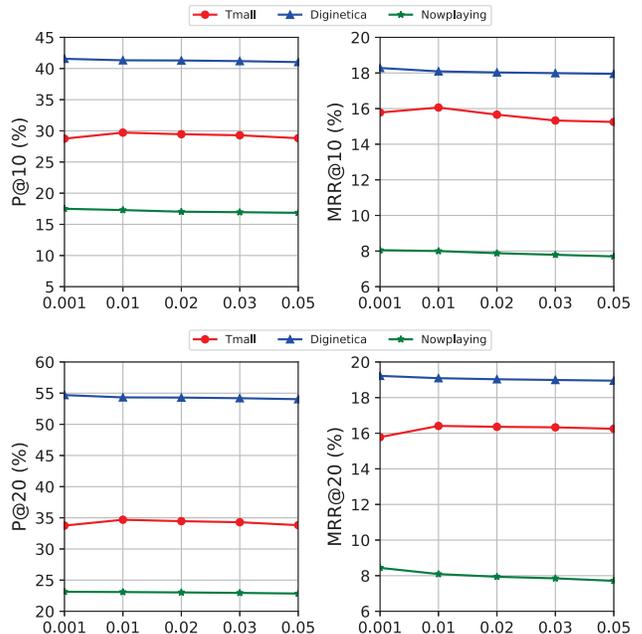


Fig. 5: The influence of the magnitude of contrastive learning.

3) *Effectiveness of Contrastive Learning (RQ4)*: Existing methods typically utilize the item prediction loss to learn the entire model. Since the length of session-based data is short, learning parameters of the model with a single optimization objective usually suffers from data sparsity. To this end, we borrow the idea of contrastive learning for alleviating the data sparsity problem and improving recommendation performance. To study the effectiveness of using contrastive learning, we design a variant version of CM-GNN named CG-POOL for comparison on all datasets with the complete CM-GNN model. Specifically, CG-POOL removes contrastive learning in CM-GNN and directly fuses local-level item representations, global-level item representations, and hyper-level item representations for learning final session representation by the attention-based fusion module.

As can be observed from table V, CM-GNN’s performance consistently outperforms the variant version CG-POOL across different datasets, which is credited to the capability of contrastive learning on learning the extra self-supervision signals and alleviating the data sparsity problem. Compared with CG-POOL, CM-GNN use contrastive learning to maximize the mutual information between different levels of session representation. And in the above process, different levels of item transition information (e.g. pairwise and high-order item transition information) can learn from each other, while CG-POOL integrates high-level information with pairwise information without differentiating and thus drops structural information in different levels of session representation.

### G. Parameter Sensitivity Analysis (RQ5)

We introduce a hyper-parameter to CM-GNN to control the magnitude of contrastive learning, i.e.,  $\beta$ . To investigate the influence of it, we report the performance with a set of representative  $\beta$  values in { 0.001, 0.01, 0.02, 0.03, 0.05 }.

As can be seen in Fig. 5, our model performs the best on Diginetica and Nowplaying when  $\beta$  is 0.001, while for Tmall, setting  $\beta$  to 0.01 achieves the best performance. Besides, it can be observed that as  $\beta$  increases, the performance of CM-GNN drops. One possible reason is that there is a gradient conflict between the recommendation task and the contrastive task, which means that when we use the contrastive task to train CM-GNN, we should choose an appropriate  $\beta$  value.

## V. CONCLUSION

In this paper, we propose a novel contrastive multi-level graph neural network (CM-GNN) for SBR. Specifically, we design a local-level graph convolutional network (L-GCN) module and a global-level graph convolutional network (G-GCN) module to capture pairwise relations of the current session and all the sessions, respectively. CM-GNN introduce an attention-based fusion module to learn the pairwise relation-based session representation by fusing the item representations generated by the above modules. Moreover, since existing methods fail to capture the high-order item transition information, CM-GNN designs a hyper-level graph convolutional network (H-GCN) module to capture the high-order item transition patterns. CM-GNN learns the high-order relation-based session representation by averaging the item representations outputted by H-GCN module. Finally, to convert the high-order item transition information into the pairwise relation-based session representation, we apply contrastive learning between different levels of session representation to enhance the training process of the network. In this paper, we have carried out a large number of experiments on CM-GNN, and comprehensive studies show that our CM-GNN outperforms eleven baselines over three benchmark datasets consistently.

In the future, we plan to study the model collapse problem caused by contrastive learning and consider applying a divergence constraint during the training process of CM-GNN to alleviate this problem.

## ACKNOWLEDGMENT

This work was supported in part by National Natural Science Foundation of China (Nos. 61976127, 61702491).

## REFERENCES

- [1] W. Min, S. Jiang, and R. Jain, "Food recommendation: Framework, existing solutions, and challenges," *IEEE Transactions on Multimedia*, vol. 22, no. 10, pp. 2659–2671, 2019.
- [2] X. Gao, F. Feng, X. He, H. Huang, X. Guan, C. Feng, Z. Ming, and T.-S. Chua, "Hierarchical attention network for visually-aware food recommendation," *IEEE Transactions on Multimedia*, vol. 22, no. 6, pp. 1647–1659, 2019.
- [3] X. Chen, D. Liu, Z. Xiong, and Z.-J. Zha, "Learning and fusing multiple user interest representations for micro-video and movie recommendations," *IEEE Transactions on Multimedia*, vol. 23, pp. 484–496, 2020.
- [4] E. Quintanilla, Y. Rawat, A. Sakryukin, M. Shah, and M. Kankanhalli, "Adversarial learning for personalized tag recommendation," *IEEE Transactions on Multimedia*, vol. 23, pp. 1083–1094, 2020.
- [5] L. Sang, M. Xu, S. Qian, M. Martin, P. Li, and X. Wu, "Context-dependent propagating based video recommendation in multimodal heterogeneous information networks," *IEEE Transactions on Multimedia*, 2020.
- [6] J. Li, P. Ren, Z. Chen, Z. Ren, T. Lian, and J. Ma, "Neural attentive session-based recommendation," in *Proceedings of the 2017 ACM Conference on Information and Knowledge Management*, 2017, pp. 1419–1428.
- [7] Q. Liu, Y. Zeng, R. Mokhosi, and H. Zhang, "Stamp: short-term attention/memory priority model for session-based recommendation," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 1831–1839.
- [8] G. Shani, D. Heckerman, R. I. Brafman, and C. Boutilier, "An mdp-based recommender system," *Journal of Machine Learning Research*, vol. 6, no. 9, 2005.
- [9] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based collaborative filtering recommendation algorithms," in *Proceedings of the 10th international conference on World Wide Web*, 2001, pp. 285–295.
- [10] S. Rendle, C. Freudenthaler, and L. Schmidt-Thieme, "Factorizing personalized markov chains for next-basket recommendation," in *Proceedings of the 19th international conference on World wide web*, 2010, pp. 811–820.
- [11] R. He and J. McAuley, "Fusing similarity models with markov chains for sparse sequential recommendation," in *2016 IEEE 16th International Conference on Data Mining (ICDM)*. IEEE, 2016, pp. 191–200.
- [12] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk, "Session-based recommendations with recurrent neural networks," *arXiv preprint arXiv:1511.06939*, 2015.
- [13] F. Yu, Q. Liu, S. Wu, L. Wang, and T. Tan, "A dynamic recurrent model for next basket recommendation," in *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, 2016, pp. 729–732.
- [14] C.-Y. Wu, A. Ahmed, A. Beutel, A. J. Smola, and H. Jing, "Recurrent recommender networks," in *Proceedings of the tenth ACM international conference on web search and data mining*, 2017, pp. 495–503.
- [15] D. Jannach and M. Ludewig, "When recurrent neural networks meet the neighborhood for session-based recommendation," in *Proceedings of the Eleventh ACM Conference on Recommender Systems*, 2017, pp. 306–310.
- [16] Y. K. Tan, X. Xu, and Y. Liu, "Improved recurrent neural networks for session-based recommendations," in *Proceedings of the 1st workshop on deep learning for recommender systems*, 2016, pp. 17–22.
- [17] B. Hidasi and A. Karatzoglou, "Recurrent neural networks with top-k gains for session-based recommendations," in *Proceedings of the 27th ACM international conference on information and knowledge management*, 2018, pp. 843–852.
- [18] C. Subakan, M. Ravanelli, S. Cornell, M. Bronzi, and J. Zhong, "Attention is all you need in speech separation," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*. IEEE, 2021, pp. 21–25.
- [19] W.-C. Kang and J. McAuley, "Self-attentive sequential recommendation," in *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2018, pp. 197–206.
- [20] J. Li, Y. Wang, and J. McAuley, "Time interval aware self-attention for sequential recommendation," in *Proceedings of the 13th international conference on web search and data mining*, 2020, pp. 322–330.
- [21] R. Ren, Z. Liu, Y. Li, W. X. Zhao, H. Wang, B. Ding, and J.-R. Wen, "Sequential recommendation with self-attentive multi-adversarial network," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 89–98.
- [22] H. Wang, G. Liu, A. Liu, Z. Li, and K. Zheng, "Dmran: A hierarchical fine-grained attention-based network for recommendation," in *IJCAI*, 2019, pp. 3698–3704.
- [23] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, and P. Jiang, "Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer," in *Proceedings of the 28th ACM international conference on information and knowledge management*, 2019, pp. 1441–1450.
- [24] Z. Pan, F. Cai, W. Chen, H. Chen, and M. de Rijke, "Star graph neural networks for session-based recommendation," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 1195–1204.
- [25] S. Wu, Y. Tang, Y. Zhu, L. Wang, X. Xie, and T. Tan, "Session-based recommendation with graph neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 346–353.
- [26] R. Qiu, J. Li, Z. Huang, and H. Yin, "Rethinking the item order in session-based recommendation with graph neural networks," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 579–588.

- [27] C. Xu, P. Zhao, Y. Liu, V. S. Sheng, J. Xu, F. Zhuang, J. Fang, and X. Zhou, "Graph contextualized self-attention network for session-based recommendation," in *IJCAI*, vol. 19, 2019, pp. 3940–3946.
- [28] T. Chen and R. C.-W. Wong, "Handling information loss of graph neural networks for session-based recommendation," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 1172–1180.
- [29] A. Zimdars, D. M. Chickering, and C. Meek, "Using temporal data for making recommendations," *arXiv preprint arXiv:1301.2320*, 2013.
- [30] M. Wang, P. Ren, L. Mei, Z. Chen, J. Ma, and M. de Rijke, "A collaborative session-based recommendation approach with parallel memory modules," in *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2019, pp. 345–354.
- [31] S. Wang, L. Hu, Y. Wang, L. Cao, Q. Z. Sheng, and M. Orgun, "Sequential recommender systems: challenges, progress and prospects," *arXiv preprint arXiv:2001.04830*, 2019.
- [32] S. Wang, L. Cao, Y. Wang, Q. Z. Sheng, M. A. Orgun, and D. Lian, "A survey on session-based recommender systems," *ACM Computing Surveys (CSUR)*, vol. 54, no. 7, pp. 1–38, 2021.
- [33] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and S. Y. Philip, "A comprehensive survey on graph neural networks," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 1, pp. 4–24, 2020.
- [34] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [35] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, "Lightgcn: Simplifying and powering graph convolution network for recommendation," in *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, 2020, pp. 639–648.
- [36] Y.-H. Chen, L. Huang, C.-D. Wang, and J.-H. Lai, "Hybrid-order gated graph neural network for session-based recommendation," *IEEE Transactions on Industrial Informatics*, vol. 18, no. 3, pp. 1458–1467, 2021.
- [37] J. Wang, K. Ding, L. Hong, H. Liu, and J. Caverlee, "Next-item recommendation with sequential hypergraphs," in *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, 2020, pp. 1101–1110.
- [38] N. Yadati, M. Nimishakavi, P. Yadav, V. Nitin, A. Louis, and P. Talukdar, "Hypergen: A new method of training graph convolutional networks on hypergraphs," *arXiv preprint arXiv:1809.02589*, 2018.
- [39] Y. Feng, H. You, Z. Zhang, R. Ji, and Y. Gao, "Hypergraph neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3558–3565.
- [40] X. Xia, H. Yin, J. Yu, Q. Wang, L. Cui, and X. Zhang, "Self-supervised hypergraph convolutional networks for session-based recommendation," *arXiv preprint arXiv:2012.06852*, 2020.
- [41] R. Ye, Q. Zhang, and H. Luo, "Cross-session aware temporal convolutional network for session-based recommendation," in *2020 International Conference on Data Mining Workshops (ICDMW)*. IEEE, 2020, pp. 220–226.
- [42] Z. Wang, W. Wei, G. Cong, X.-L. Li, X.-L. Mao, and M. Qiu, "Global context enhanced graph neural networks for session-based recommendation," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 169–178.
- [43] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [44] K. Hassani and A. H. Khasahmadi, "Contrastive multi-view representation learning on graphs," in *International Conference on Machine Learning*. PMLR, 2020, pp. 4116–4126.
- [45] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *arXiv preprint arXiv:1807.03748*, 2018.
- [46] R. D. Hjelm, A. Fedorov, S. Lavoie-Marchildon, K. Grewal, P. Bachman, A. Trischler, and Y. Bengio, "Learning deep representations by mutual information estimation and maximization," *arXiv preprint arXiv:1808.06670*, 2018.
- [47] T. H. Trinh, M.-T. Luong, and Q. V. Le, "Selfie: Self-supervised pretraining for image embedding," *arXiv preprint arXiv:1906.02940*, 2019.
- [48] L. Wu, H. Lin, Z. Gao, C. Tan, S. Li *et al.*, "Self-supervised on graphs: Contrastive, generative, or predictive," *arXiv preprint arXiv:2105.07342*, vol. 2, no. 6, p. 16, 2021.
- [49] P. Sermanet, C. Lynch, Y. Chebotar, J. Hsu, E. Jang, S. Schaal, S. Levine, and G. Brain, "Time-contrastive networks: Self-supervised learning from video," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 1134–1141.
- [50] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," *arXiv preprint arXiv:1906.00910*, 2019.
- [51] X. Ji, J. F. Henriques, and A. Vedaldi, "Invariant information clustering for unsupervised image classification and segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9865–9874.
- [52] M. Ye, X. Zhang, P. C. Yuen, and S.-F. Chang, "Unsupervised embedding learning via invariant and spreading instance feature," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6210–6219.
- [53] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.
- [54] S. Ma, Z. Zeng, D. McDuff, and Y. Song, "Active contrastive learning of audio-visual video representations," *arXiv preprint arXiv:2009.09805*, 2020.
- [55] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [56] S. Arora, H. Khandeparkar, M. Khodak, O. Plevrakis, and N. Saunshi, "A theoretical analysis of contrastive unsupervised representation learning," *arXiv preprint arXiv:1902.09229*, 2019.
- [57] H. Fang, S. Wang, M. Zhou, J. Ding, and P. Xie, "Cert: Contrastive self-supervised learning for language understanding," *arXiv preprint arXiv:2005.12766*, 2020.
- [58] D. Iter, K. Guu, L. Lansing, and D. Jurafsky, "Pretraining with contrastive sentence objectives improves discourse performance of language models," *arXiv preprint arXiv:2005.10389*, 2020.
- [59] X. Xin, A. Karatzoglou, I. Arapakis, and J. M. Jose, "Self-supervised reinforcement learning for recommender systems," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 931–940.
- [60] K. Zhou, H. Wang, W. X. Zhao, Y. Zhu, S. Wang, F. Zhang, Z. Wang, and J.-R. Wen, "S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 1893–1902.
- [61] J. Wu, X. Wang, F. Feng, X. He, L. Chen, J. Lian, and X. Xie, "Self-supervised graph learning for recommendation," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 726–735.
- [62] J. Yu, H. Yin, M. Gao, X. Xia, X. Zhang, and N. Q. V. Hung, "Socially-aware self-supervised tri-training for recommendation," *arXiv preprint arXiv:2106.03569*, 2021.
- [63] A. Sankar, Y. Wu, Y. Wu, W. Zhang, H. Yang, and H. Sundaram, "Groupim: A mutual information maximization framework for neural group recommendation," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 1279–1288.
- [64] J. Yu, H. Yin, J. Li, Q. Wang, N. Q. V. Hung, and X. Zhang, "Self-supervised multi-channel hypergraph convolutional network for social recommendation," in *Proceedings of the Web Conference 2021*, 2021, pp. 413–424.
- [65] E. Zangerle, M. Pichl, W. Gassler, and G. Specht, "# nowplaying music dataset: Extracting listening behavior from twitter," in *Proceedings of the first international workshop on internet-scale multimedia management*, 2014, pp. 21–26.