

Explainable Session-based Recommendation via Path Reasoning

Yang Cao, Shuo Shang, Jun Wang, and Wei Zhang

Abstract—This paper explores providing explainability for session-based recommendation (SR) by path reasoning. Current SR models emphasize accuracy but lack explainability, while traditional path reasoning prioritizes knowledge graph exploration, ignoring sequential patterns present in the session history. Therefore, we propose a generalized hierarchical reinforcement learning framework for SR, which improves the explainability of existing SR models via Path Reasoning, namely PR4SR. Considering the different importance of items to the session, we design the session-level agent to select the items in the session as the starting point for path reasoning and the path-level agent to perform path reasoning. In particular, we design a multi-target reward mechanism to adapt to the skip behaviors of sequential patterns in SR, and introduce path midpoint reward to enhance the exploration efficiency in knowledge graphs. To improve the completeness of the knowledge graph and to diversify the paths of explanation, we incorporate extracted feature information from images into the knowledge graph. We instantiate PR4SR in five state-of-the-art SR models (i.e., GRU4REC, NARM, GCSAN, SR-GNN, SASRec) and compare it with other explainable SR frameworks, to demonstrate the effectiveness of PR4SR for recommendation and explanation tasks through extensive experiments with these approaches on four datasets.

Index Terms—explainable recommendation, session-based recommendation, hierarchical reinforcement learning, knowledge graph



1 INTRODUCTION

BOTH session-based recommendation (SR) and explainable recommendation have gained great attention in recent years. However, most of the current SR models focus on the problem of how to improve the accuracy of recommendations, while neglecting the process of providing explainability. A large number of research works design different model structures to capture user preference information and model sequential patterns, such as recurrent neural network [1], attention mechanism [2]–[4] and graph neural network [5], [6]. However, none of these structures are explainable. And a series of studies [7]–[14] have shown that an explainable recommendation process can improve the persuasiveness and trustworthiness of recommendation systems. Therefore, some generalizable explainability frameworks are needed to improve the explainability of existing SR models.

Path reasoning [15]–[18] methods provide explainable information for recommendations by exploring paths between entities in the knowledge graph. In order to avoid enumerating possible reasoning paths in large-scale knowledge graphs [9], [19]–[21], recent research work combines reinforcement learning methods with path reasoning to improve the efficiency of exploring paths in knowledge graphs by designing suitable reward functions [22]. However, the current reward function only focuses on finding valid target points in the knowledge graph and exploring the transfer process between neighboring items in a session, but fails

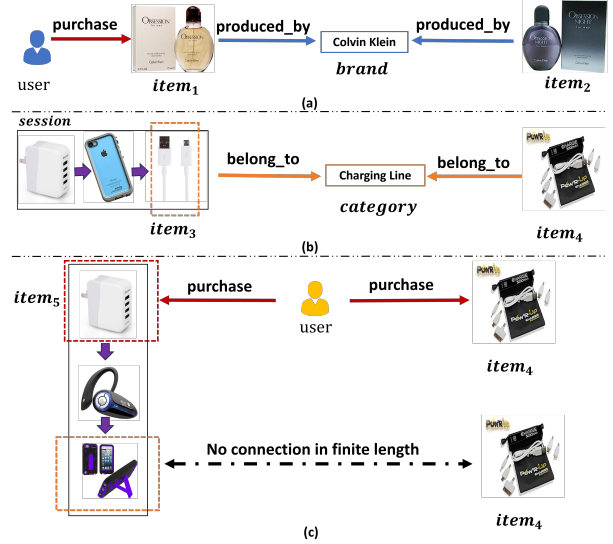


Fig. 1: Three examples of path reasoning: (a) starting with a user, (b) starting with the last item, (c) starting with the most relevant item.

to take into account the skip behaviors of sequential patterns [25] by setting corresponding reward values for them.

Currently path reasoning methods can be divided into two categories. The first part [15], [23] integrates rich heterogeneous information from the knowledge graph into the recommendation process, utilizing a path-level agent learning to navigate from a given user entity to potential items of interest to the user. Fig. 1 (a) illustrates a path of inference, $user \xrightarrow{\text{purchase}} item_1 \xrightarrow{\text{produced_by}} Colvin\ Klein \xleftarrow{\text{produced_by}} item_2$, which indicates that the reason for recommending $item_2$ to the user is that the user has purchased $item_1$

- Yang Cao, Jun Wang and Wei Zhang are with the School of Computer Science and Technology, East China Normal University, Shanghai, China. E-mail: {caoyang99775, wongjun, zhangwei.thu2011}@gmail.com.
- Shuo Shang is with School of Computer Science and Engineering and Shenzhen Institute for Advanced Study University of Electronic Science and Technology of China, Chengdu, China. E-mail: jedi.shang@gmail.com.

before, and the recommended item $item_2$ belongs to the same brand as $item_1$. However, this part of the work does not integrate session information as well as the design of the reward function only considers how to explore efficiently in the knowledge graph. The other part [24] considers combining session information and heterogeneous information in the knowledge graph to provide an explainable framework for SR, which fuses session information and heterogeneous information into the state representation for reinforcement learning and selects the last item in the session as the starting point for the path reasoning, and then plans an explainable path from the interacted item entity to the recommended item entity. For example in Fig. 1 (b), the last item in the session is chosen as the starting point for path reasoning, and a path, $item_3 \xrightarrow{\text{belong_to}} \text{Charging Line} \xleftarrow{\text{belong_to}} item_4$, is inferred from the knowledge graph, which indicates that the reason for predicting $item_4$ as the next is that both $item_4$ and the last $item_3$ in the session belong to the charging line category. Although this type of approach has taken session information into account, the design of the reward function still does not cover the skip behaviors of sequential patterns.

Moreover, some studies of SR have shown that the last item in a session does not represent the entire session sequence of interest. Caser [25] captures two sequence patterns in a sequence with horizontal filters, e.g., for a session [Airport, Hotel, Fast Food, Restaurant], the next behavior may be Greta Wall due to (Airport, Hotel) or Bar due to (Fast Food, Restaurant). NARM [2] finds that the most important items to the session appear at the beginning or in the middle of the session in some cases. These studies also suggest that simply choosing the last item in a session as a starting point for path reasoning is not sufficient to provide an explainable framework for SR. But existing explainable SR frameworks [24] only consider using the last item in the session as the starting point for path reasoning, which may lead to a failure to provide explainable paths for SR. For example in Fig. 1 (c), if the last item is chosen as the starting point of path reasoning, because the last item and $item_4$ are the cell phone holder and cell phone charging cable respectively and the connection between these two items is not strong in the knowledge graph, there does not exist a finite-length path connecting these two items in the knowledge graph; if the first item is chosen, it can be extrapolated in finite-length path extrapolation because the first item $item_5$ is the charging plug and $item_4$ is the charging cable, which are more closely related to each other. There exists such a path in the knowledge graph, $item_5 \xleftarrow{\text{purchase}} user \xrightarrow{\text{purchase}} item_4$, indicating that the exists user purchased both items.

Therefore, although the current research work demonstrates that it is possible to provide an explainable process for recommendation through path reasoning, the current path reasoning approach is not suitable for generalization in the scenario of SR. The main problem lies in the fact that previous research on path reasoning has only considered how to perform path reasoning in knowledge graphs, whereas for session recommendation scenarios, the impact of different items in a session on path reasoning should also be explored, as well as the consideration of designing appropriate reward functions for skip behaviors of sequential

patterns in session scenarios.

In addition, a series of related study [8], [23] demonstrate that the shorter the path length produced by path reasoning, the less time the reasoning takes as well as the more acceptable the explainable paths are. So that under the constraint of finite length, the closer the associations between the entities in the knowledge graph are, the higher the likelihood of exploring reasonable paths through path reasoning. In e-commerce scenarios, the images of products contain some feature information of the products, which can be used as product attributes to supplement the knowledge graph. However, in the current methods of path reasoning to construct knowledge graphs, the image information of products is not taken into account, which leads to less attribute information of products and fewer interconnections in the constructed knowledge graphs.

To address these challenges, we propose a generalized hierarchical reinforcement learning framework for SR, which improves the explainability of existing SR models through Path Reasoning (denoted as PR4SR). We design the session-level agent to select important items from the session as the starting point for path reasoning. We then make path reasoning to provide explainable paths by the proposed path-level agent. In particular, we design a multi-target reward mechanism to adapt to the skip behaviors of sequential patterns in SR, and introduce path midpoint reward to enhance the exploration efficiency in knowledge graphs. In order to make the interconnections of entities in the knowledge graph closer and increases the diversity of explainable paths, we have extracted the product feature information in the image by using image recognition method, and add them to the knowledge graph in the form of entities.

In summary, our contributions lie in three aspects:

- PR4SR is the first path reasoning approach that uses hierarchical reinforcement learning to provide a generalized and explainable framework for SR, where the session-level agent selects the important items in the session as the starting point for path reasoning and the path-level agent performs path reasoning in the knowledge graph.
- PR4SR is the first path reasoning approach that combines the skip behaviors of sequential patterns in SR into the design of the reward mechanism. The reward function is also designed to improve the exploration efficiency by considering the distance from the midpoint of the path to the goal points.
- Comparing with the method of constructing product knowledge graphs in traditional path reasoning, a new method of constructing product knowledge graphs incorporating product picture features is designed, which improves the correlation between entities and the diversity of explainable path forms.
- PR4SR generalizes well and can be combined with existing unexplainable SR models to accomplish both recommendation and explainability tasks at the same time. We compare PR4SR with state-of-the-art methods on four public datasets. The results show that PR4SR improves recommendation accuracy and model explainability for unexplainable SR models.

We also conduct a comprehensive ablation study to analyze the contributions of key components. We further illustrate that the explainability of PR4SR outperforms other explainable SR frameworks through a user survey.

2 RELATED WORK

2.1 Session-based recommendation

Early work has used MF [26] to learn users’ overall interest preferences and MC [27] to learn the dependencies between items. FPMC [28] introduces a personalized transfer matrix based on Markov chains, which can capture both temporal information and long-term user preference information, and a Matrix Decomposition model is introduced to solve the sparsity problem of the transfer matrix.

The core idea of this type of models is to mine the potential relationship between users and items by decomposing or completing the user-item scoring matrix for personalized recommendation. Their workings are not easily understood by the general user and lack intuitiveness and explainability.

Deep neural network techniques are now being applied to SR models. In order to capture the user’s interest at different points in time, GRU4Rec [1] captures the interest information through the GRU [29] layer to improve the accuracy of recommendations. To combine the user’s interest in the current session and the sequence information of the session, NARM [2] proposes a hybrid state encoder with an attention mechanism. Focusing on the fact that previous models do not consider the impact of the user’s current action on the next action, STAMP [4] proposes a short-term attention/memory priority model. SASRec [3] finds that MC performs better in sparse datasets while RNN [30] performs better in dense datasets, thus proposing a self-attention based sequential model to balance MC and RNN. SR-GNN [6] argues that previous research approaches ignore complex transformations between items, thus uses GNN [31] to capture complex transformations of items in graph-structured data. GCSAN [5] also focuses on the importance of dependencies between items and proposes a graph contextualized self-attention model.

Although these models improve the accuracy of recommendations by designing different network structures, the study of model explainability has been neglected, which may reduce the user’s satisfaction with the whole recommender system.

2.2 Path reasoning approaches

Path reasoning approaches explicitly model the paths between entities over a knowledge graph for recommendation, and improve the efficiency of exploration in the knowledge graph through reinforcement learning methods. Such current work on path reasoning combining reinforcement learning and knowledge graphs can be categorized into two groups.

One of them uses the representation of the knowledge graph as the state representation for reinforcement learning and the user entity as a starting point for path reasoning. PGPR [23] first emphasizes the importance of incorporating knowledge graphs into recommendations to provide

reasoning and explainable paths for recommendations, and proposes for the first time reasoning of explainable paths in the knowledge graph by means of reinforcement learning. ADAC [8] extracts path demonstrations from the knowledge graph as guided paths for reinforcement learning exploration, which improves the speed of model fitting. TPRec [33] uses GMM [38] to cluster the purchase relation and incorporates the time information into the reward function. SENTIMENT [7] observes that previous works do not consider the semantic information of relations in the knowledge graph, thus constructs more fine-grained types of relations by extracting comment information. Multi-level [15] introduces external knowledge to construct a multi-level knowledge graph and designs a Cascading Actor-Critic [39], [40] that uses a top-down strategy to search the space of the knowledge graph. The second type of work is to combine the session information and heterogeneous information in the knowledge graph and use the item entity as the starting point. This type of method is represented by REKS [24], which selects the last item in the session as the starting point for path reasoning and incorporates session information into the state representation.

Although REKS [21] incorporates sequential information to the state representation of path reasoning, the task of REKS is essentially in predicting possible paths between neighboring items, and does not take into account skip behaviors of sequential patterns present in the session when designing the reward function, and also selects only the last item in the session as the starting point for path reasoning. However, [2], [25] point out that skip behaviors of sequential patterns are present in the session and that items with stronger correlation with the predicted goal with not necessarily appear at the end of the session.

Therefore, in order to satisfy the properties of the SR scenario, we attempt to provide a generalized explainable framework for unexplainable SR models in the form of path reasoning using a hierarchical reinforcement learning approach combined with knowledge graphs. Specifically, the reward function consists of two components: multi-target reward, which provides feedback for skip behaviors of sequential patterns that may occur in the session; and path-midpoint reward, a mechanism designed to improve the efficiency of exploration in the knowledge graph and ensure that more valuable information can be discovered. And hierarchical reinforcement learning has two levels: the session-level agent selects important items in the session as the starting point for path reasoning and then the path-level agent performs path reasoning in the knowledge graph, where important items refer to items that are highly relevant to the predicted items.

3 PRELIMINARY

Problem Formulation Let \mathcal{V} denote the set of items and $[v_1, v_2, \dots, v_n]$ denotes session history, where $v_i \in \mathcal{V} (1 \leq i \leq n)$ and n denotes the length of the session. The SR model relies on the items present in the session to understand the user’s interests and preferences and then make accurate predictions. The model generates a prediction score \mathbf{y} for the items in the candidate pool, where $\mathbf{y} = [y_1, y_2, \dots, y_m]$ and m denotes the size of the candidate pool, and then

TABLE 1: Important notations

Symbol	Description
\mathcal{E}	The set of total entities.
\mathcal{R}	The set of total relations.
\mathcal{G}	The knowledge graph.
S_{KG^x}	Representation of entities in the knowledge graph.
S_{se}	Representations of sessions obtained through unexplainable SR models.
$Item_{Long}$	Relatively important item selected from session.
$Item_{Short}$	The last item in the session.
S_{long}^t	State representation of an exploration starting with $Item_{Long}$.
S_{short}^t	State representation of an exploration starting with $Item_{Short}$.
P^{se}	Probability of selecting item from session as $Item_{Long}$.
P^L	Starting with $Item_{Long}$, the probability of choosing the next (relation,entity).
P^S	Starting with $Item_{Short}$, the probability of choosing the next (relation,entity).
T	Select T consecutive prediction target items as the target items for path reasoning.
L^{Ce}	The cross-entropy loss function.
L^{Path}	Loss function for Path-level agent.
L^{Se}	Loss function for Session-level agent.
G_t	The discounted cumulative rewards at time step t.
K	The length of recommendation list.
$ A^{path} $	Size of the action space of the Path-level agent.

selects the top k items sorted by score $y_j(1 \leq j \leq m)$ as recommendations.

In our task, based on the SR model and the pre-trained knowledge graph representation, the session-level agent selects the starting point for path reasoning from the session, and the path-level agent reasons about explainable paths from the knowledge graph. We design this hierarchical reinforcement learning framework to provide explainable paths for SR models while improving recommendation accuracy.

Knowledge Graphs Let \mathcal{E} and \mathcal{R} represent the entity set and relation set respectively. A knowledge graph \mathcal{G} is defined as $\mathcal{G} = \{(e_h, r, e_t) | e_h, e_t \in \mathcal{E}, r \in \mathcal{R}\}$, where each triplet (e_h, r, e_t) represents a relation r from head entity e_h to tail e_t . Our work builds a knowledge graph in conjunction with recommendation scenarios, containing multiple types of entities such as item, brand, category, etc., and also multiple types of relations such as purchase, belong_to, produced_by, also_viewed, etc. A part of the relations and entities in the knowledge graph comes from the interaction between the user and the item. For example, $user_1 \xrightarrow{\text{purchase}} item_1$ indicates that $user_1$ has purchased $item_1$. Another part of the relations and entities represent some basic attributes of the item. For example, $item_1 \xrightarrow{\text{image_sim}} pearl$ indicates that $item_1$ has a pearl element in the image.

Hierarchical Reinforcement Learning Let \mathcal{S} , \mathcal{A} , and \mathcal{R} represent the state space, action space, and reward function respectively. The hierarchical reinforcement learning (HRL) problem can be formulated as a series of Markov decision processes (MDPs) $\{MDP_1, MDP_2, \dots\}$, where each markov decision process corresponds to a level or task in the hierarchy. In our task scenario, MDP_1 represents the session-level agent selecting the appropriate item from the session as the starting point for path reasoning, and MDP_2 represents the path-level agent performing path exploration in the knowledge graph.

Explainable Paths Explainable paths contain information from both session and knowledge graph. For a session: $[v_1, v_2, v_3]$, $v_i \in \mathcal{V}(1 \leq i \leq 3)$, the session-level agent chooses v_2 as the starting point of path reasoning, v_2 corresponds to the entity in the knowledge graph \mathcal{G} as $item_2$, and the path-level agent starts from $item_2$ and eventually derives multiple explainable paths, e.g., $[v_1, v_2, v_3] \xrightarrow{\text{session-level agent}} v_2 \xrightarrow{\mathcal{G}} item_2 \xrightarrow{\text{belong_to}} category_1 \xleftarrow{\text{belong_to}} item_4$.

Table 1 summarizes the main symbols.

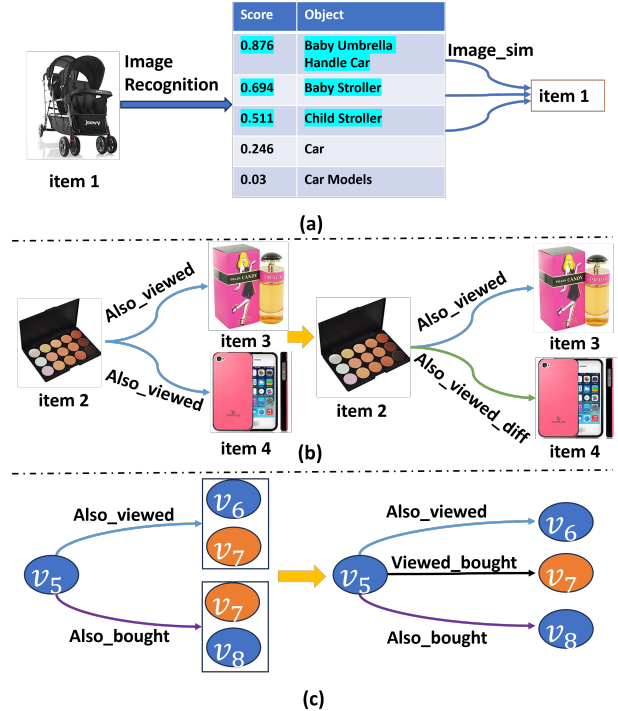


Fig. 2: Details of knowledge graph construction: (a) extract image features; (b) split relations based on different product domains; (c) merge duplicate entities.

4 KNOWLEDGE GRAPH CONSTRUCTION

As shown in Fig. 2, unlike the previous path reasoning method of constructing knowledge graphs [23], [24], we improve the completeness and clarity of the constructed e-commerce knowledge graph in the following three ways: (a) supplementing the feature information in the product images into the knowledge graph; (b) dividing the types of relations according to the entities in different product domains; (c) merging duplicate entities.

(a) Extract Image Feature As shown in Fig. 2 (a), we extract the top-5 feature information in each image through the image recognition API provided by Baidu Intelligent Cloud¹, and select the feature information with confidence score greater than 50% to be added to the knowledge graph, which connects to the items through *image_sim* relation. This part of information can supplement the missing feature information of the items, so that items with the same feature

¹ <https://cloud.baidu.com/product/imagerecognition>

TABLE 2: Number of the relations in the Amazon datasets.

Relation	Description	Beauty	Cellphones	Baby
purchase	$user \xrightarrow{\text{purchase}} product$	163,678	137,832	142,980
produced_by	$product \xrightarrow{\text{produced_by}} brand$	19,356	10,222	9,596
belong_to	$product \xrightarrow{\text{belong_to}} category$	95,832	63,178	13,720
image_sim	$product \xrightarrow{\text{image_sim}} image_feature$	414,678	330,754	359,484
title_sim	$product \xrightarrow{\text{title_sim}} title_feature$	1,567,156	1,989,166	967,518
also_bought	$product \xrightarrow{\text{also_bought}} product$	279,686	438,812	313,738
also_viewed	$product \xrightarrow{\text{also_viewed}} product$	115,218	17,402	121,936
viewed_bought	$product \xrightarrow{\text{viewed_bought}} product$	179,572	5,190	94,698
bought_together	$product \xrightarrow{\text{bought_together}} product$	17,396	16,078	10,154
also_bought_diff	$product \xrightarrow{\text{also_bought_diff}} related_product$	830,874	676,432	523,780
also_viewed_diff	$product \xrightarrow{\text{also_viewed_diff}} related_product$	554,444	61,152	378,040
bought_viewed_diff	$product \xrightarrow{\text{bought_viewed_diff}} related_product$	441,260	7,394	128,840
bought_together_diff	$product \xrightarrow{\text{bought_together_diff}} related_product$	11,072	8,018	8,956
co_occur	$product \xrightarrow{\text{co_occur}} product$	23,374	17,614	26,281

TABLE 3: Statistics of the relations on the Douban-movie.

Relation	Description	Relation number
belong_to	$movie \xrightarrow{\text{belong_to}} genre$	41,802
directed_by	$movie \xrightarrow{\text{directed_by}} director$	17,074
acted_by	$movie \xrightarrow{\text{acted_by}} actor$	95,522
described_as	$movie \xrightarrow{\text{described_as}} tag$	148,166
produced_by	$movie \xrightarrow{\text{produced_by}} region$	23,500

information can be connected, increasing the diversity of explainable paths.

(b) Split Relations As shown in Fig. 2 (b), we observe that there are different types in the items connected through *also_viewed* or *also_bought*, e.g., *item 2* and *item 3* belong to the cosmetic category from the Amazon-beauty dataset, while *item 4* is the cellphone case from the Amazon-cellphones dataset. We divide the original *also_viewed* and *also_bought* relations into *also_viewed* and *also_viewed_diff*, and *also_bought* and *also_bought_diff*, depending on whether the connected items belong to the same product domain, so that the relations in the knowledge graph can represent more specific meanings.

(c) Merge Duplicate Entities. Taking the Amazon-beauty dataset as an example, as shown in Fig. 2 (c), we find that the knowledge graph built up following the traditional approach has 51.43% of the entities connected via *also_viewed* also appearing in the entities connected via *also_bought*, and similarly, 39.19% of the nodes connected via *also_bought* also appearing in the entities connected via *also_viewed*. Thus, we take entities like v_7 that are connected by both *also_viewed* and *also_bought* and connect them via *viewed_bought*. This approach reduces the number of connected relations by about 300,000 in the knowledge graph and improves the efficiency of path reasoning.

In the e-commerce scenario, we classify entities into seven categories: user, product, brand, category, image_feature, title_feature, related product, where related product indicates that the products do not belong to the same broad category, such as Amazon-Beauty and Amazon-

Baby. And fourteen different types of relations. For the movie recommendation scenario, we classify entities into six categories: movie, genre, director, actor, tag, and region. And five different types of relations. The detailed data is summarized in Table 2, Table 3, Table 4 and Table 5.

5 PROPOSED FRAMEWORKS

The overview of the proposed framework PR4SR is illustrated in Fig. 3. PR4SR first utilizes the Session Encoder to encode sequence information as the state information of the session-level agent, which selects the appropriate item from the session history as the starting point for path reasoning, labeled as $Item_{Long}$. Subsequently, the last item in the session is selected as the other starting point for path reasoning, namely $Item_{Short}$. Next, the path-level agent integrates the sequence information encoded by the Session Encoder as well as the node information in the knowledge graph encoded by the KG encoder. The path-level agent explores the knowledge graph using $Item_{Long}$ and $Item_{Short}$ as starting points to determine the final prediction paths, respectively. The item located at the path’s endpoint serves as the predicted target, representing the sequence of items that the user is expected to interact with over the next T steps. In the rest of this section, we elaborate on the details of the framework.

5.1 Mixed state encoder

In order to capture the short-term interaction features between items, as well as the sequential information features in a session, we use a knowledge graph-level state encoder and a session-level encoder.

Knowledge Graph Level State Encoder We obtain the entity representation through a graph-based pre-training approach (termed as a static Knowledge Graph Level State Encoder (KGLSE)) [41].

$$S_{KG^x} = \text{KGLSE}(e_x) \quad (1)$$

where e_x is a certain entity in the knowledge graph.

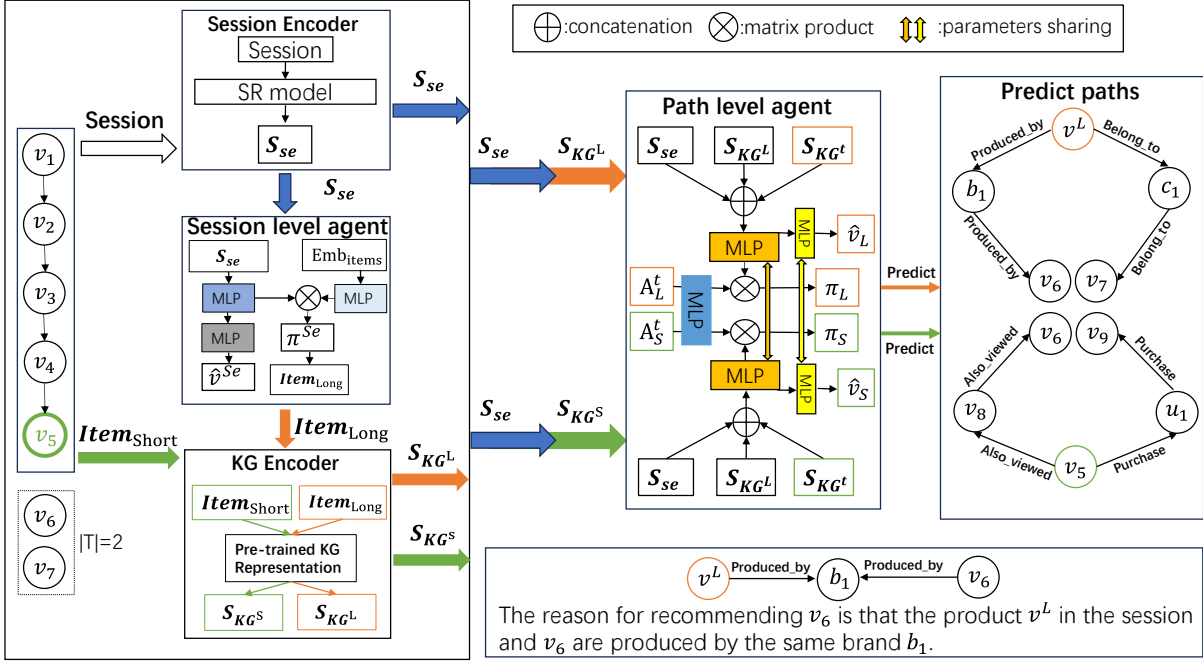


Fig. 3: The overview of the proposed framework.

Session-Level State Encoder Referring to the work of REKS [24], we use the SR model as a Session Level State Encoder (SLSE).

$$S_{se} = \text{SLSE}(\text{session}) \quad (2)$$

where $\text{session} = \{v_1, v_2, \dots, v_n\}, v_i \in \mathcal{V} (1 \leq i \leq n)$.

5.2 Hierarchical Reinforcement Learning

The key to our task is to provide an explainable recommendation process for SR, which requires selecting the important items in the session and performing path reasoning in the knowledge graph. However, the action space for selecting important items is the items within the session, whereas the action space for path reasoning encompasses the combinations of relations and entities within the knowledge graph. Therefore a single agent structure cannot solve our problem and we design the session-level agent and path-level agent.

In the process of path reasoning, the session-level agent selects the essential items from the session as the starting point of path reasoning. Considering the average length of the session, we choose an item from them as the starting point of the path, denoted as $Item_{Long}$. Meanwhile, the last item in the session is the most recent interaction behavior [4], [24], and we also choose it as the starting point for path reasoning, denoted as $Item_{Short}$.

$Item_{Long}$ and $Item_{Short}$ are passed through Eq.1 to get the state information S_{KG^L} and S_{KG^S} . The path-level agent performs path reasoning to explain the recommendation process starting from $Item_{Long}$ and $Item_{Short}$ respectively.

5.3 Definition of the MDP for PR4SR

State. For the session-level agent, the initial state is represented as S_{se} by Eq. 2. For the path-level agent, the initial state is obtained by Eq. 3 and Eq. 4.

$$S_{long}^t = S_{se} \oplus S_{KG^L} \oplus S_{KG^t} \quad (3)$$

$$S_{short}^t = S_{se} \oplus S_{KG^S} \oplus S_{KG^t} \quad (4)$$

where \oplus represents the aggregation of vectors, and S_{KG^t} denotes the entity representation of the current node in the path at time t .

Action. For session-level agent, an action space $A_{se} = \{v_i | v_i \in \text{session}\}$, which denotes the selection of an item from the current session as starting point for path reasoning. For path-level agent, at time t , an action space $A_{path}^t = \{(r_{t+1}, e_{t+1}) | (e_t, r_{t+1}, e_{t+1}) \in \mathcal{G}\}$, which represents the selection of a pair (r_{t+1}, e_{t+1}) to be used at the next step.

$$P^{se} = \text{softmax}((W_1 A_{se}) \otimes (W_2 S_{se})) \quad (5)$$

$$P^L = \text{softmax}((W_3 A_L^t) \otimes (W_4 S_{long}^t)) \quad (6)$$

$$P^S = \text{softmax}((W_3 A_S^t) \otimes (W_4 S_{short}^t)) \quad (7)$$

where W_1, W_2, W_3 and W_4 are trainable parameter matrices and \otimes is for matrix product. A_L^t and A_S^t denote the set of actions starting with $Item_{Long}$ and $Item_{Short}$ respectively. It is noteworthy that P^L and P^S share the same parameters.

Transition Function. For the session-level agent, the transition function is expressed as $p[s_{t+1}^{se} | s_t^{se}, a_t^{se} = v]$, where $v \in \mathcal{V}$. For the path-level agent, the transition function is expressed as $p[s_{t+1}^{path} | s_t^{path}, a_t^{path} = (r_{t+1}, e_{t+1})]$, where $(r_{t+1}, e_{t+1}) \in \mathcal{G}$.

Reward. We designed two types of rewards.

1) Multi-Target Reward Previous studies [2], [25] have shown the existence of skip behaviors of sequential patterns in session history. As shown in Fig. 4, [25] found that there are strong connectivity relationships in session history such as $v_3 \rightarrow v_6$ or $v_5 \rightarrow v_7$ with multiple intermediate items in between, and [2] also proved, from a visualization

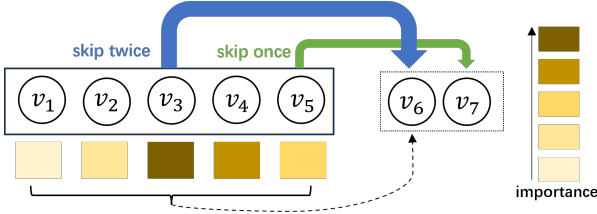


Fig. 4: An example explaining the design of Multi-Target Reward.

perspective, that the item that is most relevant to v_6 does not necessarily exist at the end of the session history. Prior work [24] has only rewarded the $v_5 \rightarrow v_6$ explanatory path with positive feedback, but in fact, the explored $v_3 \rightarrow v_6$ or $v_5 \rightarrow v_7$ paths are also consistent with the results of the session recommendation model and are explainable paths. To ensure that valuable explainable paths explored can be given positive feedback, we set up successive T targets, which are stored in T_{list} in order.

$$R_{Multi} = \begin{cases} T - \text{Index}(v_{end}) & \text{if } v_{end} \text{ in } T_{list}, \\ \log(S_{KG^{end}} * S_{KG^{tar_0}}) & \text{otherwise.} \end{cases} \quad (8)$$

where v_{end} denotes the item that appears at the end of the predicted path and $\text{Index}(v_{end})$ denotes the position of v_{end} in T_{list} , $0 \leq \text{Index}(v_{end}) < T$. tar_0 denotes the 0th item in T_{list} . $S_{KG^{end}}$ and $S_{KG^{tar_0}}$ denote the state representation by Eq. 1. The reward function indicates that: if the v_{end} is in T_{list} , $R_{Multi} = T - \text{Index}(v_{end})$; otherwise, we set the reward as the similarity between the target product tar_0 and the predicted v_{end} .

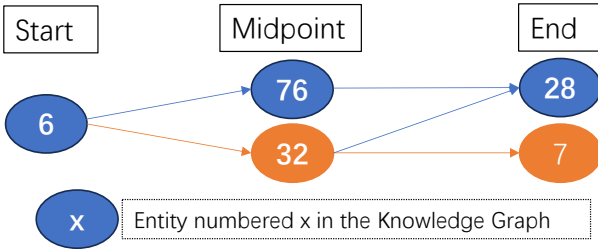


Fig. 5: An example explaining the design of Path Midpoint Reward. Entity 7 is the goal point.

2) Path Midpoint Reward As shown in Fig. 5, we find that there is a problem in relying only on the end-point of the path to calculate the reward. For example, although the path “6→32→28” and the path “6→76→28” are both incorrect paths, the path “6→32→28” has the node of 32, which is nearer to the goal point 7. We store the midpoint from the start point to the goal point in the dictionary dic_{path} , $dic_{path}[(start\ point, end\ point)] = [(r_1, e_1), \dots, (r_c, e_c)], (r_i, e_i) \in \mathcal{G}$, where c indicates that there are several different middle (r_i, e_i) . Then we define path midpoint reward as follows:

$$R_{PMid} = \max \left(0, \max_{i \in \text{range}(T)} [T - i \mid (r, e) \in dic_{path}^i] \right) \quad (9)$$

where T denotes T successive targets, $dic_{path}^i = dic_{path}[(start\ point, tar_i)]$, $tar_i \in T_{list}$. (r, e) indicates the current action.

Optimization. Our loss function consists of three components: session-level state encoder loss, path-level agent loss, and session-level agent loss.

$$L = L^{Ce} + \alpha L^{Path} + \beta L^{Se} \quad (10)$$

where α and β are hyper-parameters to balance the three parts of loss.

The path-level agent’s loss function is defined as:

$$G_t = \sum_{i=t+1}^M \gamma^{i-t-1} R_i \quad (11)$$

$L_{\pi}^{Path}(\theta_{\pi}^{path}) = E_{\pi}[(\hat{v}^{path}(s_t) - G_t) \ln \pi(a_t | s_t, \theta_{\pi}^{path})]$ (12) where G_t represents the discounted cumulative rewards at time step t , M denotes the number of iteration steps, γ indicates the discount factor, $R_i = R_{PMid}$ or R_{Multi} . $\hat{v}^{path}(s_t)$ represents the estimate of the state value function and θ_{π}^{path} denotes the set of path level agent parameters.

The loss function of the session-level agent is defined as:

$$L_{\pi}^{Se}(\theta_{\pi}^{se}) = E_{\pi}[(\hat{v}^{se}(s) - G_{M-1}) \ln \pi(a | s, \theta_{\pi}^{se})] \quad (13)$$

where θ_{π}^{se} denotes the set of session level agent parameters and $\hat{v}^{se}(s)$ represents the estimate of the state value function.

The loss function for the session-level state encoder is defined as:

$$L^{Ce} = - \sum (y_j \log(\hat{y}_j) + (1 - y_j) \log(1 - \hat{y}_j)) \quad (14)$$

where y_j is the ground-truth score of item v_j and \hat{y}_j is the predicted score.

Algorithm 1 summarizes the process of training our model. After training the model, a probabilistic beam search algorithm [8] is used to derive the process of recommending items and provide explainable paths.

6 EXPERIMENTS

We conduct experiments on four real-world datasets to study RP4SR. We aim to answer the following research questions:

RQ1: How does PR4SR improve the performance of the session-based recommendation methods and the well-performed explainable framework?

RQ2: How do different components of PR4SR affect the model performance?

RQ3: How do different hyper-parameter settings affect the performance of PR4SR?

RQ4: How well does PR4SR perform on the explanation task compared to the existing explainable session-based framework?

6.1 Experimental Setup

Datasets. We consider the four real-world datasets, including three datasets Beauty, Cellphones, Baby from the Amazon e-commerce platform [42], and douban-movie from Douban². We extract multiple types of relations and entities from the dataset. We consider a user’s interactions that

² <https://movie.douban.com/>

Algorithm 1 Overview of the PR4SR Framework

```

1: Pre-training knowledge graph
2: Initialize model parameters  $\theta_\pi$  and  $\theta_{\hat{v}}$  of each agent and
   parameters of the session-level state encoder
3: for each session do
4:   Generate  $S_{se}$  by Eq. 2
5:   Select the  $Item_{Long}$  as  $e_0$  according to  $\pi^{se}$ 
6:   done=True, t=0
7:   while done do
8:     Generate  $S_{long}^t$  by Eq. 1 and action Space  $A_L^t$ 
9:     Select an action  $(r_{t+1}, e_{t+1})$  according to  $\pi_L$ 
10:    Get reward  $R_t = R_{PMid}$  or  $R_{Multi}$ .
11:    Store  $(e_t, R_t, (r_{t+1}, e_{t+1}))$  into path
12:    t=t+1
13:    if t==2 then
14:      done=False
15:    end if
16:  end while
17:  Select the  $Item_{Short}$  as  $e_0$ 
18:  done=True, t=0
19:  while done do
20:    Generate  $S_{short}^t$  by Eq. 1 and action Space  $A_S^t$ 
21:    Select an action  $(r_{t+1}, e_{t+1})$  according to  $\pi_S$ 
22:    Get reward  $R_t = R_{PMid}$  or  $R_{Multi}$ .
23:    Store  $(e_t, R_t, (r_{t+1}, e_{t+1}))$  into path
24:    t=t+1
25:    if t==2 then
26:      done=False
27:    end if
28:  end while
29:  Compute  $L^{Se}(\theta_{\hat{v}}^{se}, \theta_\pi^{se}), L^{Path}(\theta_{\hat{v}}^{path}, \theta_\pi^{path})$  and  $L^{Ce}$ 
30:  Update parameters with Adam Gradient Descent
31: end for

```

occurred in one day as a session, and filter out items with less than 5 interactions and sessions with lengths smaller than 2. We randomly sampled 75% of the session data as the training set, 10% of the data as the validation set, and 15% of the data as the test set. Detailed data are described in Table 4, Table 5 and Table 6.

TABLE 4: Number of the entities in the Amazon datasets.

Entity	Beauty	Cellphones	Baby
user	15,438	17,933	13,655
product	11,673	9,805	6,860
brand	2,008	904	716
category	238	107	1
image_feature	2,677	1,256	2,136
title_feature	11,647	10,202	6,270
related_product	160,281	96,674	68,168

TABLE 5: Number of Entities in the Douban-movie dataset.

entity	moive	genre	director	actor	tag	region
numbers	19,114	40	6,975	34,497	16,761	373

Metrics. Regarding the recommendation task, We chose two metrics, HR@k (Hit Ratio), NDCG@k (Normalized Discounted Cumulative Gain), and the values of k are 5, 10, and 20. HR@k measures the proportion of recommendation lists where at least one relevant item in the top-k list

TABLE 6: Statistics of Amazon and Douban-movie datasets.

Dataset	Beauty	Cellphones	Baby	Douban
#entities	204,007	136,811	97,851	77,760
#relations	4,566,296	3,779,244	3,099,721	326,064
#sessions	198,502	194,439	160,792	216,992
average length	2.96	2.75	2.014	2.17

TABLE 7: Hyper-parameter Settings.

Methasod	Dataset	learning rate	α	β
PR4SR_GRU4Rec	Beauty	0.0001	0.005	0.005
	Cellphones	0.0001	0.05	0.01
	Baby	0.00005	0.01	0.005
	Douban	0.001	0.01	0.005
PR4SR_NARM	Beauty	0.0001	0.01	0.005
	Cellphones	0.0001	0.005	0.01
	Baby	0.0005	0.01	0.005
	Douban	0.001	0.005	0.0075
PR4SR_GCSAN	Beauty	0.0001	0.01	0.005
	Cellphones	0.0001	0.005	0.01
	Baby	0.0005	0.005	0.005
	Douban	0.001	0.01	0.0075
PR4SR_SR-GNN	Beauty	0.0001	0.01	0.005
	Cellphones	0.0001	0.005	0.01
	Baby	0.0005	0.005	0.005
	Douban	0.001	0.01	0.0075
PR4SR_SASRec	Beauty	0.0001	0.01	0.005
	Cellphones	0.0001	0.05	0.01
	Baby	0.0005	0.01	0.005
	Douban	0.001	0.05	0.01

is included. NDCG@k takes into account the position of the correctly recommended item in the top-k list, the higher the position, the higher the score. A larger value indicates better performance for both metrics. Regarding the explanation task, we use a questionnaire approach.

Baselines. We apply PR4SR on five representative and widely-used SR baselines, i.e., one RNN-based method GRU4REC [1], one hybrid method that combines attention and RNN NARM [2], two GNN-based models GCSAN [5] and SR-GNN [6], and one attention-based method SAS-Rec [3]. An explainable SR framework REKS³ [24].

Reproducibility Settings. We followed the original settings suggested by the authors [24] to train baseline model (i.e., GRU4Rec, NARM, GCSAN, SR-GNN) on Amazon dataset. For the recommendation process, the path length is set to 2, the maximum action space is set to 200, and the action dropout is set to 0.7. For the reinforcement learning design, the discount factor γ is set to 0.99, and $W_1 \in R^{400 \times 400}$, $W_2 \in R^{400 \times 400}$, $W_3 \in R^{400 \times 400}$, $W_4 \in R^{800 \times 400}$. For the Amazon dataset, Beauty, Cellphones, and Baby, our model is trained for 150 epochs using Adam optimization. For the Douban dataset, our model is trained for 50 epochs using Adam optimization. The batch size is set to 256, and T is set to 5. In the test phase, the sample sizes of the two steps are set to 100 and 1 respectively. More details in Table 7.

6.2 Results and Analysis

RQ1: Overall Results.

The experimental results with all baseline methods with REKS or with our method PR4SR are illustrated in Table 8, where NONE denotes no use of explainable frameworks and “REKS-Avg Improv.” denotes the average improvement

3. To the best of our knowledge, REKS is the only explainable SR framework before.

TABLE 8: Overall performance comparison on the four datasets.

Dataset	Metric	GRU4Rec			NARM			GCSAN			SR-GNN			SASRec			REKS-Avg. Improv.	NONE-Avg. Improv.
		NONE	REKS	PR4SR	NONE	REKS	PR4SR	NONE	REKS	PR4SR	NONE	REKS	PR4SR	NONE	REKS	PR4SR		
Beauty	HR@5	8.70	9.96	10.98	9.48	11.14	12.57	7.87	9.19	9.56	9.43	9.78	10.41	9.00	11.38	12.95	9.46%	26.92%
	HR@10	10.98	13.57	15.41	11.88	14.82	17.01	10.49	12.88	14.02	11.62	12.89	15.41	14.11	15.34	17.26	13.84%	34.40%
	HR@20	14.00	17.62	20.13	14.67	19.11	21.93	13.09	17.30	19.91	14.21	17.10	21.26	19.56	20.00	22.16	15.83%	41.64%
	NDCG@5	6.49	6.82	7.55	7.12	7.82	8.85	5.76	6.19	6.46	6.92	6.82	7.08	5.97	7.91	9.07	9.36%	21.41%
	NDCG@10	7.23	7.99	8.98	7.89	9.00	10.28	6.61	7.36	7.90	7.63	7.82	8.70	7.61	9.19	10.46	11.79%	25.08%
	NDCG@20	7.99	9.01	10.17	8.59	10.09	11.52	7.26	8.48	9.40	8.28	8.89	10.17	8.99	10.36	11.77	13.19%	28.93%
Cellphones	HR@5	7.22	7.06	9.68	8.16	7.77	9.97	6.85	8.08	9.16	7.16	6.95	9.39	6.12	8.11	10.11	27.70%	37.27%
	HR@10	9.62	10.76	13.72	10.98	10.94	14.16	9.48	11.54	13.09	9.87	10.21	13.50	11.07	11.10	14.24	26.15%	35.00%
	HR@20	12.91	15.70	18.70	14.30	15.25	19.62	12.81	16.12	18.16	13.23	13.82	19.03	17.59	15.26	20.02	25.86%	36.29%
	NDCG@5	5.10	4.51	6.47	5.78	5.20	6.76	4.66	5.19	6.17	5.06	4.63	6.29	3.58	5.43	6.65	30.14%	37.28%
	NDCG@10	5.87	5.70	7.78	6.71	6.21	8.11	5.50	6.31	7.43	5.95	5.67	7.62	5.16	6.40	8.01	28.84%	34.33%
	NDCG@20	6.70	6.95	9.04	7.55	7.30	9.48	6.34	7.47	8.70	6.79	6.58	9.01	6.80	7.45	9.39	27.91%	33.69%
Baby	HR@5	4.83	5.17	5.83	5.15	5.67	6.05	3.56	5.59	5.72	3.90	5.37	5.47	5.01	5.67	6.14	6.34%	32.30%
	HR@10	7.41	7.53	8.47	7.37	8.05	8.86	5.97	7.97	8.39	6.22	7.90	8.28	7.22	8.09	8.91	8.54%	26.30%
	HR@20	10.92	10.87	12.25	11.01	11.06	12.71	9.63	11.04	12.24	9.31	10.89	11.93	10.79	11.50	12.32	11.06%	19.42%
	NDCG@5	2.86	3.47	4.07	3.25	3.94	4.14	2.43	3.82	3.91	2.44	3.57	3.72	3.34	3.93	4.15	6.99%	41.53%
	NDCG@10	3.69	4.23	4.91	3.96	4.71	5.05	3.01	4.59	4.76	3.18	4.38	4.62	4.04	4.71	5.05	7.96%	37.78%
	NDCG@20	4.57	5.07	5.87	4.89	5.46	6.02	3.92	5.36	5.73	3.96	5.14	5.54	4.94	5.56	5.90	9.36%	31.39%
Douban	HR@5	3.26	4.69	4.94	3.04	4.98	5.32	2.80	5.07	5.07	3.40	4.58	4.89	4.00	4.87	5.48	6.33%	57.76%
	HR@10	6.32	7.10	7.38	6.78	7.25	8.28	6.00	7.35	7.62	6.81	7.04	7.53	7.23	7.05	8.55	10.00%	18.96%
	HR@20	8.97	9.48	10.68	10.53	9.17	11.66	8.33	8.97	10.56	10.78	9.46	11.11	11.17	8.94	11.80	21.42%	13.07%
	NDCG@5	1.81	2.69	3.01	1.71	2.88	3.31	1.52	3.22	3.14	1.92	2.71	3.04	2.36	2.84	3.30	10.60%	72.97%
	NDCG@10	2.72	3.41	3.78	2.84	3.57	4.26	2.52	3.92	3.93	2.97	3.46	3.85	3.35	3.49	4.27	12.83%	40.46%
	NDCG@20	3.09	3.60	4.43	3.61	3.55	4.86	2.76	3.71	4.47	3.61	3.88	4.60	4.17	3.73	4.93	26.23%	37.20%

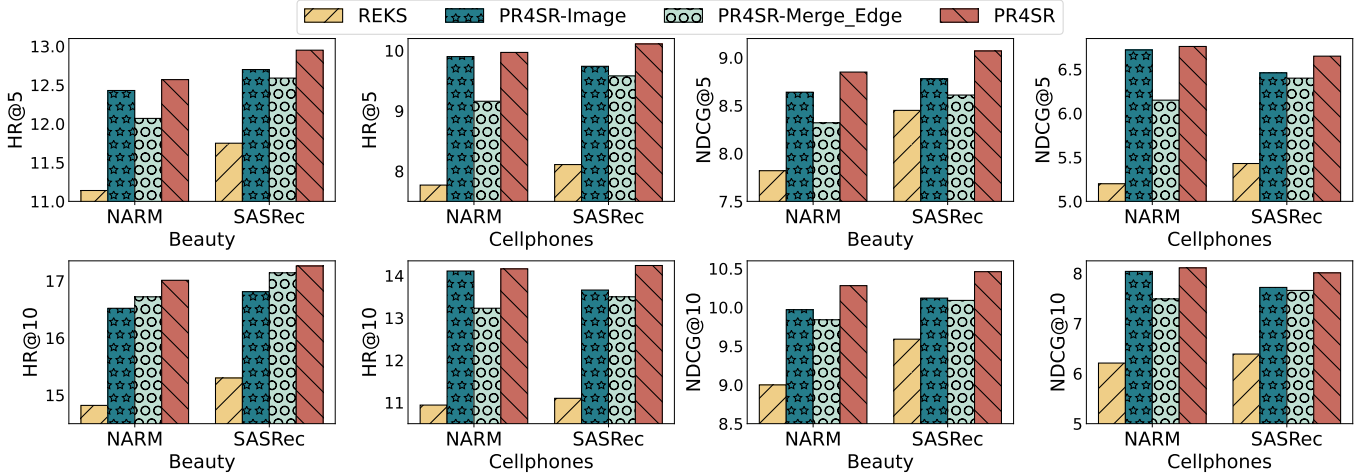


Fig. 6: Ablation performance of different components in the knowledge graph.

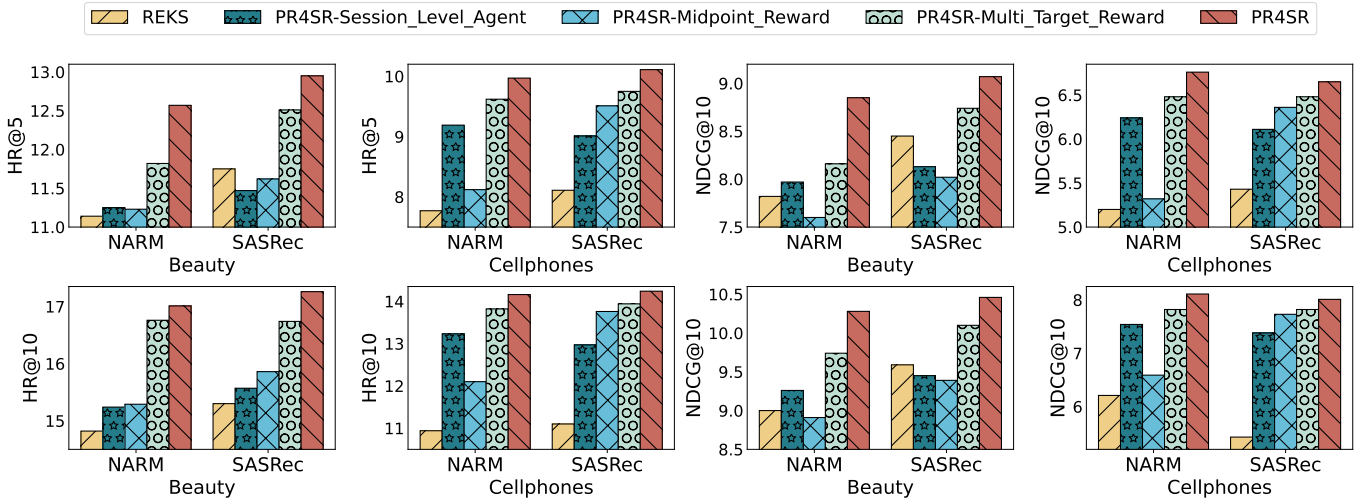


Fig. 7: Ablation performance of different components in the Hierarchical Reinforcement Learning.

in the unexplainable SR model with PR4SR relative to the unexplainable SR model with REKS and “NONE-Avg Improv.” denotes the average improvement in the unexplainable SR model with PR4SR relative to the unexplainable SR models. PR4SR improves the recommendation performance

of unexplainable SR models with all metrics in all datasets. It also outperforms the well-performed explainable SR framework. This proves the effectiveness of PR4SR. Moreover, on the three datasets Cellphones, Baby, and Douban, when the k in top- k is smaller, PR4SR improves more obviously

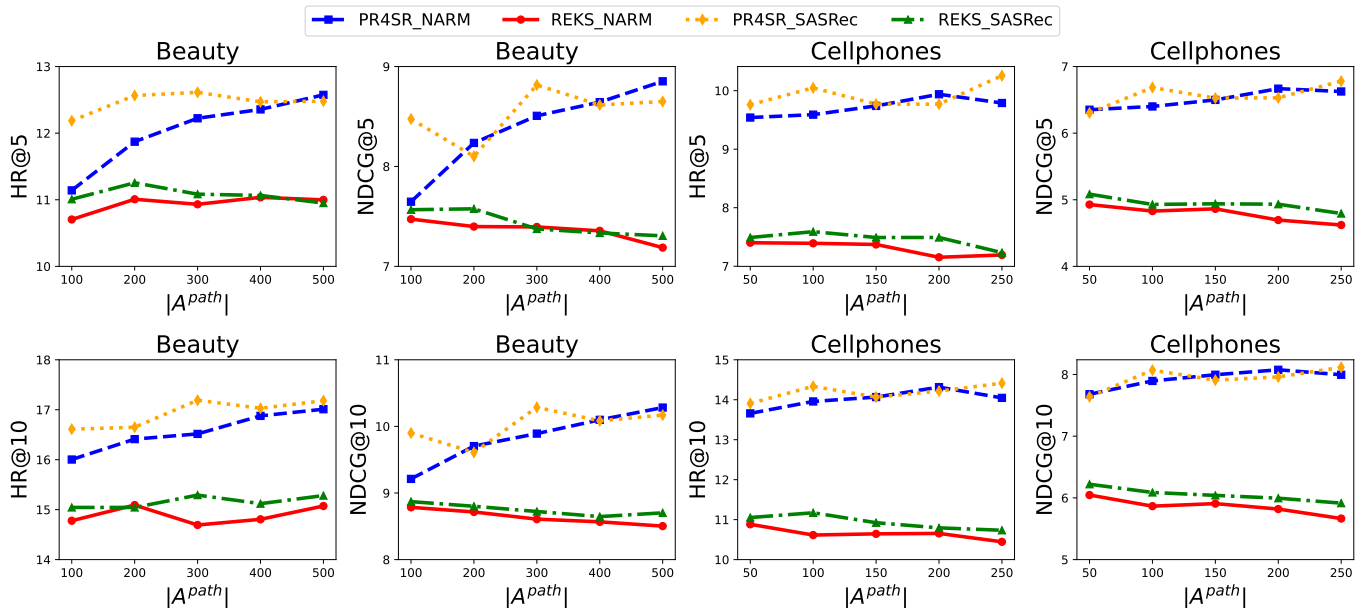


Fig. 8: Impact of path-level agent’s action space size on performance (K = 5,10).

than unexplainable SR models. This is due to the fact that customers pay more attention to the nearest few items when shopping, further indicating that PR4SR is suitable for practical session-based scenarios.

We find that NARM and SASRec achieve the best results on multiple datasets via PR4SR compared to other SR models. And GRU4Rec also has relatively good results compared to SR-GNN and GCSAN on different datasets via PR4SR. This suggests that PR4SR is more suitable for the attention mechanism and RNN-based SR models, compared to the GNN-based SR models.

RQ2: Ablation Study.

As shown in Fig. 6 and Fig. 7, based on NAMR and SASRec, we design ablation experiments to test the design of innovations in the model from two aspects: knowledge graph (KG) and hierarchical reinforcement learning (HRL).

KG: (1) PR4SR-Image denotes the removal of image_feature entity and image_sim relation. (2) PR4SR-Merge_Edge indicates not merging the entities appearing in also_bought and also_viewed. From the results, we can see that all components have an impact on the improvement of the model performance, even though the addition of image information increases exploration space for the path-level agent. Although adding information about the image has little impact on the accuracy of the recommendation, image_feature entity and image_sim relation increase the diversity to the explanation paths. The ablation experiments after dividing also_viewed and also_bought according to different product domains are not considered, as the number of entities and relations in the knowledge graph does not change before and after the division, which produces a relatively small impact on the experiments, but this division method would have provided diversity of explainable paths.

HRL: (1) PR4SR-Session_Level_Agent indicates that the session-level agent is not applied, and only the last item of the session is considered to do path reasoning. (2) PR4SR-

Midpoint_Reward indicates no consideration of Path Midpoint Reward. (3) PR4SR-Multi_Target_Reward indicates no consideration of Multi-Target Reward, but only a single target point. The effect of removing the session-level agent on the model is relatively large, which suggests that we have selected the relatively important items in the session through the session-level agent to be the starting point of the path reasoning. The results also show that the two reward designs also help to improve the performance of the model.

RQ3: Sensitivity of Hyper-parameters

Action space size. Since the product knowledge graph increases the number of entities and relations as new products are added, it is important to increase the RL action space while ensuring that RL can explore valuable information. From Fig. 8, it can be seen that in most cases the performance of REKS decreases or slightly improves as the action space increases, while the performance of PR4SR shows a steady increase. It further shows that PR4SR is able to adapt to changing knowledge graphs compared to REKS.

TABLE 9: The effect of the length $|T|$ on Beauty.

Setting	NARM		SASRec	
	HR@5	NDCG@5	HR@5	NDCG@5
T=1	11.82	8.16	12.81	8.80
T=2	12.20	8.53	12.54	8.60
T=3	12.04	8.32	12.84	8.84
T=4	12.05	8.42	12.75	8.91
T=5	12.50	8.83	12.32	8.50

Length of successive targets. As shown in Table 9, we find that PR4SR_NARM and PR4SR_SASRec achieve good results at T=5 and T=3, respectively. The experimental results indicate that the optimal T is greater than 2, further illustrating the effectiveness of the multi-target reward.

Impact of α, β in loss function. We further validate the role of α and β in the loss function. Fig. 9 shows the results for different parameters, including α and β . The framework performance is relatively stable when α is between 0.005

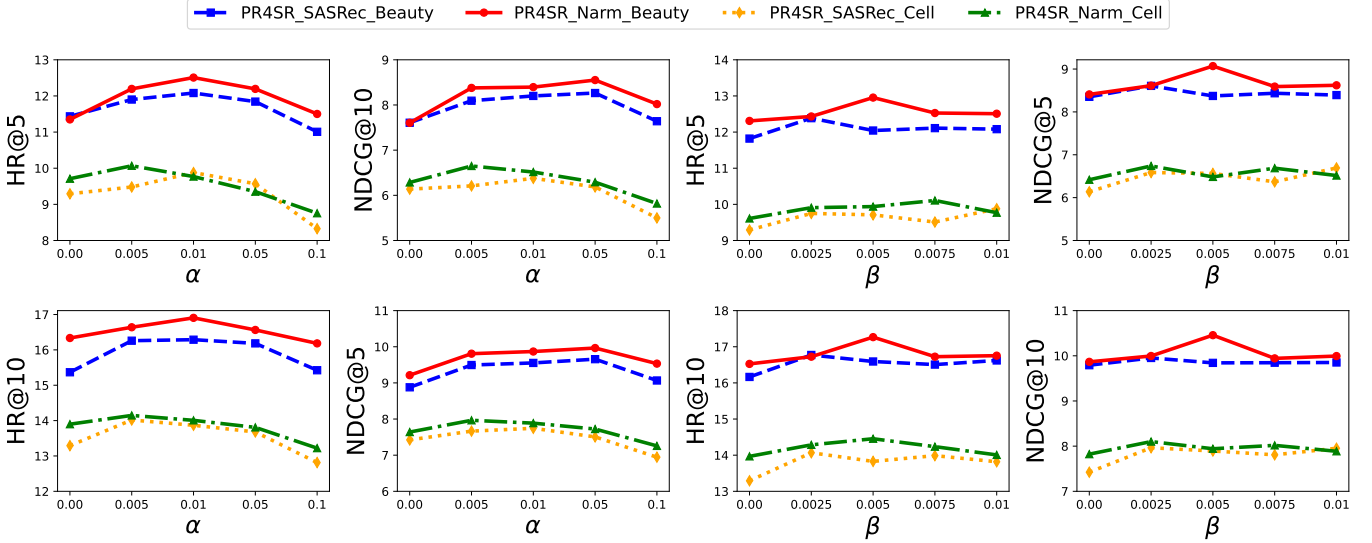


Fig. 9: Impact of different hyper-parameters ($K = 5, 10$).

and 0.1 and β is between 0.0025 and 0.01, which further illustrates the robustness of our framework.

RQ4: Model Performance on Explanation Task

Case Study. The case study in Fig. 10 demonstrates that our proposed hierarchical reinforcement learning framework can better provide interpretability for SR models as well as strategies for constructing knowledge graphs making connections between entities stronger, where the session record is $item_1, item_2, item_3, item_4$ and the predicted item is $item_5$.

First, our proposed model PR4SR improves the accuracy of prediction. PR4SR selects $item_4$ as the starting point of path reasoning from the session by the session-level agent, and there are multiple inference paths of length 2 between $item_3$ and the predicted entity $item_5$, e.g., $item_3 \xrightarrow{\text{belong_to}} phone\ case \xrightarrow{\text{belong_to}} item_5$, which indicates that the reason for recommending $item_5$ is because $item_5$ and $item_3$ both have a distinct image feature—“pearl”. Whereas the REKS only considered $item_4$ as the starting point for path reasoning, under the restriction of exploring paths of length 2, there does not exist a path connecting $item_4$ to the predicted entity $item_5$ and exists only $item_4 \xrightarrow{\text{purchase}} user \xrightarrow{\text{purchase}} item_8$.

Second, the method of constructing the knowledge graph complements the information characterizing the items and increases the diversity of explanations. As in the Fig. 10 “pearl” indicates that the picture of the item contains pearls. If the knowledge graph is constructed without considering the feature information of the image, due to the lack of attribute information of $item_1$ in the dataset, there is no connecting relation between $item_1$ and the category “phone case”, and $item_5$ cannot be inferred from $item_1$. By adding the newly added feature information of the image “pearl”, we can infer that there are two explainable paths, $item_1 \xrightarrow{\text{image_sim}} pearl \xrightarrow{\text{image_sim}} item_5$ and $item_3 \xrightarrow{\text{image_sim}} pearl \xrightarrow{\text{image_sim}} item_5$. There exists a path, $item_3 \xrightarrow{\text{also_bought_diff}} item_6 \xrightarrow{\text{also_bought_diff}} item_5$,

which states that after splitting relations based on different product domains we can determine which product class the entity belongs to by the type of the relation, increasing the explainability of the path.

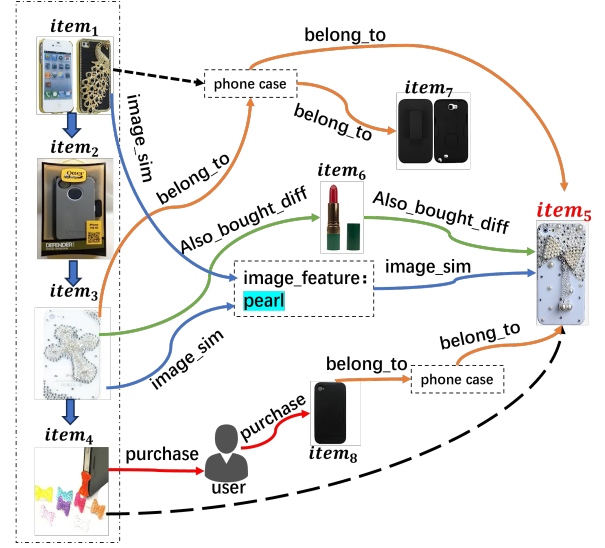


Fig. 10: An example of explainable paths generated by our method. The direction of the arrow on the right side of the picture indicates the direction of path reasoning.

User Study. In this study, we verify the superiority of our explanation from a human perspective. A group of 50 participants was chosen for the study. We randomly selected 30 cases from datasets Beauty, Cellphones, and Baby. Each case consisted of two different explanatory content from PR4SR and REKS. We provided detailed data in the questionnaire, including the user’s session history, the starting point for path reasoning, and a description of the explainable path. Participants evaluated the choice of the starting point for path reasoning and the explainable paths based on their understanding of the session history.

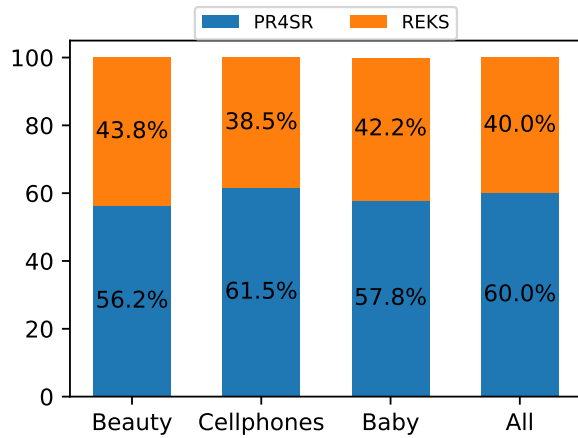


Fig. 11: Percentage of user choices.

The results in Fig. 11 indicate that PR4SR outperforms REKS in terms of explainability in the case of all three datasets. This is more related to the introduction of the `image_feature` and `title_feature` entities, as well as `image_sim` and `title_sim`, which are in line with the actual buying experience of users in the shopping system.

7 CONCLUSION

In this paper, we introduce PR4SR, the first framework to utilize hierarchical reinforcement learning for providing explainable path reasoning for existing unexplainable SR models. Within this framework, the session-level agent selects key items from the session history as the path starting point and the path-level agent performs path reasoning in the knowledge graph. In particular, we motivate the learning of skip behaviors of sequential patterns in the session scenario by multi-target reward and use path-midpoint reward to improve the exploration efficiency. Meanwhile, incorporating image information into the knowledge graph enhances its completeness and the diversity of explanation paths. Extensive experiments on four datasets show that our framework outperforms current unexplainable SR models and the explainable SR framework in terms of both recommendation performance and explainability. In future research work, we will explore generic explainable frameworks for the currently existing unexplainable cross-domain recommendation models as well as constructing product knowledge graphs with more diverse explainability.

REFERENCES

- [1] B. Hidasi, A. Karatzoglou, L. Baltrunas, and D. Tikk, "Session-based recommendations with recurrent neural networks," *arXiv preprint arXiv:1511.06939*, 2015.
- [2] J. Li, P. Ren, Z. Chen, Z. Ren, T. Lian, and J. Ma, "Neural attentive session-based recommendation," in *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 2017, pp. 1419–1428.
- [3] W.-C. Kang and J. McAuley, "Self-attentive sequential recommendation," in *2018 IEEE international conference on data mining (ICDM)*. IEEE, 2018, pp. 197–206.
- [4] 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.
- [5] 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.
- [6] 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.
- [7] S.-J. Park, D.-K. Chae, H.-K. Bae, S. Park, and S.-W. Kim, "Reinforcement learning over sentiment-augmented knowledge graphs towards accurate and explainable recommendation," in *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, 2022, pp. 784–793.
- [8] K. Zhao, X. Wang, Y. Zhang, L. Zhao, Z. Liu, C. Xing, and X. Xie, "Leveraging demonstrations for reinforcement recommendation reasoning over knowledge graphs," in *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, 2020, pp. 239–248.
- [9] Z. Fu, Y. Xian, R. Gao, J. Zhao, Q. Huang, Y. Ge, S. Xu, S. Geng, C. Shah, Y. Zhang *et al.*, "Fairness-aware explainable recommendation over knowledge graphs," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 69–78.
- [10] Y. Xian, Z. Fu, H. Zhao, Y. Ge, X. Chen, Q. Huang, S. Geng, Z. Qin, G. De Melo, S. Muthukrishnan *et al.*, "Cafe: Coarse-to-fine neural symbolic reasoning for explainable recommendation," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 1645–1654.
- [11] H. Lu, W. Ma, Y. Wang, M. Zhang, X. Wang, Y. Liu, T.-S. Chua, and S. Ma, "User perception of recommendation explanation: Are your explanations what users need?" *ACM Transactions on Information Systems*, vol. 41, no. 2, pp. 1–31, 2023.
- [12] C. Liu, W. Wu, S. Wu, L. Yuan, R. Ding, F. Zhou, and Q. Wu, "Social-enhanced explainable recommendation with knowledge graph," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [13] Z. Lyu, Y. Wu, J. Lai, M. Yang, C. Li, and W. Zhou, "Knowledge-enhanced graph neural networks for explainable recommendation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 35, no. 5, pp. 4954–4968, 2022.
- [14] Y. Wei, X. Qu, X. Wang, Y. Ma, L. Nie, and T.-S. Chua, "Rule-guided counterfactual explainable recommendation," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [15] X. Wang, K. Liu, D. Wang, L. Wu, Y. Fu, and X. Xie, "Multi-level recommendation reasoning over knowledge graphs with reinforcement learning," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 2098–2108.
- [16] X. V. Lin, R. Socher, and C. Xiong, "Multi-hop knowledge graph reasoning with reward shaping," *arXiv preprint arXiv:1808.10568*, 2018.
- [17] Z. Sun, J. Yang, J. Zhang, A. Bozzon, L.-K. Huang, and C. Xu, "Recurrent knowledge graph embedding for effective recommendation," in *Proceedings of the 12th ACM conference on recommender systems*, 2018, pp. 297–305.
- [18] S. Geng, Z. Fu, J. Tan, Y. Ge, G. De Melo, and Y. Zhang, "Path language modeling over knowledge graphs for explainable recommendation," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 946–955.
- [19] W. Ma, M. Zhang, Y. Cao, W. Jin, C. Wang, Y. Liu, S. Ma, and X. Ren, "Jointly learning explainable rules for recommendation with knowledge graph," in *The world wide web conference*, 2019, pp. 1210–1221.
- [20] X. Wang, D. Wang, C. Xu, X. He, Y. Cao, and T.-S. Chua, "Explainable reasoning over knowledge graphs for recommendation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 5329–5336.
- [21] Q. Zhu, X. Zhou, J. Wu, J. Tan, and L. Guo, "A knowledge-aware attentional reasoning network for recommendation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, 2020, pp. 6999–7006.
- [22] R. T. Icarte, T. Q. Klassen, R. Valenzano, and S. A. McIlraith, "Reward machines: Exploiting reward function structure in reinforcement learning," *Journal of Artificial Intelligence Research*, vol. 73, pp. 173–208, 2022.
- [23] Y. Xian, Z. Fu, S. Muthukrishnan, G. De Melo, and Y. Zhang, "Reinforcement knowledge graph reasoning for explainable recommendation," in *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*, 2019, pp. 285–294.

- [24] H. Wu, H. Fang, Z. Sun, C. Geng, X. Kong, and Y.-S. Ong, "A generic reinforced explainable framework with knowledge graph for session-based recommendation," in *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 2023, pp. 1260–1272.
- [25] J. Tang and K. Wang, "Personalized top-n sequential recommendation via convolutional sequence embedding," in *Proceedings of the eleventh ACM international conference on web search and data mining*, 2018, pp. 565–573.
- [26] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [27] 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.
- [28] 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.
- [29] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [30] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.
- [31] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [32] Z. Cui, H. Chen, L. Cui, S. Liu, X. Liu, G. Xu, and H. Yin, "Reinforced kgs reasoning for explainable sequential recommendation," *World Wide Web*, vol. 25, no. 2, pp. 631–654, 2022.
- [33] Y. Zhao, X. Wang, J. Chen, Y. Wang, W. Tang, X. He, and H. Xie, "Time-aware path reasoning on knowledge graph for recommendation," *ACM Transactions on Information Systems*, vol. 41, no. 2, pp. 1–26, 2022.
- [34] C. Ma, P. Kang, and X. Liu, "Hierarchical gating networks for sequential recommendation," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 825–833.
- [35] H. Hou and C. Shi, "Explainable sequential recommendation using knowledge graphs," in *Proceedings of the 5th International Conference on Frontiers of Educational Technologies*, 2019, pp. 53–57.
- [36] Y. Li, H. Chen, Y. Li, L. Li, S. Y. Philip, and G. Xu, "Reinforcement learning based path exploration for sequential explainable recommendation," *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [37] X. Li, Y. Shen, and L. Chen, "Mcore: Multi-agent collaborative learning for knowledge-graph-enhanced recommendation," in *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2021, pp. 330–339.
- [38] C. Stauffer and W. E. L. Grimson, "Adaptive background mixture models for real-time tracking," in *Proceedings. 1999 IEEE computer society conference on computer vision and pattern recognition (Cat. No PR00149)*, vol. 2. IEEE, 1999, pp. 246–252.
- [39] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel *et al.*, "Soft actor-critic algorithms and applications," *arXiv preprint arXiv:1812.05905*, 2018.
- [40] M. M. Afsar, T. Crump, and B. Far, "Reinforcement learning based recommender systems: A survey," *ACM Computing Surveys*, vol. 55, no. 7, pp. 1–38, 2022.
- [41] Q. Ai, V. Azizi, X. Chen, and Y. Zhang, "Learning heterogeneous knowledge base embeddings for explainable recommendation," *Algorithms*, vol. 11, no. 9, p. 137, 2018.
- [42] R. He and J. McAuley, "Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering," in *proceedings of the 25th international conference on world wide web*, 2016, pp. 507–517.
- [43] H. Jiang, C. Li, J. Cai, and J. Wang, "Rcen: A reinforced and contrastive heterogeneous network reasoning model for explainable news recommendation," in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2023, pp. 1710–1720.



Yang Cao is currently pursuing the master's degree with the Department of Computer Science and Technology, East China Normal University, China. His research interests include recommendation system, knowledge graph and reinforcement learning.



Shuo Shang is currently a professor of computer science with the University of Electronic Science and Technology of China. He was a senior scientist with the Inception Institute of Artificial Intelligence (IIAI), leading its data mining research group. His research interests include big data, data mining, and machine learning.



Jun Wang received the Ph.D. degree in electrical engineering from Columbia University, New York, NY, USA, in 2011. Currently, he is a professor at the School of Computer Science and Technology, East China Normal University and an adjunct faculty member of Columbia University. From 2010 to 2014, he was a research staff member at IBM T. J. Watson Research Center, Yorktown Heights, NY, USA. His research interests include machine learning, data mining, mobile intelligence, and computer vision.



Wei Zhang received his Ph.D. degree in computer science and technology from Tsinghua university, Beijing, China, in 2016. He is currently a professor in the School of Computer Science and Technology, East China Normal University, Shanghai, China. His research interests mainly include user data mining and machine learning applications. He is a senior member of China Computer Federation.