

Knowledge Enhanced Multi-intent Transformer Network for Recommendation

Ding Zou*

CCIIP Lab
School of Computer Science and
Technology
Huazhong University of Science and
Technology
Joint Laboratory of HUST and Pingan
Property & Casualty Research (HPL)
Wuhan, China
Taotian Group
Hangzhou, China
m202173662@hust.edu.cn

Wei Wei†

CCIIP Lab
School of Computer Science and
Technology
Huazhong University of Science and
Technology
Joint Laboratory of HUST and Pingan
Property & Casualty Research (HPL)
Wuhan, China
weiw@hust.edu.cn

Feida Zhu

Singapore Management University
Singapore, Singapore
fdzhu@smu.edu.sg

Chuanyu Xu

Taotian Group
Hangzhou, China
tracy.xcy@taobao.com

Tao Zhang

Taotian Group
Hangzhou, China
guyan.zt@taobao.com

Chengfu Huo

Taotian Group
Hangzhou, China
chengfu.huocf@taobao.com

ABSTRACT

Incorporating Knowledge Graphs (KGs) into Recommendation has attracted growing attention in industry, due to the great potential of KG in providing abundant supplementary information and interpretability for the underlying models. However, simply integrating KG into recommendation usually brings in negative feedback in industry, mainly due to the ignorance of the following two factors: i) users' multiple intents, which involve diverse nodes in KG. For example, in e-commerce scenarios, users may exhibit preferences for specific styles, brands, or colors. ii) knowledge noise, which is a prevalent issue in Knowledge Enhanced Recommendation (KGR) and even more severe in industry scenarios. The irrelevant knowledge properties of items may result in inferior model performance compared to approaches that do not incorporate knowledge. To tackle these challenges, we propose a novel approach named Knowledge Enhanced Multi-intent Transformer Network for Recommendation (KGTN), which comprises two primary modules: Global Intents Modeling with Graph Transformer, and Knowledge Contrastive Denoising under Intents. Specifically, Global Intents with Graph Transformer focuses on capturing learnable user intents, by incorporating global signals from user-item-relation-entity interactions with a well-designed graph transformer,

and meanwhile learning intent-aware user/item representations. On the other hand, Knowledge Contrastive Denoising under Intents is dedicated to learning precise and robust representations. It leverages the intent-aware user/item representations to sample relevant knowledge, and subsequently proposes a local-global contrastive mechanism to enhance noise-irrelevant representation learning. Extensive experiments conducted on three benchmark datasets show the superior performance of our proposed method over the state-of-the-arts. And online A/B testing results on Alibaba large-scale industrial recommendation platform also indicate the real-scenario effectiveness of KGTN. The implementations are available at: <https://github.com/CCIPLab/KGTN>.

CCS CONCEPTS

• **Information systems** → **Recommender systems**.

KEYWORDS

Knowledge Enhanced Recommendation, Graph Transformer, Graph Neural Networks

ACM Reference Format:

Ding Zou, Wei Wei, Feida Zhu, Chuanyu Xu, Tao Zhang, and Chengfu Huo. 2024. Knowledge Enhanced Multi-intent Transformer Network for Recommendation. In *Companion Proceedings of the ACM Web Conference 2024 (WWW '24 Companion)*, May 13–17, 2024, Singapore, Singapore. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3589335.3648296>

1 INTRODUCTION

Knowledge graphs (KGs) have emerged as a promising approach to enhance the accuracy and interpretability of recommender systems in both academic and industry scenarios. By incorporating entities and relations, KGs provide a rich source of information for user/item representation learning, which not only captures the diverse relationships among items (such as the same item brand),

*Work done during internship at Taotian Group.

†Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
WWW '24 Companion, May 13–17, 2024, Singapore, Singapore.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0172-6/24/05...\$15.00
<https://doi.org/10.1145/3589335.3648296>



Figure 1: (a) A simple case for illustrating multiple user intents with global information; (b) Performance comparison.

but also allows for the interpretation of user preferences (such as attributing a user’s selection of a clothing to its fashionable style).

In an effort to effectively integrate the item-side KG information into recommendation, considerable research efforts have been devoted to Knowledge Enhanced Recommendation (*aka*. KGR). Early studies [7, 19, 33] directly integrate knowledge graph embeddings with items to enhance their representations. Some subsequent studies [6, 16, 25] enrich the interactions via meta-paths that capture relevant connectivities between users and items with KG. They either select prominent paths over KG [17], or represent the interactions with multi-hop paths from users to items [6, 25]. Nevertheless, most of them heavily rely on manually designed meta-paths, which makes it hard to optimize in reality. As a result, later methods have embraced Graph Neural Networks (GNNs) [22, 23] to automatically aggregate high-order information over KG, which iteratively integrate multi-hop neighbors into representations and have demonstrated promising performance for recommendation. Most recently, there have been efforts to incorporate Contrastive Learning (CL) into KGR for addressing noisy knowledge and long-tail problems [27, 29, 37] via contrasting the user-item (collaborative part) and item-entity (knowledge part) graphs.

However, current KGR methods usually bring poor performance in large-scale industry scenarios, due to their commonly overlooking two crucial factors: 1) Users’ multiple intents underlying interaction behavior. For instance, as depicted in Figure 1(a), users may have diverse intentions when shopping in Alibaba E-commerce platform, such as long-term interest, passing time, or social reason, *etc.* 2) Redundant Knowledge information. In the context of user intents, some knowledge facts in the KG may be irrelevant noise [3], which can potentially disrupt the learning process of user/item representations. As shown in Figure 1(b), incorporating KGs may result in a worse model performance than the models without KG utilization (the details of comparison could refer to Section 4.2).

But still, it’s not trivial to model user intents in KGR, since user intents may be composed of multiple heterogeneous information, including items, relations, and entities. Previous multi-intent modeling methods usually define the intents as a linear combination of either interacted items [24] or entire relation sets [23], then update the intent representations through local aggregation in the user-intent-item heterogeneous graph. Nevertheless, such a multi-intent learning paradigm may not fully meet the requirements for KGR, as it neglects the global information in intent defining and learning. To illustrate this, we present an example in Figure 1(a). In this example, user u_1 may purchase the item i_1 for the intent c_1 of long-term

interest, resulting in a focus on clothing style (*e.g.*, whether it is fashionable), which means intent c_1 is associated with KG relation r_1 and entity e_1 ; while u_1 may buy the item i_n for the intent c_k of social reason (such as friend u_2 recommend), which means intent c_k is associated with user u_2 and item i_k .

In this paper, we focus on modeling user intents behind interaction behaviors with global collaborative (user-item) and knowledge (item-relation-entity) information, and exploiting these modeled intents to guide knowledge sampling, facilitating fine-grained and accurate user/item representation learning. We propose a novel model, KGTN, which comprises two essential components for solving the foregoing limitations: i) Global Intents Modeling with Graph Transformer. We predefine K intent representations for user/item, then learn these intents with global information from collaborative and knowledge graphs. Specifically, it first merges knowledge information into items, then propose a novel graph transformer in the user-item graph to learn global intents and generate intent-aware user/item representations. ii) Knowledge Contrastive Denoising under Intents. KGTN first exploits the intent-aware user/item representations to guide the knowledge sampling, effectively pruning the irrelevant knowledge. Then a novel local-global contrastive mechanism is proposed here to denoise the user/item representations. Empirically, KGTN outperforms the state-of-the-art models on three benchmark datasets in offline testing, and achieves significant improvements in online A/B testing.

Our contributions of this work can be summarized as follows:

- **General Aspects:** We emphasize the importance of intent modeling with global information, which plays a crucial role in fine-grained representation learning and knowledge denoising.
- **Novel Methodologies:** We propose a novel model KGTN, which models user intents from global signals with a novel graph transformer; and denoises item representations with i) knowledge denoising under intents, and ii) local-global graph contrastive learning.
- **Multifaceted Experiments:** We conduct extensive offline experiments on three benchmark datasets and online A/B testing on Alibaba recommendation platform. The results demonstrate the advantages of our KGTN in better representation learning.

2 PROBLEM FORMULATION

In this section, we begin by formulating the structural data of CF (user-item interactions) and KG (item-relation-entity knowledge) in KGR, then present the problem statement.

Interaction Data. In a typical recommendation scenario, let $\mathcal{U} = \{u_1, u_2, \dots, u_M\}$ be a set of M users and $\mathcal{V} = \{v_1, v_2, \dots, v_N\}$ a set of N items. Let $Y \in \mathbf{R}^{M \times N}$ be the user-item interaction matrix, where $y_{uv} = 1$ indicates that user u engaged with item v , such as behaviors like clicking or purchasing; otherwise $y_{uv} = 0$.

Knowledge Graph. A KG stores luxuriant real-world facts associated with items, encompassing item attributes or external commonsense knowledge, in the form of a heterogeneous graph [16]. Let $\mathcal{G} = \{(h, r, t) \mid h, t \in \mathcal{E}, r \in \mathcal{R}\}$ be the KG, where h, r, t represent the head, relation, tail of a knowledge triple, respectively; \mathcal{E} and \mathcal{R} denote the sets of entities and relations in \mathcal{G} . In many recommendation scenarios, an item $v \in \mathcal{V}$ corresponds to one entity $e \in \mathcal{E}$. We hence establish a set of item-entity alignments

$\mathcal{A} = \{(v, e) | v \in \mathcal{V}, e \in \mathcal{E}\}$, where (v, e) indicates that item v can be aligned with an entity e in KG. With the alignments between items and KG entities, KG is able to profile items and offer complementary information to the interaction data.

Problem Statement. Given the user-item interaction matrix Y and the KG \mathcal{G} , KGR aims to learn a function that can predict how likely a user would adopt an item.

3 METHODOLOGY

We now present the proposed Knowledge Enhanced Multi-intent Transformer Network for Recommendation (KG_{TN}). KG_{TN} aims at modeling user intents with global information and exploiting user intents to denoise KG for accurate and robust user/item representation learning. Figure 2 displays the framework of KG_{TN}, which mainly consists of two key components: 1) Global Intent Modeling with graph transformer. Initially, KG_{TN} defines a set of K learnable global intents for users and items. It then models these intents and learns intent-aware user/item representations, via integrating global signals with a graph transformer in the user-item graph, where knowledge information has been encoded into items. 2) Knowledge Contrastive Denoising under intents. It first exploits the learned intent-aware user/item representations to sample intent-relevant knowledge, then designs a contrastive self-supervised task between the local aggregation and global aggregation features within the sampled graph to facilitate robust representation learning.

3.1 Global Intents Modeling with Graph Transformer

3.1.1 Intent Initialization with Global signals. When interacting with items, users often have diverse intents, such as preferences for specific clothing brands and styles, friends recommending, or passing time with randomly clicking [14, 23]. To capture these diverse intents, we assume K different intents c_u and c_v from the user and item sides, respectively, where the intents on the item side can also be understood as the theme or context of the item, for example, a user who intends to purchase a fashionable dress may like clothes of “young” topic. Our predictive objective of user-item preference can be presented as follows:

$$\int_{c_u} \int_{c_v} P(y, c_u, c_v | u, v) dc_v dc_u = \sum_k P(y, c_u^k, c_v^k | u, v). \quad (1)$$

Specifically, we define K global intent prototypes $\{c_u^k \in \mathbb{R}^d\}_{k=1}^K$ and $\{c_v^k \in \mathbb{R}^d\}_{k=1}^K$ for user and item, respectively. With these pre-defined intent prototypes, we then are supposed to integrate them into user/item representations, and update them with related global signals.

3.1.2 Intent Modeling with graph transformer. Towards accurately modeling user intents with global information and learning intent-aware user/item representations, we perform an intent-aware information propagation with these learnable intents. Specifically, intent-aware user/item embeddings are acquired by an attentive sum of the intent prototypes, and user/item embeddings of each layer are updated by aggregating the global user/item/relation/entity signals.

Formally, we could get intent-aware user/item representations at the l -th user/item embedding layer, by aggregating information across different K learnable intent prototypes (including c_u and c_v), using the following design:

$$\mathbf{e}_u^l = \sum_k^K c_u^k P(c_u^k | \mathbf{e}_u^l), \quad (2)$$

$$P(c_u^k | \mathbf{e}_u^l) = \frac{\eta(\mathbf{e}_u^{l-1\top} c_u^k)}{\sum_{k'}^K \eta(\mathbf{e}_u^{l-1\top} c_{k'}^k)}, \quad (3)$$

where the $P(c_u^k | \mathbf{e}_u^l)$ and $P(c_v^k | \mathbf{e}_v^l)$ denotes the importance score of c_u^k for l -th user embeddings that has encodes the global signals. Similarly, the $P(c_v^k | \mathbf{e}_v^l)$ denotes the importance score of c_v^k for l -th item embeddings.

As for the way of calculating the l -th user/item embeddings, we propose to adopt a two-step process to encode the global user/item/relation/entity information in the whole heterogeneous graph. The first step is to merge the knowledge information (including both relation and entity) into item embeddings with a proposed relation-aware graph aggregation, making the item representation more comprehensive and informative. It injects the relational context into the embeddings of the neighboring entities, and weighting them with the knowledge rationale scores (It's worth noting that items are a subset of knowledge entities), as follows:

$$\begin{aligned} \mathbf{e}_i^{(l+1)} &= \frac{1}{|N_i|} \sum_{(r,v) \in N_i} \beta(i, r, v) \mathbf{e}_r \odot \mathbf{e}_v^{(l)}, \\ \beta(i, r, v) &= \text{softmax} \left((\mathbf{e}_i || \mathbf{e}_r)^T \cdot (\mathbf{e}_v || \mathbf{e}_r) \right) \\ &= \frac{\exp \left((\mathbf{e}_i || \mathbf{e}_r)^T \cdot (\mathbf{e}_v || \mathbf{e}_r) \right)}{\sum_{(v',r) \in \hat{N}(i)} \exp \left((\mathbf{e}_i || \mathbf{e}_r)^T \cdot (\mathbf{e}_{v'} || \mathbf{e}_r) \right)}, \end{aligned} \quad (4)$$

where $||$ denotes concat operation, N_i denotes the set of neighboring entities.

Then the second step is to apply a novel graph transformer among user-item graph, which encodes global user/item/entity information into user/item representations. By doing so, the user/item representations of each layer are integrated with global signals, which would be exploited into intent modeling and representation updating, as follows:

$$\begin{aligned} \mathbf{e}_u^{l+1} &= \sum_i \prod_{h=1}^H m_{u,i} \beta_{u,i}^h \mathbf{W}_V^h \mathbf{e}_i^l; m_{u,i} = \begin{cases} 1 & \text{if } (u, i) \in Y \\ 0 & \text{otherwise} \end{cases} \\ \beta_{u,i}^h &= \frac{\exp \beta_{u,i}^h}{\sum_i \exp \beta_{u,i}^h}; \bar{\beta}_{u,v'}^h = \frac{(\mathbf{W}_Q^h \cdot \mathbf{e}_u^l)^\top \cdot (\mathbf{W}_K^h \cdot \mathbf{e}_i^l)}{\sqrt{d/H}}, \end{aligned} \quad (5)$$

where H denotes the number of attention heads (indexed by h). $m_{u,v'}$ is the binary indicator to decide whether to calculate the attentive relations between user u and item i . $\beta_{u,i}^h$ denotes the attention weight for user-item interaction pair (u, i) w.r.t. the h -th head representation space. $\mathbf{W}_Q^h, \mathbf{W}_K^h, \mathbf{W}_V^h \in \mathbb{R}^{d/H \times d}$ denotes the query, key, the value embedding projection for the h -th head, respectively.

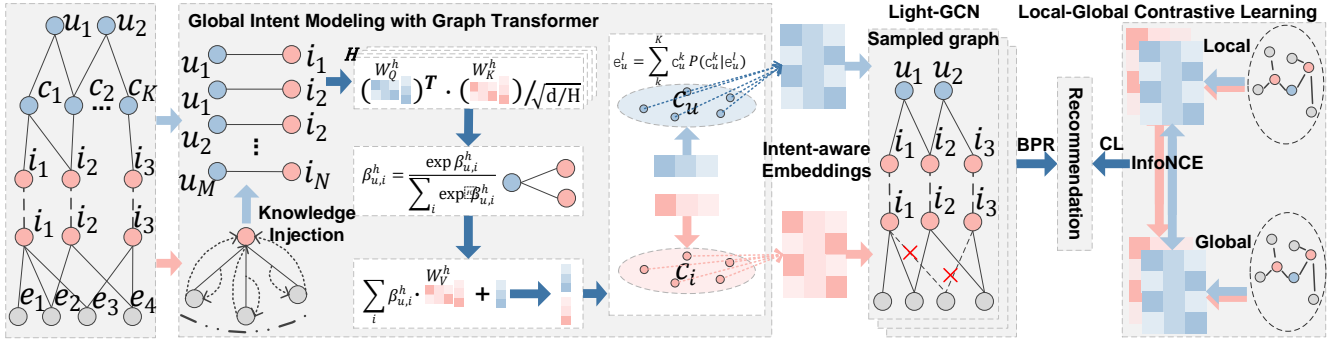


Figure 2: Overall framework illustration of the proposed KGTN model. Best viewed in color.

By integrating global information into users/items, we could learn intent-aware user/item representations and update the learnable intents according to Equation 2.

3.2 Knowledge Contrastive Denoising under Intents

It is intuitive that noisy or irrelevant connections between entities in knowledge graphs can lead to suboptimal representation learning, which is opposite to original purpose of introducing the KG. To eliminate the noise effect in the KG and distill informative signals that benefit the recommendation task, we propose to highlight important connections consistent to user intents, while removing the irrelevant ones.

3.2.1 Knowledge Sampling under intents. With the intent-aware user/item representations, we then try to denoise the item-entity graph by removing the irrelevant edges and nodes and sampling the important ones. We first exploit the intent-aware representations to calculate the importance score of knowledge triplets (*i.e.*, the item-relation-entity pairs) same as Equation 4, then add the Gumbel noise [8] to the learned importance scores to improve the sampling robustness, as follows:

$$\begin{aligned} \beta(i, r, v) &= \text{softmax} \left((\mathbf{e}_i || \mathbf{e}_r)^T \cdot (\mathbf{e}_v || \mathbf{e}_r) \right) \\ \beta(i, r, v) &= \beta(i, r, v) - \log(-\log(\epsilon)); \quad \epsilon \sim \text{Uniform}(0, 1), \end{aligned} \quad (6)$$

where ϵ is a random variable sampled from a uniform distribution. Then it follows a top-k sampling strategy for generating the new item-entity graph that removes the irrelevant edges and nodes:

$$\widehat{\beta}(i, r, v) = \begin{cases} \beta(i, r, v), & \beta(i, r, v) \in \text{top-k}(\beta(i, r, v)), \\ 0, & \text{otherwise,} \end{cases} \quad (7)$$

where $\widehat{\beta}(i, r, v)$ is the sampled triples in item-entity graph, which would be used to replace the original graph structure in the following user/item representation learning.

3.2.2 Local-Global Knowledge Contrastive Learning. With the sampled item-entity graph, we then propose to iteratively update the intent-aware representations in it. And inspired by previous contrastive learning based methods that align the item representations

from KG and CF to denoise, we further propose a local-global contrastive mechanism to improve the robustness of representation learning.

Specifically, we exploit the user-item graph and sampled item-entity graph to perform light information aggregation with intent-aware user/item representations $\mathbf{e}_u, \mathbf{e}_i$ as input $\mathbf{z}_u^{(0)}, \mathbf{z}_i^{(0)}$, for acquiring a robust and effective intent-aware user/item representations, as follows:

$$\begin{aligned} \mathbf{z}_i^{(l+1)} &= \frac{1}{|N_i|} \sum_{(r,v) \in N_i} \mathbf{e}_r \odot \mathbf{z}_v^{(l)}, \\ \mathbf{z}_u^{(l+1)} &= \frac{1}{|N_u|} \sum_{i \in N_u} \mathbf{z}_i^{(l)}, \end{aligned} \quad (8)$$

where $\mathbf{z}_u^{(0)}, \mathbf{z}_i^{(0)}$ memorize the global signals, and we hence get final representations of user/item $\mathbf{z}_u^{(l)}, \mathbf{z}_i^{(l)}$ ($l \in L$).

Besides the supervised user/item representation learning, we propose to perform a contrastive learning between the nodes embeddings that encode global signals and local signals, which is different from traditional cl-based methods that contrast the CF and KG parts. We perform information aggregation in the sampled graph with the initial user/item representations $\mathbf{e}_u, \mathbf{e}_i$ to acquire the local results $\mathbf{z}_{u,local}^{(l)}, \mathbf{z}_{i,local}^{(l)}$ ($l \in L$), while utilizing the intent-aware user/item representations $\mathbf{e}_u, \mathbf{e}_i$ that contains global signals to acquire the global results $\mathbf{z}_u^{(l)}, \mathbf{z}_i^{(l)}$ ($l \in L$). Then perform layer-wise contrastive learning between local and global results.

The local aggregation layer embeddings $\mathbf{z}_{u,local}^{(l)}, \mathbf{z}_{i,local}^{(l)}$ and global aggregation layer embeddings $\mathbf{z}_u^{(l)}, \mathbf{z}_i^{(l)}$ are made to be contrasted in a layer-wise way. We generate each positive pair using the embeddings of the same user (item) from the local view and each of the global view, and other nodes form the negative pairs. We could get the contrastive loss of users as follows:

$$\mathcal{L}_c^u = \frac{1}{L} \sum_{l=0}^L -\log \frac{\exp(s(\mathbf{z}_u^{(l)}, \mathbf{z}_{u,local}^{(l)})/\tau)}{\sum_{k \neq u} \exp(s(\mathbf{z}_u^{(l)}, \mathbf{z}_k^{(l)})/\tau) + \sum_{k \neq u} \exp(s(\mathbf{z}_u^{(l)}, \mathbf{z}_{k,local}^{(l)})/\tau)}, \quad (9)$$

where $s(\cdot)$ denotes the cosine similarity calculating, and τ denotes a temperature parameter. And similarly we could get the contrastive loss of item \mathcal{L}_c^i . By summing the two contrastive losses we hence have the total local-global contrastive loss \mathcal{L}_c .

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

Table 1: Statistics for the three datasets.

3.3 Model Prediction

After learning intent-aware user/item representations with global signals and performing contrastive learning between local and global information, we have multi-layer intent-aware representations for user/item. By summing all the layers' representations, we have the final user/item representations and predict their matching score through inner product, as follows:

$$\mathbf{z}_u = \mathbf{z}_u^{(0)} + \dots + \mathbf{z}_u^{(K)}, \quad \mathbf{z}_i = \mathbf{z}_i^{(0)} + \dots + \mathbf{z}_i^{(K)}. \quad (10)$$

$$\hat{y}(u, i) = \mathbf{z}_u^\top \mathbf{z}_i.$$

By adopting a BPR loss [15] to reconstruct the historical data, which encourages the prediction scores of a user's historical items to be higher than the unobserved items, we acquire the supervised loss:

$$\mathcal{L}_{\text{BPR}} = \sum_{(u,i,j) \in \mathcal{O}} -\ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}), \quad (11)$$

where $\mathcal{O} = \{(u, i, j) \mid (u, i) \in \mathcal{O}^+, (u, j) \in \mathcal{O}^-\}$ is the training dataset consisting of the observed interactions \mathcal{O}^+ and unobserved counterparts \mathcal{O}^- ; σ is the sigmoid function.

3.4 Multi-task Training

To combine the recommendation task with the self-supervised task, we optimize the whole model with a multi-task training strategy. We combine the local-global contrastive loss with BPR loss, and learn the model parameter via minimizing the following objective function:

$$\mathcal{L}_{\text{KGTN}} = \mathcal{L}_{\text{BPR}} + \alpha \mathcal{L}_c + \lambda \|\Theta\|_2^2, \quad (12)$$

where Θ is the model parameter set, α is a hyperparameter to determine the local-global contrastive loss ratio, β and λ are two hyperparameters to control the contrastive loss and L_2 regularization term, respectively.

4 EXPERIMENT

Aiming to answer the following research questions, we conduct both offline experiments and online A/B tests on three public datasets and Alibaba online platform:

- **RQ1:** How does KGTN perform, compared to present models?
- **RQ2:** How do the main components in KGTN affect its effectiveness?
- **RQ3:** How do different hyper-parameter settings affect KGTN?
- **RQ4:** How does KGTN perform with noisy injection?

- **RQ5:** How does KGTN perform in a live system serving billions of users?

4.1 Experiment Settings

4.1.1 Dataset and Metrics. Three benchmark datasets are utilized to evaluate the effectiveness of KGTN: Last.FM¹, Book-Crossing², and MovieLens-1M³. The detailed statistics of them are summarized in Table 1, which vary in size and sparsity and make our experiments more convincing. As for the data pre-process, we first follow RippleNet [18] to transform their explicit feedback into implicit one, and randomly sample negative samples from his unwatched items with the size equal to his positive ones to construct the negative parts. As for the sub-KG construction, we follow RippleNet [18] and use Microsoft Satori⁴ to construct it for MovieLens-1M, Book-Crossing, and Last.FM datasets. Each sub knowledge graph that follows the triple format is a subset of the whole KG with a confidence level greater than 0.9.

We evaluate our method in two experimental scenarios: (1) In click-through rate (CTR) prediction, we apply the trained model to predict each interaction in the test set. We adopt two widely used metrics [18, 21] *AUC* and *F1* to evaluate CTR prediction. (2) In top- K recommendation, we use the trained model to select K items with the highest predicted click probability for each user in the test set, and we choose *Recall@K* to evaluate the recommended sets.

4.1.2 Baselines. To demonstrate the effectiveness of our proposed KGTN, we compare it with four types of KGR methods: CF-based methods (BPRMF [15]), embedding-based method (CKE [33]), RippleNet [18]), path-based method (PER [32]), GNN-based methods (KGCN [21], KGNN-LS [20], KGAT [22], CKAN [26], KGIN [23], CG-KGR [3]), CL-based methods (KGCL[29], MCCLK [37]).

4.1.3 Parameter Settings. We implement our KGTN and all baselines in Pytorch and carefully tune the key parameters. For a fair comparison, we fix the embedding size to 64 for all models, and the embedding parameters are initialized with the Xavier method [4]. We optimize our method with Adam [9] and set the batch size to 2048. A grid search is conducted to confirm the optimal settings, we tune the learning rate η among {0.0001, 0.0003, 0.001, 0.003} and λ of L_2 regularization term among $\{10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$. Other hyper-parameter settings are provided in Table 1. The best settings for hyper-parameters in all comparison methods are researched by either empirical study or following the original papers.

4.2 Performance Comparison (RQ1)

We report the empirical results of all methods in Table 2 and Table 3. The improvements and statistical significance test are performed between KGTN and the strongest baselines (highlighted with underline). Analyzing such performance comparison, we have the following observations:

- **Our proposed KGTN achieves the best results.** KGTN consistently outperforms all baselines across three datasets in terms of all measures, which achieves significant improvements over the

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

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

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

⁴<https://searchengineland.com/library/bing/bing-satori>

Model	Book-Crossing		MovieLens-1M		Last.FM	
	AUC	F1	AUC	F1	AUC	F1
BPRMF	0.6583(−13.18%)	0.6117(−7.59%)	0.8920(−4.52%)	0.7921(−7.21%)	0.7563(−13.41%)	0.7010(−9.95%)
CKE	0.6759(−11.42%)	0.6235(−6.41%)	0.9065(−3.07%)	0.8024(−6.18%)	0.7471(−14.33%)	0.6740(−12.65%)
RippleNet	0.7211(−6.90%)	0.6472(−4.04%)	0.9190(−1.82%)	0.8422(−2.20%)	0.7762(−11.42%)	0.7025(−9.80%)
PER	0.6048(−18.53%)	0.5726(−11.50%)	0.7124(−22.48%)	0.6670(−19.72%)	0.6414(−24.90%)	0.6033(−19.72%)
KGCN	0.6841(−10.60%)	0.6313(−5.63%)	0.9090(−2.82%)	0.8366(−2.76%)	0.8027(−8.77%)	0.7086(−9.19%)
KGNN-LS	0.6762(−11.39%)	0.6314(−5.62%)	0.9140(−2.32%)	0.8410(−2.32%)	0.8052(−8.52%)	0.7224(−7.81%)
KGAT	0.7314(−5.87%)	0.6544(−3.32%)	0.9140(−2.32%)	0.8440(−2.02%)	0.8293(−6.11%)	0.7424(−5.81%)
CKAN	0.7420(−4.81%)	0.6671(−2.05%)	0.9082(−2.90%)	0.8410(−2.32%)	0.8418(−4.86%)	0.7592(−4.13%)
KGIN	0.7273(−6.28%)	0.6614(−2.62%)	0.9190(−1.82%)	0.8441(−2.01%)	0.8486(−4.18%)	0.7602(−4.03%)
CG-KGR	0.7498(−4.03%)	0.6689(−1.87%)	0.9110(−2.62%)	0.8359(−2.83%)	0.8336(−5.68%)	0.7433(−5.72%)
KGCL	0.7453(−4.48%)	0.6679(−1.97%)	0.9184(−1.88%)	0.8437(−2.05%)	0.8455(−4.49%)	0.7596(−4.00%)
MCCLK	0.7625(−2.76%)	0.6777(−0.99%)	0.9252(−1.20%)	0.8559(−0.83%)	0.8663(−2.41%)	0.7753(−2.43%)
KGTN	0.7901*	0.6876*	0.9372*	0.8642*	0.8904*	0.7996*

Table 2: The result of AUC and F1 in CTR prediction. The best results are in boldface and the second best results are underlined. * denotes statistically significant improvement by unpaired two-sample t -test with $p < 0.001$.

Model	Book-Crossing		MovieLens-1M		Last.FM	
	R@10	R@20	R@10	R@20	R@10	R@20
BPRMF	0.0334	0.0525	0.0939	0.1512	0.0923	0.1740
CKE	0.0421	0.0562	0.0867	0.1364	0.0780	0.1532
RippleNet	0.0507	0.0622	0.1082	0.1766	0.0942	0.1520
PER	0.0322	0.0481	0.0523	0.1204	0.0540	0.1167
KGCN	0.0496	0.0540	0.0965	0.1720	0.1416	0.1776
KGNN-LS	0.0422	0.0526	0.1286	0.1757	0.1312	0.1933
KGAT	0.0522	0.0670	0.1468	0.2296	0.1640	0.2313
CKAN	0.0462	0.0566	0.1511	0.2400	0.1412	0.2465
KGIN	0.0555	0.0699	0.1511	0.2404	0.1758	0.2487
CG-KGR	0.0612	0.0781	0.1621	0.2495	0.1578	0.2106
KGCL	0.0679	0.0845	0.1633	0.2499	0.1759	0.2471
MCCLK	0.0769	0.0936	0.1642	0.2503	0.1835	0.2598
KGTN	0.1060*	0.1275*	0.1841*	0.2826*	0.2104*	0.3106*

Table 3: The result of Recall@10 and Recall@20 in top-K recommendation.

strongest baselines *w.r.t.* AUC by 2.76%, 1.20%, and 2.41% in Book, Movie, and Music respectively, and demonstrates its effectiveness. We attribute such improvements to the following aspects: (1) By modeling user intents with global signals, KGTN is able to learn user/item representations in a more fine-grained and comprehensive manner; (2) The knowledge sampling strategy under intents could remove less relevant knowledge information for a robust representation learning; (3) The local-global contrastive learning improves the representation learning in a self-supervised manner, via contrasting the local and global information.

- **Incorporating KG not always benefits recommender system.** Comparing CKE with BPRMF, leaving KG untapped limits the performance of BPRMF, which shows the effectiveness of KG information. While PER gets a worse performance than BPRMF, which means that only incorporating suitable knowledge could benefit the model. This fact stresses the importance of knowledge sampling and knowledge denoising.

- **GNN has a strong power of graph learning.** Most of the GNN-based methods perform better, suggesting the importance of modeling long-range connectivity for graph representation learning. This fact inspires us to go beyond the local aggregation paradigm, and to consider the global signals.
- **Contrastive Learning is effective.** The most recently proposed CL-based methods have the best performance, which shows the effectiveness of incorporating a self-supervised task for improving representation learning. It inspires us to design proper contrastive mechanisms to denoise the knowledge and improve the model performance.

4.3 Ablation Studies (RQ2)

As shown in Figure 3, here we examine the contributions of main components in our model to the final performance by comparing KGTN with the following three variants: 1) KGTN_{w/o S}: In this variant, the knowledge sampling under intents module is removed. 2) KGTN_{w/o C}: This variant removes local-global contrastive mechanism. 3) KGTN_{w/o I}: This variant removes the multi-intent modeling, which means both global intent modeling and knowledge contrastive denoising do not exist in this variant. The results of two variants and KGTN are reported in Figure 3, from which we have the following observations:

- Removing both knowledge sampling and local-global contrasting would degrade model performance, which shows their effectiveness in representation learning.
- Ablating the multi-intent modeling brings the worst performance, which shows the importance of incorporating global signals and considering multiple intents.

4.4 Sensitivity Analysis (RQ3)

4.4.1 Impact of graph transformer depth. To study the influence of graph transformer depth, we vary L in range of {1, 2, 3} on book, movie, and music datasets. As shown in Table 4, KGTN performs best when $L = 1$. It convinces that one iteration is enough for

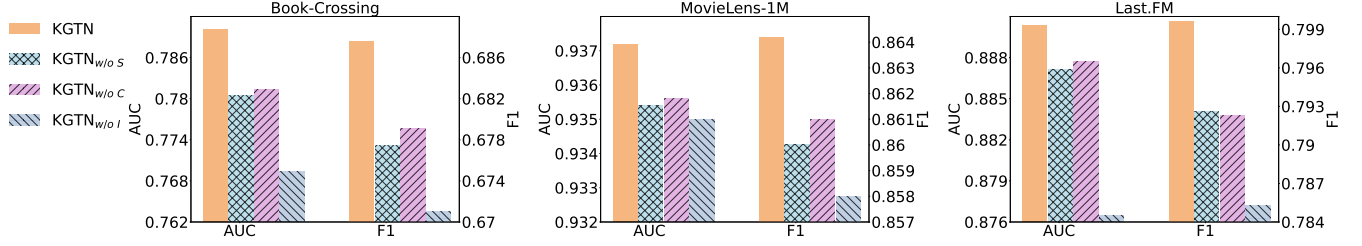
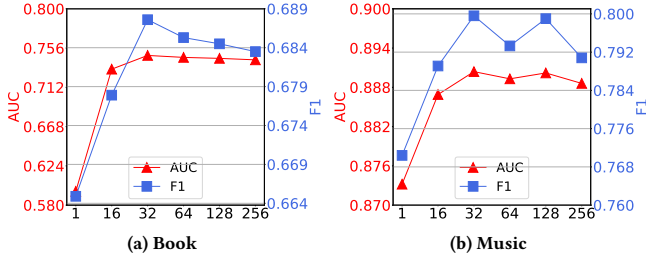
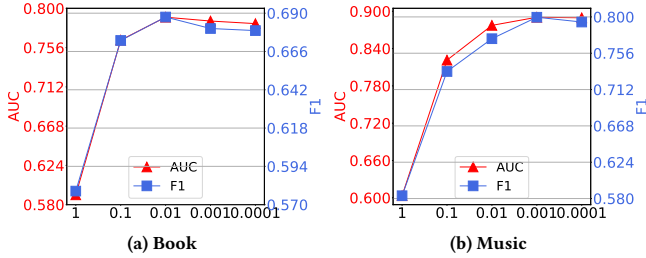


Figure 3: Effect of ablation study.

	Book		Movie		Music	
	Auc	F1	Auc	F1	Auc	F1
$L=1$	0.7901	0.6876	0.9372	0.8642	0.8904	0.7996
$L=2$	0.7743	0.6783	0.9349	0.8623	0.8834	0.8068
$L=3$	0.7603	0.6709	0.9278	0.8481	0.8785	0.7951

Table 4: Impact of graph transformer depth.


 Figure 4: Impact of intent number K .

 Figure 5: Impact of contrastive loss ratio α .

integrating the global signals into user/item representations, which shows its low reliance on model depth.

4.4.2 Impact of intent number K . To investigate the impact of the intent number, we vary it from the range {16, 32, 64, 128, 256} and the model performance is shown in Figure 4, from which we could draw the following conclusions: i) Ignoring the multiple intents ($K = 1$) results in the worst performance, which convinces the effectiveness of incorporating multi-intent modeling. ii) The model performance first arises then drops with the intent increasing. A suitable intents number boosts the model performance with fine-grained preference learning, while too many intents mean too fine-grained modeling and inversely introduce noise into representation learning.

4.4.3 Impact of contrastive loss ratio α . The parameter α determines the importance of the contrastive loss during the multi-task

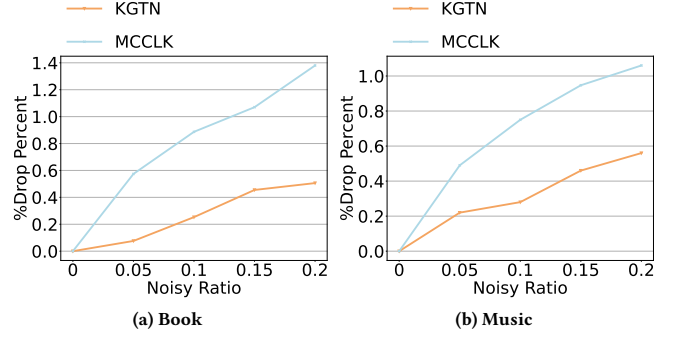


Figure 6: Noise Analysis in Music and Book datasets.

training. Hence we vary it in range {1, 0.1, 0.01, 0.001} to study its influence. As shown in Figure 5, we observe that: $\alpha = 0.1$ brings the best model performance, the main reason is that changing the contrastive loss to a fairly equal level to recommendation task loss could boost the model performance.

4.5 Denoising Analysis (RQ4)

We additionally conduct a denoising experiment here, for checking the model robustness under noisy interactions. Specifically, we contaminate the training set by adding a certain proportion of adversarial examples (*i.e.*, 5%, 10%, 15%, 20% negative user-item interactions), meanwhile keeping the validation and testing sets unchanged, following SGL [28]. This experiment could show the model ability of noise-irrelevant representation learning, from an overall perspective. The experimental results on Baby are shown in Figure 6, where the Noisy Ratio means the percentage of noisy interactions added for the model, and the %Drop Percentage represents the percentage of performance degradation.

From the experimental results, we could clearly find that: Although adding noise degrades the model performance of both KGTN and MCCLK, the proposed KGTN is clearly less influenced than the GNN-based and CL-based MCCLK. It is more apparent with a bigger noise ratio, since the performance dropping gap between KGTN and MCCLK becomes larger and larger as the noise ratio increases, which suggests that KGTN is more robust to noisy perturbation.

Metric	Relative Improvement
Item page view per user	+2.3%
Unique visitor list to order	+2.22%
Unique visitor click through rate	+2.3%

Table 5: Results of online A/B testing.

4.6 Online A/B Testing (RQ5)

We conduct online A/B testing in our recommender system in Alibaba, to validate the benefits of KGTTN in the real business scenario. In our online A/B testing, users are randomly divided into A group and B group. When group A browses the website, the system recommends results provided by the previous model, while the group B is recommended with results provided by our KGTTN model. Experiments are run for three weeks, during which both the control and experiment models are trained continuously with new interactions and feedback being used as training data. As shown in Table 5, KGTTN improves item page view (ipv) per user by 2.3%, unique visitor list to order (uv l2o) by 2.22%, and unique visitor click through rate (uv ctr) by 2.3% relatively compared with DMR [12], which is the last version of CTR model in our system. This reveals the practical application value of KGTTN.

5 RELATED WORK

5.1 Knowledge Enhanced Recommendation

Existing Knowledge Enhanced recommendation methods can be roughly categorized into three lines: embedding-based, path-based, and GNN-based methods. **Embedding-based methods** [7, 19, 34] pre-train the KG entity embeddings with knowledge graph embeddings methods (KGE) [1, 10], for enriching item representations, such as CKE [33] and KTUP [2], and RippleNet [18]. **Path-based methods** [16, 31, 35] explore various patterns of connections among items in KG to provide additional guidance for the recommendation, such as PER [32] and MCRec[6]. KPRN [25] further automatically extracts paths between users and items, modeling these paths with RNNs. **GNN-based methods** are founded on the information aggregation mechanism of graph neural networks (GNNs) [5, 30], which integrates multi-hop neighbors into node representations, modeling long-range connectivity. KGCCN [21] and KGNN-LS [20] firstly utilize GNN on KG side, then KGAT[22] propose to utilize GAT on the unified user-item-entity heterogeneous graph. Then CKAN [26] separately models collaborative signals and knowledge signals, and CG-KGR [3] exploits the collaborative signals to guide the aggregation on KG. KGIN [23] models user-item interactions at an intent level, which reveals user intents behind the KG interactions and performs GNN on the user-intent-item-entity graph. More recently, **CL-based Methods** such as MCCLK [37], KGIC [38], and KGCL [29] combine a contrastive learning paradigm with GNN-based methods, and build cross-view contrastive frameworks as additional self-discrimination supervision signals to enhance robustness.

5.2 Multi-intent Modeling

Current multi-intent modeling usually adopts a disentangled representation learning paradigm, which splits the user embedding

into multiple chunks and tries to learn each chunked embedding respectively for representing each user intent. In graph representation learning area, DisenGCN [13], IPGDN [11] and ADGCN [36] adopt such a paradigm and utilize the Hilbert-Schmidt Independence Criterion (HSIC) and adversarial learning for ensuring the effectiveness of intent modeling. As for recommendation scenarios, DGCF [24] proposes to disentangle the user representations and adopts a mutual information restraint for the independence of all intents. And KGIN [23] considers each intent as an attentive combination of KG relations, then use a local aggregation manner for the intent modeling. MIDGN [39] performs fine-grained intent disentanglement from the hierarchical structure in bundle recommendation, together with an intent contrastive mechanism. Despite the success of disentangled learning attempts in previous methods, they commonly learn the multi-intent representations with local information, while ignore the importance of global signals. Hence, our work focuses on learning intent-aware representations with global information, and exploits the intent features to denoise the knowledge information.

6 CONCLUSION

In this paper, we focus on modeling multiple intents with global information, and leveraging intents to denoise the knowledge information. We propose a novel framework, KGTTN, which achieves fine-grained and robust user/item representation learning from two dimensions: 1) KGTTN models global intents and learns intent-aware user/item representations with a proposed graph transformer, by integrating global signals into learnable intents. 2) KGTTN exploits the user intents to sample the relevant knowledge, and designs a local-global contrastive mechanism within the sampled graph, for robust representation learning. Extensive experiments on three public datasets demonstrate that KGTTN significantly improves the recommendation performance over baselines on both Click-Through rate prediction and Top-K recommendation tasks. Furthermore, the online A/B testing on recommender system of Alibaba demonstrates the practical application value of KGTTN.

ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China under Grant No. 62276110, No. 62172039 and in part by the fund of Joint Laboratory of HUST and Pingan Property & Casualty Research (HPL). The authors would also like to thank the anonymous reviewers for their comments on improving the quality of this paper.

REFERENCES

- [1] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Neural Information Processing Systems (NIPS)*. 1–9.
- [2] Yixin Cao, Xiang Wang, Xiangnan He, Zikun Hu, and Tat-Seng Chua. 2019. Unifying knowledge graph learning and recommendation: Towards a better understanding of user preferences. In *WWW*. 151–161.
- [3] Yankai Chen, Yaming Yang, Yujing Wang, Jing Bai, Xiangchen Song, and Irwin King. 2022. Attentive Knowledge-aware Graph Convolutional Networks with Collaborative Guidance for Personalized Recommendation. In *ICDE*. 299–311.
- [4] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*. 249–256.
- [5] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*. 1025–1035.

- [6] Binbin Hu, Chuan Shi, Wayne Xin Zhao, and Philip S Yu. 2018. Leveraging meta-path based context for top-n recommendation with a neural co-attention model. In *SIGKDD*. 1531–1540.
- [7] Jin Huang, Wayne Xin Zhao, Hongjian Dou, Ji-Rong Wen, and Edward Y Chang. 2018. Improving sequential recommendation with knowledge-enhanced memory networks. In *SIGIR*. 505–514.
- [8] Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparametrization with gumble-softmax. In *International Conference on Learning Representations (ICLR 2017)*. OpenReview. net.
- [9] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [10] Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning entity and relation embeddings for knowledge graph completion. In *AAAI*.
- [11] Yanbei Liu, Xiao Wang, Shu Wu, and Zhitao Xiao. 2020. Independence promoted graph disentangled networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 4916–4923.
- [12] Ze Lyu, Yu Dong, Chengfu Huo, and Weijun Ren. 2020. Deep match to rank model for personalized click-through rate prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 156–163.
- [13] Jianxin Ma, Peng Cui, Kun Kuang, Xin Wang, and Wenwu Zhu. 2019. Disentangled graph convolutional networks. In *International conference on machine learning*. PMLR, 4212–4221.
- [14] Xubin Ren, Lianghao Xia, Jiashu Zhao, Dawei Yin, and Chao Huang. 2023. Disentangled Contrastive Collaborative Filtering. *arXiv preprint arXiv:2305.02759* (2023).
- [15] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv* (2012).
- [16] Chuan Shi, Binbin Hu, Wayne Xin Zhao, and S Yu Philip. 2018. Heterogeneous information network embedding for recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2018), 357–370.
- [17] Zhu Sun, Jie Yang, Jie Zhang, Alessandro Bozzon, Long-Kai Huang, and Chi Xu. 2018. Recurrent knowledge graph embedding for effective recommendation. In *RecSys*. 297–305.
- [18] Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo. 2018. Ripplenet: Propagating user preferences on the knowledge graph for recommender systems. In *CIKM*. 417–426.
- [19] Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi Guo. 2018. DKN: Deep knowledge-aware network for news recommendation. In *WWW*. 1835–1844.
- [20] Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and Zhongyuan Wang. 2019. Knowledge-aware graph neural networks with label smoothness regularization for recommender systems. In *SIGKDD*. 968–977.
- [21] Hongwei Wang, Miao Zhao, Xing Xie, Wenjie Li, and Minyi Guo. 2019. Knowledge graph convolutional networks for recommender systems. In *WWW*. 3307–3313.
- [22] Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. 2019. Kgat: Knowledge graph attention network for recommendation. In *SIGKDD*. 950–958.
- [23] Xiang Wang, Tinglin Huang, Dingxian Wang, Yancheng Yuan, Zhenguang Liu, Xiangnan He, and Tat-Seng Chua. 2021. Learning Intents behind Interactions with Knowledge Graph for Recommendation. In *WWW*. 878–887.
- [24] Xiang Wang, Hongye Jin, An Zhang, Xiangnan He, Tong Xu, and Tat-Seng Chua. 2020. Disentangled graph collaborative filtering. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*. 1001–1010.
- [25] Xiang Wang, Dingxian Wang, Canran Xu, Xiangnan He, Yixin Cao, and Tat-Seng Chua. 2019. Explainable reasoning over knowledge graphs for recommendation. In *AAAI*, Vol. 33. 5329–5336.
- [26] Ze Wang, Guangyan Lin, Huobin Tan, Qinghong Chen, and Xiyang Liu. 2020. CKAN: Collaborative Knowledge-aware Attentive Network for Recommender Systems. In *SIGIR*. 219–228.
- [27] Ziyang Wang, Wei Wei, Ding Zou, Yifan Liu, Xiao-Li Li, Xian-Ling Mao, and Minghui Qiu. 2024. Exploring global information for session-based recommendation. *Pattern Recognition* 145 (2024), 109911.
- [28] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *SIGIR*. 726–735.
- [29] Yuhao Yang, Chao Huang, Lianghao Xia, and Chenliang Li. 2022. Knowledge graph contrastive learning for recommendation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1434–1443.
- [30] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *SIGKDD*. 974–983.
- [31] Xiao Yu, Xiang Ren, Quanquan Gu, Yizhou Sun, and Jiawei Han. 2013. Collaborative filtering with entity similarity regularization in heterogeneous information networks. *IJCAI HINA* (2013).
- [32] Xiao Yu, Xiang Ren, Yizhou Sun, Quanquan Gu, Bradley Sturt, Urvasi Khandelwal, Brandon Norick, and Jiawei Han. 2014. Personalized entity recommendation: A heterogeneous information network approach. In *WSDM*. 283–292.
- [33] Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing Xie, and Wei-Ying Ma. 2016. Collaborative knowledge base embedding for recommender systems. In *SIGKDD*. 353–362.
- [34] Yongfeng Zhang, Qingyao Ai, Xu Chen, and Pengfei Wang. 2018. Learning over knowledge-base embeddings for recommendation. *arXiv preprint arXiv:1803.06540* (2018).
- [35] Huan Zhao, Quanming Yao, Jianda Li, Yangqiu Song, and Dik Lun Lee. 2017. Meta-graph based recommendation fusion over heterogeneous information networks. In *SIGKDD*. 635–644.
- [36] Shuai Zheng, Zhenfeng Zhu, Zhizhe Liu, Shuiwang Ji, Jian Cheng, and Yao Zhao. 2021. Adversarial graph disentanglement. *arXiv preprint arXiv:2103.07295* (2021).
- [37] Ding Zou, Wei Wei, Xian-Ling Mao, Ziyang Wang, Minghui Qiu, Feida Zhu, and Xin Cao. 2022. Multi-level cross-view contrastive learning for knowledge-aware recommender system. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1358–1368.
- [38] Ding Zou, Wei Wei, Ziyang Wang, Xian-Ling Mao, Feida Zhu, Rui Fang, and Danyang Chen. 2022. Improving knowledge-aware recommendation with multi-level interactive contrastive learning. In *Proceedings of the 31st ACM international conference on information & knowledge management*. 2817–2826.
- [39] Ding Zou, Sen Zhao, Wei Wei, Xian-ling Mao, Ruixuan Li, Danyang Chen, Rui Fang, and Yuanyan Fu. 2023. Towards Hierarchical Intent Disentanglement for Bundle Recommendation. *IEEE Transactions on Knowledge and Data Engineering* (2023).