

Towards Multi-Subsession Conversational Recommendation

Yu Ji*

Tongji University
Shanghai, China
2230779@tongji.edu.cn

Hang Yu

Tongji University
Shanghai, China
2053881@tongji.edu.cn

Qi Shen*

Tongji University
Shanghai, China
1653282@tongji.edu.cn

Yiming Zhang

Tongji University
Shanghai, China
2030796@tongji.edu.cn

Shixuan Zhu

Tongji University
Shanghai, China
2130768@tongji.edu.cn

Chuan Cui

Tongji University
Shanghai, China
cuichuan@tongji.edu.cn

Zhихua Wei[†]

Tongji University
Shanghai, China
zhихua_wei@tongji.edu.cn

ABSTRACT

Conversational recommendation systems (CRS) could acquire dynamic user preferences towards desired items through multi-round interactive dialogue. Previous CRS works mainly focus on the single conversation (**subsession**) that the user quits after a successful recommendation, neglecting the common scenario where the user has multiple conversations (**multi-subsession**) over a short period. Therefore, we propose a novel conversational recommendation scenario named **Multi-Subsession Multi-round Conversational Recommendation (MSMCR)**, where the user would still resort to CRS after several subsessions and might preserve vague interests, and the system would proactively ask attributes to activate user interests in the current subsession. To fill the gap in this new CRS scenario, we devise a novel framework called **Multi-Subsession Conversational Recommender with Activation Attributes (MSCAA)**. Specifically, we first develop a context-aware recommendation module, comprehensively modeling user interests from historical interactions, previous subsessions, and feedback in the current subsession. Furthermore, an attribute selection policy module is proposed to learn a flexible strategy for asking appropriate attributes to elicit user interests. Finally, we design a conversation policy module to manage the above two modules to decide actions between asking and recommending. Extensive experiments on four datasets verify the effectiveness of our MSCAA framework for the proposed MSMCR setting.

*Both authors contributed equally to this research.

[†]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 ACM 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](https://permissions.acm.org).
Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/XXXXXXX.XXXXXXX>

CCS CONCEPTS

• **Information systems** → **Users and interactive retrieval; Recommender systems.**

KEYWORDS

Conversational Recommendation; Human-in-the-Loop Learning; Recommender Systems

ACM Reference Format:

Yu Ji, Qi Shen, Shixuan Zhu, Hang Yu, Yiming Zhang, Chuan Cui, and Zhихua Wei. 2018. Towards Multi-Subsession Conversational Recommendation. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Recent years have witnessed increasing attention and research effort in the conversational recommendation system (CRS), which aims to elicit dynamic user preferences and make successful recommendations through a real-time multi-round conversation with the user [9]. To tackle this trending research topic, researchers have proposed various methods based on different conversational recommendation problem settings [3, 17, 20, 45]. Among these studies, the multi-round conversational recommendation (MCR) setting is envisioned as the most realistic CRS setting so far [8, 39]. In this work, we focus on the MCR setting where the system takes actions based on the user's current needs in each turn, either asking questions about user preferences on attributes or recommending items, with the purpose of recommending successfully within fewer turns.

Although existing MCR studies have made significant progress, they all focus on a single conversation episode that the user quits after receiving the satisfactory item recommendation [8, 17, 18]. We argue that this single conversation setting overlooks the prevalence of **multiple** conversations in the real-world CRS scenario. The user would continue the dialogue with the system after a successful recommendation to browse items on other topics aimlessly or obtain additional system suggestions. Under the situation of multiple conversations, the user might also preserve possible **vague** interest

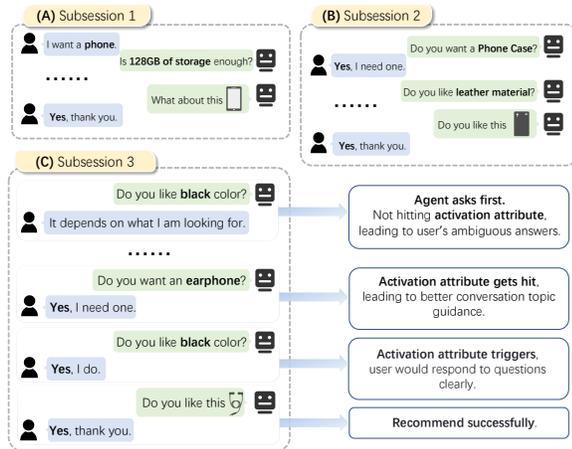


Figure 1: A toy example of Multi-Subsession Multi-round Conversational Recommendation. The lower right side explains the multi-round process of the 3rd sub-session (C) of one session in detail.

rather than clearly express his/her preference, which is more realistic in the world and breaks the assumption of the previous MCR works [17, 18].

Figure 1 illustrates an interaction process (**session**) that the user has three conversations (**subsessions**). In the first sub-session of one session, the user successfully seeks his/her target item *Phone*. At this time, he/she might require other items related to *Phone* but be unaware of the related attributes to find them. After several system suggestions, he/she further gets satisfactory items (*Phone Case* and *Earphone*) in the subsequent subsessions. More specifically considering vague user interest shown in the 3rd sub-session (C) of Figure 1, the user would be overwhelmed by over-specific questions (e.g., *black color*) before the main attributes (e.g., *earphone*) of the target item, i.e., **activation attributes** are determined, which leads to ambiguous answers (e.g., *It depends on what I am looking for.* for attribute question *black color*). Once the activation attribute *earphone* is confirmed by the user as the topic of the current sub-session, he/she would activate his/her interest and respond explicitly on system-asked questions towards the current topic, naturally including over-specific questions in the early turns.

To this end, we extend MCR [17, 18, 39] to a general setting, i.e., **Multi-Subsession Multi-round Conversational Recommendation (MSMCR)**. In this scenario, the user has several subsessions with the system over a short period. And the system would proactively ask questions on attributes to elicit dynamic user interests towards desired items in the subsequent subsessions (i.e., after the first sub-session of one session). Since our MSMCR scenario involves the user’s previous subsessions and activation attributes of the current sub-session, surpassing the previous MCR works, we summarize two main challenges to handle this scenario listed as follows:

- **How to model user interests comprehensively?** Except for historical user-item interactions and feedback in the current sub-session, the additional previous subsessions information should be further considered in our MSMCR scenario. As such, how to organize and aggregate these multiple information sources effectively for modeling user interests is a problem.

- **How to ask appropriate questions?** In our scenario, the system is required to ask appropriate questions to elicit the user’s interest for recommending his/her desired item. Concretely, the user preserves unclear interest at the beginning of the current sub-session, where the system should activate his/her current interest as early as possible. After that, the system needs to ask attributes to eliminate the uncertainty of candidate items, following existing MCR works [17, 18]. As a result, it is essential to learn a flexible attribute selection strategy to activate user interests and then reduce the uncertainty of candidate items.

To effectively address the aforementioned challenges for MSMCR, we propose a novel method named **Multi-Subsession Conversational Recommender with Activation Attributes (MSCAA)**, which organically combines the two policy modules with a recommendation module. In detail, *Context-aware Recommendation* module models multiple information for item and attribute prediction. More specifically, it extracts user long-term interest representations from the user-item-attribute graph via a relational graph convolution network. Then for previous online subsessions information, we adopt dual gated recurrent units to learn user short-term interest representations. Besides, the user feedback, i.e., accepted attributes, rejected attributes and items in the current sub-session, is also aggregated to represent the user’s current interest. Finally, a unified transformer encoder is applied to fuse the above three different interests for prediction. The *Attribute Selection Policy* module automatically learns the strategy for asking activation attributes or others based on predicted user-like score and entropy of candidate attributes. And eventually, the *Conversation Policy* module further utilizes the above two modules and then decides actions between asking and recommending.

Our main contributions of this work are summarized below:

- We introduce a novel multi-subsession conversation recommendation scenario named MSMCR, a natural extension of the multi-round conversational recommendation, which considers previous online subsessions and system-guided user interest activation.
- Building upon our new scenario, we devise a novel framework MSCAA that combines the user’s long- and short-term interaction and real-time feedback to model user interests comprehensively and learn a flexible asking policy to activate and crystallize user demands.
- We evaluate MSCAA on four adapted datasets and conduct a human evaluation to demonstrate the effectiveness of our method.

2 RELATED WORKS

Conversational Recommendation. Conversational recommendation systems (CRS) [9, 16] leverage immediate user feedback information to capture dynamic user preferences and make precise recommendations in conversational scenarios. Existing CRS research can be categorized into four directions. (1) Multi-round Conversational Recommendation (MCR) [17, 29]. In each dialogue turn, MCR chooses between attribute asking and item recommendation to achieve an early and accurate hit on the target item. (2) Exploration-Exploitation Trade-offs [21, 40]. These methods focus on the cold-start situation and propose a paradigm to balance users’ exploration and exploitation trade-offs. (3) Question-based User Preference Elicitation works [37, 45] ask “clarification/clarifying” questions to figure out how to gain more preference information

within limited conversations. (4) Dialogue Understanding and Generation [20, 44] focuses on how to understand and provide a real conversation and make a smoother response.

Multi-round Conversational Recommendation. Among all the above CRS research scenarios, we focus on the MCR task, in which the system alternates between asking attribute-based questions and recommending items several times to seek the user’s desired item within fewer interaction turns. The conversation strategy is the core design in MCR, for deciding whether to ask or recommend, and the strategy is typically modeled by deep reinforcement learning (DRL) [23, 30] to tackle this multi-step decision-making process. For instance, CRM [29] brings up the conversation setting starting with the agent’s query, which gathers user feedback information and gives recommendations in the end. EAR [17] extends the CRM framework where the agent can recommend many times. SCPR [18] proposes an interactive path reasoning method to find appropriate candidates for items and attributes. FPAN [34] utilizes online user feedback to update user embedding via the gate mechanism based on EAR framework. UNICORN [8] contributes a dynamic weighted graph-based unified policy learning framework. Furthermore, CRIF [14] explicitly uses item feedback to capture the implicit user preference for the recommendation and employs inverse reinforcement learning to learn a better conversation strategy. Considering RL often suffers from unstable learning, some recent works attempt non-RL alternative methods, such as a simple rule [38] or decision tree [10], to achieve comparable performance. In addition to the above classical MCR task, more and more studies start to focus on exploring the new MCR scenario applied in many fields. For example, MIMCR [39] considers the incompleteness and diversity of user interests and extends MCR to a more realistic setting named multiple choice questions, which is adopted by follow-ups [15, 42]. Bundle-MCR [12] incorporates bundle generation with MCR mechanisms. GPR [41] proposes a two-level graph path reasoning method that also recommends items when asking attributes. MetaCRS [5] contributes a meta reinforcement learning framework for the scenario of cold-start users. Different from these works, and some studies [31, 43] that utilize the user’s historical interaction sequence to improve recommendations in MCR, we concentrate on a **novel scenario** where the user has several dialogues with the system in a short time and might preserve unclear interests.

3 MSMCR SCENARIO

3.1 Definition

In this scenario, we define the set of users \mathcal{U} and items \mathcal{V} . We collect all attributes \mathcal{A} corresponding to items, and each item v is associated with a set of attributes \mathcal{A}_v . Unlike previous MCR works that call a conversation session, we define a conversation episode as a *subsession* s . And a *session* $S_u = [s_u^1, s_u^2, \dots, s_u^{n-1}, s_u^n]$ consists of multiple subsessions where n is variational length for each session. Moreover, a session can be divided into two parts: the previous subsessions $P_{S_u} = [s_u^1, s_u^2, \dots, s_u^{n-1}]$ and the current subsession s_u^n . The goal of MSMCR is to recommend the desired item v^n of user u in the current subsession s_u^n within fewer turns, based on P_{S_u} and the current subsession information.

For a user $u \in \mathcal{U}$, the workflow of MSMCR is listed as follows:

(1) After several subsessions (i.e., P_{S_u}) successfully end, the user

resorts to the subsequent conversation (i.e., s_u^n) instead of quitting the session. In this state, the user might initially have no explicit attribute query, and hence the current subsession is started from the system side. (2) Then, the system is free to *ask* questions about an attribute from the candidate attribute set \mathcal{A}_{cand} (e.g., $\mathcal{A}_{cand} = \mathcal{A}$ at first) or to *recommend* a certain number of items (e.g., top- K) from the candidate item set \mathcal{V}_{cand} . (3) Next, user u provides feedback, i.e., accept, **unknown**, or reject for the asked attribute, and accept or reject for the recommended items. (4) After that, the system updates the sets of candidate attributes and items based on user feedback. (5) Within multiple iterations of step (2)-(4), the system elicits clearer user interests and provides more accurate recommendations. The current subsession will terminate when the system recommends successfully, or the interaction turn reaches the maximum T .

Different from the user setting that has clear attribute preferences in previous MCR works [17, 18], the user’s preference might be vague and should be proactively guided [44] by the system in the MSMCR scenario. That is, in step (3), the user will respond “**unknown**” for the other attributes when the **activation attributes** $\mathcal{A}_{v^n}^*$ (i.e., several main attributes that enable the user to trigger his/her current demands) have not been clarified. And after one of the activation attributes is hit by the system, the user would generate explicit feedback (accept or reject) for all attributes.

3.2 General Framework

Similar to existing MCR works, we formulate our MSMCR problem as a Markov Decision Process (MDP) of interaction between the user and the recommendation system agent. The goal of our framework is to learn the policy network π_τ (cf. Section 4.3) to maximize the expected cumulative rewards for the overall conversations. We decompose one conversation turn into four steps: **state**, **action**, **transition** and **reward** under our framework as follows.

State. The system maintains the conversation state s_t at each turn t in the subsession, which encodes the user feedback and candidate item information. Specially, user u ’s feedback includes the accepted attributes \mathcal{A}_{acc}^t , rejected attributes \mathcal{A}_{rej}^t and rejected items \mathcal{V}_{rej}^t .

Action. The agent takes actions according to the conversation state, where the action can be asking an attribute or recommending items. If the system decides to ask, it will select an attribute from the decision of an attribute agent using policy π_a (cf. Section 4.2). While for recommendation action, the system will select top- K items based on the recommendation score (cf. Section 4.1).

Transition. When the agent takes actions, and the user provides corresponding feedback to the attribute or items, the current state will transition to the next s_{t+1} . In this step, we update the sets \mathcal{A}_{acc} , \mathcal{A}_{rej} , \mathcal{V}_{rej} , \mathcal{V}_{cand} and \mathcal{A}_{cand} based on user feedback following the transition operation of [32, 39]. Note that in our MSMCR setting, user feedback for the asked attribute will be “unknown” until an activation attribute is hit. In this case, we do not update them.

Reward. We design six kinds of rewards in our framework: (1) r_{rec_acc} and (2) r_{rec_rej} denote positive and negative rewards when the user accepts or rejects the recommended item, respectively; (3) r_{ask_acc} , (4) r_{ask_rej} and (5) r_{ask_unk} denote positive, negative and negative rewards when the user accepts, rejects and responds “unknown” to the asked attribute, respectively; (6) r_{quit} is a strongly negative reward if the maximum subsession turn T is reached.

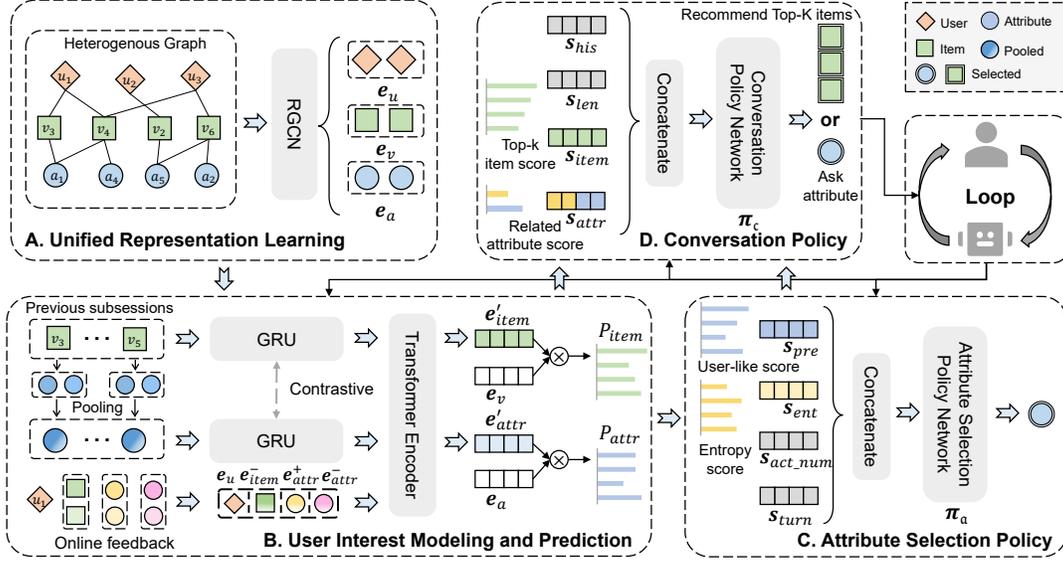


Figure 2: The overview of MSCAA. Our framework includes recommendation modules (A and B) and policy modules (C and D).

4 METHODOLOGY

As shown in Figure 2, our framework contains three main components: (1) *Context-aware Recommendation* module including Unified Representation Learning (A) and User Interest Modeling and Prediction (B), which comprehensively models user interests based on the information from historical interactions, previous subsessions, and the current subsession; (2) *Attribute Selection Policy* module (C), which aims to select an appropriate asked attribute based on attribute-related information; (3) *Conversation Policy* module (D), which decides actions between asking and recommending based on the overall conversation.

4.1 Context-aware Recommendation

A well-learned recommendation module will improve the performance of the online system agent [14]. To this end, we carefully design the following three parts in the context-aware recommendation: *unified representation learning* to extract long-term user interest, *user interest modeling* to learn short-term interest from subsessions (i.e., context) and then integrate different interests, and finally, *prediction* to estimate probabilities that how the user likes items and attributes.

4.1.1 Unified Representation Learning. We first construct a unified heterogeneous graph [25, 39] based on historical user-item interactions and static attribute-item relations, including three kinds of nodes: users, items, and attributes. This graph can be denoted as $G = \{N = \{\mathcal{U}, \mathcal{V}, \mathcal{A}\}, E = \{\mathcal{E}_{uv}, \mathcal{E}_{av}\}\}$, where the user-item edge $(u, v) \in \mathcal{E}_{uv}$ denotes that user u has interacted with item v and the attribute-item edge $(a, v) \in \mathcal{E}_{av}$ means that attribute a belongs to item v . Then, we introduce a L_g -layer relational graph convolutional network (RGCN) [28] to learn the node representations in the graph G . At first, each user, item, and attribute node is assigned with a unique node index and then initialized with the node embedding matrix $E^0 \in \mathbb{R}^{|N| \times D}$, as the input embeddings of the first layer of RGCN. And for each layer until L_g , RGCN will first model different edge types based on node neighbors, and then the aggregated edge features will be integrated into node features. We

take item node v which involves two kinds of edges $\mathcal{R}_v = \{r_{uv}, r_{av}\}$ as an example as follows:

$$e_v^{l+1} = \text{ReLU}\left(\sum_{r \in \mathcal{R}_v} \sum_{i \in \mathcal{N}_v^r} \frac{1}{\sqrt{|\mathcal{N}_v^r| |\mathcal{N}_i^r|}} \mathbf{W}_r^l e_i^l + \mathbf{W}_0^l e_v^l\right), \quad (1)$$

where \mathcal{N}_v^r denotes the neighbor nodes of node v under relation r . \mathbf{W}_0^l , $\mathbf{W}_{r_{uv}}^l$ and $\mathbf{W}_{r_{av}}^l \in \mathbb{R}^{D \times D}$ are trainable parameters. We can obtain the refined user and attribute embeddings e_u^{l+1} , e_a^{l+1} in the same way. After L_g layers information propagation, we combine node embeddings of each layer to capture different semantics [11]:

$$e_n = \frac{1}{L_g + 1} \sum_{j=0}^{L_g} e_n^j, \quad \forall n \in N = \{\mathcal{U}, \mathcal{V}, \mathcal{A}\}. \quad (2)$$

In the rest of this paper, we will use e_u , e_v and e_a to represent refined user, item and attribute embedding, respectively. Note that e_u can be regarded as **long-term** user interest [25, 34].

4.1.2 User Interest Modeling. In our scenario, the user also has previous subsessions that occur recently with the system online, which are probably relevant, and hence crucial for user short-term interest modeling. For each subsession s_u^i in P_{S_u} , much abundant information should be considered, including the accepted item v^i and its associated attributes \mathcal{A}_{v^i} , which reflects short-term user preference on item and attributes. Therefore, we have two sequences from P_{S_u} : the previous desired item sequence $P_{S_u}^{\mathcal{V}} = [v^1, v^2, \dots, v^{n-1}]$ and a sequence of the previous desired attribute sets $P_{S_u}^{\mathcal{A}} = [\mathcal{A}_{v^1}, \mathcal{A}_{v^2}, \dots, \mathcal{A}_{v^{n-1}}]$. Following the classic sequential recommendation method [13, 36], we separately model the above two sequences via gated recurrent units (GRU) [6] to capture the temporal interactions between each subsession. Specifically, for item sequence $P_{S_u}^{\mathcal{V}}$, we use the final hidden state of the GRU as the representation of the user's **short-term item-level** interest:

$$e'_{item} = \text{GRU}(\{e_v^1, e_v^2, \dots, e_v^{n-1}\}), \quad (3)$$

where e_v^i is the representation of i -th item in $P_{S_u}^{\mathcal{V}}$. Similarly, we can get the user's **short-term attributes-level** interest e'_{attr} via another GRU. Note that $P_{S_u}^{\mathcal{A}}$ is a attribute set sequence. Therefore,

we first generate the overall representation of each attribute set by mean pooling operation: $\mathbf{g}_a^i = \text{MEAN}(\{\mathbf{e}_a | a \in \mathcal{A}_{v^i}\})$ [36].

For the current subsession, the online feedback including the accepted and rejected attributes and the rejected items, is also pivot information for revealing the **current** user interest. To comprehensively integrate the user’s long- and short-term interest as well as current interest, we adopt a unified L_t -layer Transformer encoder [33] to fuse these multiple kinds of information. In our case, the input sequence is $[\mathbf{e}_u, \mathbf{e}'_{item}, \mathbf{e}'_{attr}, \mathbf{e}^-_{item}, \mathbf{e}^+_{attr}, \mathbf{e}^-_{attr}]$, where $\mathbf{e}^-_{item}, \mathbf{e}^+_{attr}, \mathbf{e}^-_{attr} \in \mathbb{R}^{1 \times D}$ denotes the aggregated representation of rejected items, accepted and rejected attributes. Here, we use the mean pooling operation to generate these three kinds of information, e.g., $\mathbf{e}^-_{item} = \text{MEAN}(\{\mathbf{e}_v | v \in \mathcal{V}_{rej}\})$. For the first few turns that may have no rejected items, accepted attributes or rejected attributes, the corresponding $\mathbf{e}^-_{item}, \mathbf{e}^+_{attr}, \mathbf{e}^-_{attr}$ is masked for Transformer input. Finally, we can obtain the user’s **final item-level and attribute-level interest** from different kinds of information, i.e., \mathbf{e}'_{item} and \mathbf{e}'_{attr} in the output embeddings of Transformer encoder. Here we maintain the notations of transformed output embeddings unchanged for the sake of simplicity.

4.1.3 Prediction. Based on the learned representations above, we can estimate the probability that the user likes items and attributes:

$$P_{item}(u, v) = \tanh(\mathbf{e}'_{item} \top \mathbf{e}_v); P_{attr}(u, a) = \tanh(\mathbf{e}'_{attr} \top \mathbf{e}_a), \quad (4)$$

where v and a are candidate item in \mathcal{V}_{cand} and attribute in \mathcal{A}_{cand} , respectively.

4.2 Attribute Selection Policy

In our MSMCR scenario, due to the existence of special activation attributes, we expect these activation attributes that the user might well like should be hit as soon as possible. After that, other target attributes asked can better eliminate the uncertainty of candidate items. However, the previous attribute selection rules, like *Max Entropy* [14, 18] and *Max User-like* [39], only consider one aspect of the above issues, e.g., *Max User-like* is mainly suitable for asking attributes that match user interests. Furthermore, these rules cannot deal with the dynamic attribute response setting well in MSMCR, especially the more complex environment in the real-world scenario. Hence we employ an attribute selection policy to adaptively select the appropriate asked attribute to handle the above challenge.

Specifically, the policy network is implemented by a two-layer MLP, which can output the asked attribute based on the user’s conversation state \mathbf{s}_a . The state vector \mathbf{s}_a involves 4 parts of information as follows:

$$\mathbf{s}_a = \mathbf{s}_{pre} \oplus \mathbf{s}_{ent} \oplus \mathbf{s}_{act_num} \oplus \mathbf{s}_{turn}, \quad (5)$$

where \oplus is the concatenation operation, \mathbf{s}_{pre} encodes user interests on all attributes, in which each dimension is the estimated probability (i.e., preference score) P_{attr} calculated by Eq. (4). \mathbf{s}_{ent} is the entropy of each attribute among candidate items \mathcal{V}_{cand} , where each dimension is the entropy w_a of attribute a calculated by $w_a = -p(a) \log p(a) - (1 - p(a)) \log(1 - p(a))$, where $p(a) = |\mathcal{V}_{cand} \cap \mathcal{V}_a| / |\mathcal{V}_{cand}|$, following [14, 17]. Note that only the attributes in \mathcal{A}_{cand} are preserved for the attribute selection, and all other attributes are masked in \mathbf{s}_{pre} and \mathbf{s}_{ent} . These two vectors provide useful user-side and item-side information for attribute selection, following [17]. \mathbf{s}_{act_num} and \mathbf{s}_{turn} denote the

accepted attributes number and the current turn number, to perceive the changes during the conversation. This agent receives the ask-related rewards (r_{ask_acc} , r_{ask_rej} and r_{ask_unk}) and user-quit reward (r_{quit}) only, and updates the state with corresponding state transitions.

4.3 Conversation Policy

The conversation policy is responsible for deciding whether to ask an attribute or recommend items with the purpose of hitting the user’s desired item of the current subsession in fewer turns. We use a simple two-layer MLP to implement the interaction policy network, which maps the state vector \mathbf{s}_c to two actions, i.e., ask or recommend. The action with a higher expected reward will be selected by agent. The state vector \mathbf{s}_c is the concatenation of the following components:

$$\mathbf{s}_c = \mathbf{s}_{his} \oplus \mathbf{s}_{len} \oplus \mathbf{s}_{item} \oplus \mathbf{s}_{attr}, \quad (6)$$

where \mathbf{s}_{his} encodes the conversation history in the subsession and \mathbf{s}_{len} encodes the size of the candidate item set following previous CRS works [17, 18]. \mathbf{s}_{item} is a K -dimension vector which records the top- K candidate item score from the context-aware recommendation module. \mathbf{s}_{attr} represents the preference and entropy score of the attribute by the attribute selection policy. Intuitively, the greater value is in vector \mathbf{s}_{item} , the more probably the agent selects the “recommend” action. The same is for \mathbf{s}_{attr} . By comparing the last two vectors, the agent can make full use of detailed information to decide the appropriate action.

4.4 Model Training

4.4.1 Recommendation Pre-Training. Since the performance of the online agent largely depends on the context-aware recommendation module, we first pre-train this module based on training interaction data for item and attribute prediction.

For item prediction, given a session $S_u = \{P_{S_u}, s_u^n\}$, the target item v^n is considered as a positive ground-truth item, which is expected to rank higher than other negative candidate items. Following [17], We employ pairwise Bayesian Personalized Ranking (BPR) [26] loss:

$$\mathcal{L}_{item} = \sum_{(u, v^n, v') \in \mathcal{D}_1} -\ln \sigma(P_{item}(u, v^n) - P_{item}(u, v')), \quad (7)$$

where $\mathcal{D}_1 := \{(u, v^n, v') | v' \in \mathcal{V} \setminus \mathcal{V}_u\}$ denotes the training set of item pairs where \mathcal{V}_u is the set of items interacted by user u and v' is sampled from non-interacted items of user u . Similar to item prediction, attribute prediction needs to rank the attributes of target item $a \in \mathcal{A}_{v^n}$ ahead of others. To achieve this, we also adopt BPR loss as \mathcal{L}_{attr1} . Besides, it is noted that in our MSMCR scenario, activation attributes are more vital than any other attributes for triggering user interests, which should be ranked higher. Therefore, the loss w.r.t. activation attribute prediction is defined as:

$$\mathcal{L}_{attr2} = \sum_{(u, a, a') \in \mathcal{D}_2} -\ln \sigma(P_{attr}(u, a) - P_{attr}(u, a')), \quad (8)$$

where $\mathcal{D}_2 := \{(u, a, a') | a \in \mathcal{A}_{v^n}^*, a' \in \mathcal{A}_{v^n} \setminus \mathcal{A}_{v^n}^*\}$ denotes the pairwise attribute training data. The final loss of attribute prediction is: $\mathcal{L}_{attr} = \mathcal{L}_{attr1} + \mathcal{L}_{attr2}$.

Meanwhile, we additionally introduce a contrastive loss (InfoNCE [24]) to constrain the uniformity of user item-level and

attribute-level interest:

$$\mathcal{L}_{cont} = \sum_k -\log \frac{\exp(\mathbf{e}'_{item,k} \top \mathbf{e}'_{attr,k}/\tau)}{\sum_{k'} \exp(\mathbf{e}'_{item,k} \top \mathbf{e}'_{attr,k'}/\tau)}, \quad (9)$$

where τ is a temperature hyper-parameter to control the concentration of features. For k -th item-level interest representation $\mathbf{e}'_{item,k}$ within the same minibatch, $\mathbf{e}'_{attr,k}$ is a positive sample that reflects the same user's attribute-level interest, while all others are negative samples. By this constraint term, user interests from different views can complement each other to enhance the final interest representation. Finally, we optimize the recommendation model with the above losses jointly:

$$\mathcal{L}_{rec} = \mathcal{L}_{item} + \mathcal{L}_{attr} + \omega \mathcal{L}_{cont}, \quad (10)$$

where ω controls the weight of constraint loss.

4.4.2 Policy Pre-Training. Then we pre-train the above two policies π_a and π_c to accelerate training by using DAGger [27] with expert demonstrations. It iteratively trains policy via supervised learning on observation-action pairs (s^*, a^*) from the online expert rule: $\hat{\pi}_i = \arg \min_{\pi_i} \sum -\ln \pi_i(a^*|s^*)$, where $i \in \{a, c\}$. For attribute selection policy π_a , the expert online rule is defined that if the number of accepted attributes is less than 1, the attribute agent asks the attribute with *Max User-like* otherwise with *Max Entropy*. And for conversation policy π_c , the expert rule is defined that at turn t , the conversation agent chooses the *ask* action with probability $1 - \frac{t}{T}$ or *recommends* top- K items. The intuition behind this is that attributes are expected to be asked in the first few turns, and later the recommendation is presented.

4.4.3 Policy Training. After the recommendation and policy are pre-trained, we further train two policy modules with online user interaction meanwhile the recommendation module is fixed. Concretely, we utilize each session collected online to optimize these two policy networks by Policy Gradient [30] following [12, 34].

5 EXPERIMENTS

In this section, we conduct experiments on the MSMCR scenario to evaluate the performance of our method compared with other adapted state-of-the-art (SOTA) models¹.

Table 1: Statistics of datasets used in experiments, where n is the length of one session.

Statistic	#Users	#Items	#Interactions	#Attributes	Avg. n
LastFM/LastFM*	1,801	7,432	76,693	33/8438	2.82/2.82
Yelp/Yelp*	27,675	70,311	1,368,606	29/590	2.85/2.85

5.1 Experimental Setup

5.1.1 Dataset. To evaluate the proposed method, we adopt four existing MCR benchmark datasets. The details of these datasets are given in Table 1.

- **LastFM and LastFM*** [18]. The LastFM dataset is used in assessing music artist recommendations. To facilitate modeling, original features are merged into 33 coarse-grained groups. Considering the unfeasibility of utilizing knowledge graphs to manually merge attributes in real situation, the initial features are used in the LastFM* dataset.

- **Yelp and Yelp*** [18]. The Yelp dataset is for the business recommendation field, involving two-tier taxonomy with 29 first-tier categories [17]. While Yelp* is a fine-grained version with 590 second-tier category attributes.

5.1.2 User Simulator. Due to the difficulty of interacting with real users, we design a user simulator to train and evaluate MCR frameworks. Following previous works [8, 17, 18], the classical user simulator is adapted to our MSMCR scenario from the construction of a session and the generation of activation attributes. Note that this paper mainly focuses on the scenario itself and its solution, and a more reasonable user simulator can be designed for MSMCR.

- **The Construction of a Session.** We simulate a session composed of ordered subsessions from a chronological user-item interaction sequence. i -th subsession s_u^i is constructed by i -th interacted item v^i of user u correspondingly. For each subsession s_u^i in session S_u , we regard the item v^i as the ground-truth target item, and its attribute set \mathcal{A}_{v^i} as the oracle attributes preferred by the user. The session length n is defined as $\min(\text{Random}(N_{min}, N_{max}), M_u)$, where M_u is the interaction length of user u , N_{min} and N_{max} are the threshold values of session length n to align with realistic conversation session size. We report the mean value of n resulted from our simulation in Table 1.

- **The Generation of Activation Attributes.** The activation attributes $\mathcal{A}_{v^n}^*$ of the current subsession within a session are the subset of oracle attributes \mathcal{A}_{v^n} . Specifically, it is defined as the top-two-ranked attributes based on user-attribute affinity score: $\mathbf{U}_u^T \mathbf{A}_a + (\sum_{j=1}^{n-1} \mathbf{V}_{v^j}^T \mathbf{A}_a) / (n-1)$, $\forall a \in \mathcal{A}_{v^n}$ for each session, following the design of related works [2, 5]. In this formula, v^j represents the target item of the j -th subsession in session S_u , and $\mathbf{U}, \mathbf{V}, \mathbf{A}$ are pre-trained embeddings of users, items and attributes, respectively, which are obtained via TransE [1] in [8]. Note that these activation attributes are applied uniformly across all methods, and their validity will be confirmed in Section 5.5.

5.1.3 Baseline Models. To verify model performance, our proposed model is compared with the following five representative multi-round CRS methods.

- **Max Entropy (MaxE)** [17] chooses to ask the attribute with maximum entropy among the candidate items or recommend the top-ranked items with a certain probability.
- **EAR** [17] proposes an Estimation-Action-Reflection three-stage solution, which achieves more accurate results than traditional algorithms like Abs Greedy [3] and CRM [29].
- **SCPR** [18] is designed to use interactive path reasoning on the graph to trim candidate attributes and adopt the DQN [23] framework for action selecting.
- **FPAN** [34] extends EAR by adapting item and user embedding via the gate mechanism based on online user feedback.
- **UNICORN (UNI)** [8] constructs a unified policy learning framework based on a dynamic weighted graph.

As discussed in Section 2, several recent works (e.g., MIMCR [39], Bundle-MCR [12], GPR [41], and MetaCRS [5]) primarily focus on the different novel MCR scenarios, which go beyond the scope of our MSMCR setting. Therefore, we do not adopt them. Also, we do not use the recent work CRIF [14], since main contributions of CRIF are based on Inverse Reinforcement Learning [4] for conversation

¹Our code and data will be released for research purposes.

Table 2: Experimental results on four datasets. The bold number indicates the improvements over the best baseline (underlined) are statistically significant ($p < 0.01$) with paired t -test. SR, hN and AR stand for SR@10, hN@(10, 10) and AR@10, respectively. All the baselines with \dagger and \ddagger -superscript are adapted to our MSMCR setting (cf. Section 5.1.3 for more details).

Models	LastFM			LastFM*			Yelp			Yelp*						
	SR \uparrow	AT \downarrow	hN \uparrow	AR \uparrow	SR \uparrow	AT \downarrow	hN \uparrow	AR \uparrow	SR \uparrow	AT \downarrow	hN \uparrow	AR \uparrow	SR \uparrow	AT \downarrow	hN \uparrow	AR \uparrow
MaxE \ddagger	0.061	9.94	0.015	0.451	0.067	9.91	0.017	0.199	0.846	6.52	0.278	0.953	0.015	9.98	0.004	0.324
EAR \ddagger	0.117	9.78	0.037	0.254	0.059	9.94	0.015	0.197	0.804	6.91	0.273	0.956	0.054	9.95	0.009	0.353
SCPR \ddagger	0.112	9.77	0.036	0.262	0.072	9.88	0.023	0.208	0.782	6.89	0.272	0.977	0.058	9.94	0.015	0.362
SCPR \ddagger	0.128	9.69	0.042	0.273	0.075	9.82	0.027	0.248	0.826	6.85	0.301	<u>0.985</u>	0.091	9.84	0.023	0.371
FPAN \ddagger	0.167	9.74	0.039	0.266	0.067	9.90	0.022	0.214	0.880	5.94	0.297	0.966	0.071	9.91	0.017	0.354
FPAN \ddagger	0.264	9.54	0.052	0.313	0.069	9.89	0.028	<u>0.365</u>	<u>0.895</u>	<u>5.87</u>	0.304	0.972	0.082	9.88	0.022	0.361
UNI \ddagger	0.243	9.58	0.050	0.413	0.158	9.64	0.042	0.129	0.805	6.57	0.289	0.950	0.124	9.80	0.032	0.319
UNI \ddagger	<u>0.286</u>	<u>9.48</u>	<u>0.068</u>	<u>0.475</u>	<u>0.196</u>	<u>9.45</u>	<u>0.062</u>	0.178	0.830	6.42	<u>0.312</u>	0.954	<u>0.162</u>	<u>9.76</u>	<u>0.052</u>	<u>0.382</u>
Ours	0.547	8.88	0.170	0.855	0.291	9.47	0.089	0.439	0.924	5.81	0.349	0.968	0.277	9.75	0.081	0.778

policy learning, which differs from the general fixed reward setting in other baselines and our MSCAA (cf. Section 5.1.5).

For a fair comparison, we revise all the above baselines as follows: (1) We utilize previous subsessions information to complement the user presentation for the item and attribute scoring. Specifically, we first apply GRU to obtain the representation of user item-level interest from previous subsessions, aligned with our method (cf. Section 4.1.2). Then, the original user representation and item-level interest presentation are further fused as the updated user representation via the gating mechanism [22, 25]. (2) We employ our user simulator for all baselines.

Note that all baseline implementations are modified to adapt to the MSMCR scenario, and therefore these methods are named with superscript \ddagger . In addition, to further verify the necessity of previous subsessions, we remove the implementation of (1) and only preserve the setting of activation attributes (2) for several baselines which are named with \ddagger -.

5.1.4 Evaluation Metrics. Following previous studies on MCR [17, 18], the accumulative ratio of successful conversational recommendation assessed by SR@ T (success rate at turn T) and the average number of turns for all subsessions assessed by AT (average turn) are adopted. Similar to SR@ T , we introduce AR@ T (activation rate at turn T) to measure the accumulative ratio of hitting activation attributes in the session. Moreover, we also adopt hN@(T, K) to measure the recommendation accuracy following [8], which considers the rank performance in both list- and turn-level. We omit the @ T and K terms in the following experiments for simplicity. The higher value of SR, AR, and hN indicates better performance, while a lower AT indicates higher effectiveness. Notably, SR, AT and our proposed AR only consider a certain aspect of CRS [8], while hN is a comprehensive metric for evaluating the overall framework.

5.1.5 Implementation Details. We implement the proposed method based on PyTorch. We randomly divide interactions across each dataset into training, validation, and test parts with the ratio of 7 : 1.5 : 1.5 and then separately generate the session samples based on by user simulator described in Section 5.1.2. The maximum/minimum session length N_{max}, N_{min} is set to 4, 2 respectively. We set the maximum turn T and the size of the recommendation list K as 10. The embedding dimension D is set as 64. We employ the Adam/SGD optimizer to train the recommendation/policy module with the learning rate $5e^{-4}$ and $1e^{-3}$ separately. The heterogeneous

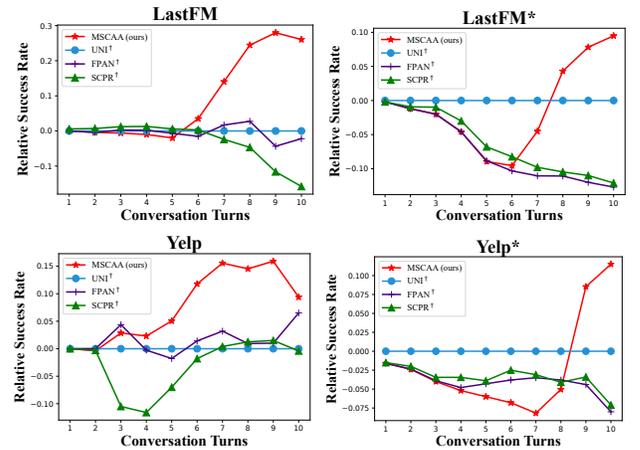


Figure 3: Comparisons of Relative Success Rate at Different Conversation Turns.

graph G is constructed by the training dataset. The number of RGCN layers L_g and Transformer layers L_t are all set to be 2. Discount factor γ is set to be 0.7. The temperature parameter τ is set with 0.5. The weight ω of InfoNCE constraint is set to be 0.01. We set six rewards $[r_{rec_acc}, r_{rec_rej}, r_{ask_acc}, r_{ask_rej}, r_{ask_unk}, r_{quit}] = [1, -0.1, 0.01, -0.1, -0.1, -0.3]$ following reward settings in [18].

5.2 Performance Comparison

5.2.1 Overall Comparison. To demonstrate the overall performance of the proposed model, we compare it with the adapted SOTA methods for the MSMCR scenario. Table 2 shows the comparison of experimental results between the baseline models and our proposed method. We can obtain the following notable observations:

- In most cases, our proposed method exhibits significant advantages among all baselines on four datasets. The main reason behind this is that MSCAA carefully combines three dedicated modules, especially the recommendation module that comprehensively models different user interests from heterogeneous information sources, to hit the user desired item in a few turns of subsession.
- Interestingly, several baselines (e.g., SCPR \ddagger and UNI \ddagger) sometimes outperform MSCAA in AR and AT metrics. It could be attributed to the policy bias in these models that prefer to ask or recommend. For instance, more asking times will result in better accumulative

performance on AR. However, our method still outperforms them in the comprehensive metric (i.e., hN) with a scalability trade-off between asking and recommendation.

- Our method’s improvement is relatively limited on the Yelp dataset among all four datasets. In addition, the performance of both LastFM* and Yelp* is comparatively inferior to the other two datasets. These could be explained that the setting of the Yelp dataset is to ask enumerated questions [17], which makes it easier to hit the activation attributes and then sharply diminish the candidate item space for all methods. However, for the LastFM* and Yelp* datasets, it is harder to ask appropriate attributes to activate user interests due to the numerous candidate attributes.
- Compared to “†-”, “†” version method achieves a better overall performance, which indicates that it is effective to utilize the information of previous subsessions via sequential modeling in the MSMCR scenario.

5.2.2 Comparison at Different Conversation Turns. Besides SR@10, we also present a detailed performance comparison of success rate at each turn in Figure 3. To facilitate the observation, we only report the relative success rate of four representative methods compared with the SOTA baseline UNI† in our MSMCR setting. For example, in Figure 3, the blue line of UNI† is set to $y = 0$. As can be seen, UNI† or even other baselines may outperform MSCAA in the first few turns, while our MSCAA achieves a pretty performance at the end stage of the conversation. This phenomenon can be explained that MSCAA is more inclined to ask appropriate attributes to elicit user interest at first, and after that, can more confidently recommend items to satisfy the user’s current demand. Conversely, the other baselines may recommend items successfully instead of asking more attributes in the earlier turns since they do not consider the activation of user interest. However, due to the lack of acquisition of the current user interest, they are inferior to MSCAA lastly.

5.3 Ablation Study

To verify the design effectiveness of MSCAA, we remove or replace key modules one-by-one and report results in Table 3.

Impact of recommendation modules. The results in (a-c) show that the missing of any type of information (i.e., historical interactions, previous subsessions, or the online user feedback information) causes the degradation of model performance. Especially, the long-term interest is the most significant among these three for highlighting the personality of users. Moreover, (d) demonstrates the significance of contrastive loss for the recommendation. (e) suggests that unified graph modeling is pivotal for long-term interest extraction. (f) represents the GRU modules in Section 4.1.2 are replaced with “mean” operation, which indicates the importance of sequential encoder of previous subsessions for recommendation in our MSMCR scenario. Based on the above observations, we manifest that any input information and main components are indeed essential in the context-aware recommendation module.

Impact of Attribute and Conversation Policy. (g) and (h) indicate that our attribute selection policy is replaced by Max Entropy and Max User-like score rule, respectively. The results show the effectiveness of our attribute agent for selecting an appropriate attribute based on the conversation state. And (i) means the fine-grained score information $s_{item} \oplus s_{attr}$ in conversation state is

Table 3: Results of the Ablation Study.

Models	LastFM			Yelp*		
	SR↑	hN↑	AR↑	SR↑	hN↑	AR↑
(a) w/o Long-term Interest	0.255	0.078	0.490	0.065	0.019	0.450
(b) w/o Short-term Interest	0.463	0.143	0.844	0.146	0.043	0.550
(c) w/o Current Interest	0.537	0.167	0.845	0.178	0.052	0.720
(d) w/o InfoNCE Loss	0.538	0.168	0.852	0.265	0.078	0.759
(e) w/o RGCN	0.517	0.161	0.838	0.261	0.076	0.760
(f) w/ Mean	0.537	0.166	0.853	0.262	0.077	0.764
(g) w/ Max Entropy	0.102	0.030	0.333	0.032	0.009	0.438
(h) w/ Max User-like	0.422	0.132	0.851	0.138	0.041	0.775
(i) w/o $s_{item} \oplus s_{attr}$	0.538	0.169	0.848	0.270	0.079	0.770
MSCAA	0.547	0.170	0.855	0.277	0.081	0.778

masked. From the result, this kind of information is necessary for a high-level action decision.

5.4 Case Study

To show the process of activating user interest from the system side and then recommending successfully, we further present a session case generated by our framework MSCAA and UNI† based on the same test instance from Yelp* dataset in Figure 4. This session consists of three subsessions where in the previous two subsessions, the user u_{12019} has already obtained his/her desired items *The Melting Pot* and *Wigle Whiskey* which are related to food and drink in *city_17*, and he/she continues the third subsession to seek more advice. For the current subsession 3, two main attributes *city_17* and *Restaurants* of the target item *Smallman Galley* are considered as the activation attributes of this dialogue.

It can be observed that based on the comprehensive information especially the previous subsessions, our MSCAA precisely hits the activation attribute *city_17* to elicit user interest for subsequent specific attribute-asking questions. As the subsession turn progresses, the user demand is gradually crystallized for assisting the system in seeking a satisfactory item (also a restaurant) successfully. While the baseline UNI† recommends items at the beginning of the subsession since it utilizes user interest from previous subsessions to be confident to choose this action. However, it fails because of the missing of current user interest in this subsession. And then, UNI† can present an attribute *price_2* that satisfies the user to a certain extent, but this over-specific question may confuse the user before his/her current interest is activated. In summary, compared to UNI†, MSCAA can not only ask appropriate attributes to elicit user interest quickly in the early turns but make a comprehensive decision of recommendation or asking, beneficial for succeeding.

5.5 Human Evaluation

Although the above experiments demonstrate the effectiveness of our framework in the proposed user settings, we further conduct a human evaluation to answer two questions: (1) Are the activation attributes appropriate in our user simulator? (2) Does our method outperform other baselines in the realistic scenario? Specifically, we randomly select 100 session samples from Yelp* from our framework and the competitive baseline UNI†. The experiment involves 20 post-graduate volunteers who evaluate 10 samples each, with each sample being evaluated by 2 different volunteers. For the first question, volunteers need to review the user’s previous subsessions

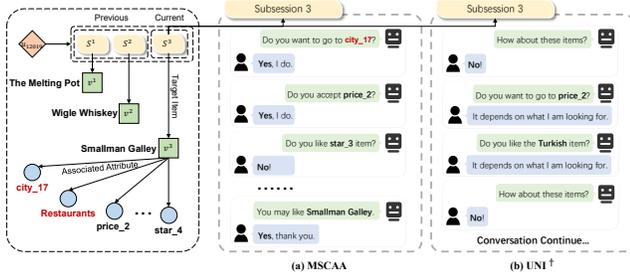


Figure 4: A session case generated by our framework MSCAA and UNI†. The left side illustrates the sampled user and his/her interacted three subsessions with corresponding target items and associated attributes, in which activation attributes are red bold fonts. We present the conversation process of the current subsession (i.e., subsession 3) in detail.

and choose 2 main attributes of the target item in the current subsession as human-labeled activation attributes. We then calculates Jaccard Similarity [19] for each sample between human-labeled and our simulated activation attribute set, and the average value is 0.73. For the second question, the volunteers are required to browse the conversation process of the current subsession and select the better conversation from $\langle \text{MSCAA}, \text{UNI}^\dagger \rangle$ episode pairs in their subjective view, following the human evaluation method in [12]. Lastly, we collect 200 results: $\langle \text{MSCAA}, \text{UNI}^\dagger \rangle$ votes are 127 : 73. Overall, these results validate the quality of our user simulator and the effectiveness of our method in the realistic scenario.

6 CONCLUSION

In this paper, we extend the MCR to a general CRS setting, MSMCR in which the user continues the dialogue with the system after several subsessions and might preserve no clear interest in the current subsession, and the system would proactively take actions to activate the user’s dynamic interest. For this scenario, a novel framework called MSCAA is introduced to model user interests comprehensively for the recommendation, learn the flexible strategy for asking the appropriate attributes, and manage actions between asking and recommending adaptively. Extensive experimental results on four adapted datasets verify the effectiveness and superiority of our framework in the proposed scenario.

REFERENCES

- [1] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. *NeurIPS* (2013).
- [2] Bo Chen, Wei Guo, Ruiming Tang, Xin Xin, Yue Ding, Xiuqiang He, and Dong Wang. 2020. TGCN: Tag graph convolutional network for tag-aware recommendation. In *CIKM*.
- [3] Konstantina Christakopoulou, Filip Radlinski, and Katja Hofmann. 2016. Towards conversational recommender systems. In *SIGKDD*.
- [4] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *NIPS* (2017).
- [5] Zhendong Chu, Hongning Wang, Yun Xiao, Bo Long, and Lingfei Wu. 2023. Meta Policy Learning for Cold-Start Conversational Recommendation. *WSDM*.
- [6] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).
- [7] Chuan Cui, Qi Shen, Shixuan Zhu, Yitong Pang, et al. 2022. Intention Adaptive Graph Neural Network for Category-aware Session-based Recommendation. In

DASFAA.

- [8] Yang Deng, Yaliang Li, Fei Sun, Bolin Ding, and Wai Lam. 2021. Unified Conversational Recommendation Policy Learning via Graph-based Reinforcement Learning. In *SIGIR*.
- [9] Chongming Gao, Wenqiang Lei, Xiangnan He, Maarten de Rijke, and Tat-Seng Chua. 2021. Advances and challenges in conversational recommender systems: A survey. *AI Open* (2021).
- [10] ASM Ahsan-Ul Haque and Hongning Wang. 2022. Rethinking Conversational Recommendations: Is Decision Tree All You Need?. In *CIKM*.
- [11] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *SIGIR*.
- [12] Zhankui He, Handong Zhao, Tong Yu, Sungchul Kim, Fan Du, and Julian McAuley. 2022. Bundle MCR: Towards conversational bundle recommendation. In *RecSys*.
- [13] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2016. Session-based recommendations with recurrent neural networks. *ICLR* (2016).
- [14] Chenhao Hu, Shuhua Huang, Yansen Zhang, and Yubao Liu. 2022. Learning to Infer User Implicit Preference in Conversational Recommendation. In *SIGIR*.
- [15] Heeseon Kim, Hyeonjun Yang, and Kyong-Ho Lee. 2023. Confident Action Decision via Hierarchical Policy Learning for Conversational Recommendation. In *WWW*.
- [16] Wenqiang Lei, Xiangnan He, Maarten de Rijke, and Tat-Seng Chua. 2020. Conversational recommendation: Formulation, methods, and evaluation. In *SIGIR*.
- [17] Wenqiang Lei, Xiangnan He, Yisong Miao, Qingyun Wu, Richang Hong, Min-Yen Kan, and Tat-Seng Chua. 2020. Estimation-action-reflection: Towards deep interaction between conversational and recommender systems. In *WSDM*.
- [18] Wenqiang Lei, Gangyi Zhang, Xiangnan He, Yisong Miao, Xiang Wang, Liang Chen, and Tat-Seng Chua. 2020. Interactive path reasoning on graph for conversational recommendation. In *SIGKDD*.
- [19] Michael Levandowsky and David Winter. 1971. Distance between sets. *Nature* (1971).
- [20] Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018. Towards deep conversational recommendations. *NIPS* (2018).
- [21] Shijun Li, Wenqiang Lei, Qingyun Wu, Xiangnan He, Peng Jiang, and Tat-Seng Chua. 2021. Seamlessly unifying attributes and items: Conversational recommendation for cold-start users. *TOIS* (2021).
- [22] Chen Ma, Peng Kang, and Xue Liu. 2019. Hierarchical gating networks for sequential recommendation. In *KDD*.
- [23] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, et al. 2015. Human-level control through deep reinforcement learning. *nature* (2015).
- [24] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748* (2018).
- [25] Yitong Pang, Lingfei Wu, Qi Shen, Yiming Zhang, Zhihua Wei, et al. 2021. Heterogeneous Global Graph Neural Networks for Personalized Session-based Recommendation. In *WSDM*.
- [26] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).
- [27] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. 2011. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*.
- [28] Michael Schlichtkrull, Thomas N Kipf, Peter Bloem, Rianne Van Den Berg, and other. 2018. Modeling relational data with graph convolutional networks. In *ESWC*.
- [29] Yueming Sun and Yi Zhang. 2018. Conversational recommender system. In *SIGIR*.
- [30] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*.
- [31] Xintao Tian, Yongjing Hao, Pengpeng Zhao, Deqing Wang, Yanchi Liu, and Victor S Sheng. 2021. Considering Interaction Sequence of Historical Items for Conversational Recommender System. In *DASFAA*. Springer.
- [32] Quan Tu, Shen Gao, Yanran Li, Jianwei Cui, Bin Wang, and Rui Yan. 2022. Conversational Recommendation via Hierarchical Information Modeling. In *SIGIR*.
- [33] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *NIPS* (2017).
- [34] Kerui Xu, Jingxuan Yang, Jun Xu, Sheng Gao, Jun Guo, and Ji-Rong Wen. 2021. Adapting User Preference to Online Feedback in Multi-round Conversational Recommendation. In *WSDM*.
- [35] Liu Yang, Hamed Zamani, Yongfeng Zhang, Jiafeng Guo, and W Bruce Croft. 2017. Neural matching models for question retrieval and next question prediction in conversation. *arXiv preprint arXiv:1707.05409* (2017).
- [36] Feng Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. 2016. A dynamic recurrent model for next basket recommendation. In *SIGIR*.
- [37] Yongfeng Zhang, Xu Chen, Qingyao Ai, Liu Yang, and W Bruce Croft. 2018. Towards conversational search and recommendation: System ask, user respond.

- In *CIKM*.
- [38] Yinan Zhang, Boyang Li, Yong Liu, Hao Wang, and Chunyan Miao. 2022. Minimalist and High-performance Conversational Recommendation with Uncertainty Estimation for User Preference. *arXiv preprint arXiv:2206.14468* (2022).
 - [39] Yiming Zhang, Lingfei Wu, Qi Shen, Yitong Pang, Zhihua Wei, Fangli Xu, Bo Long, and Jian Pei. 2022. Multiple Choice Questions based Multi-Interest Policy Learning for Conversational Recommendation. In *WWW*.
 - [40] Canzhe Zhao, Tong Yu, Zhihui Xie, and Shuai Li. 2022. Knowledge-aware Conversational Preference Elicitation with Bandit Feedback. In *WWW*.
 - [41] Rongmei Zhao, Shenggen Ju, Jian Peng, Ning Yang, Fanli Yan, and Siyu Sun. 2022. Two-Level Graph Path Reasoning for Conversational Recommendation with User Realistic Preference. In *CIKM*.
 - [42] Sen Zhao, Wei Wei, Yifan Liu, Ziyang Wang, Wendi Li, Xian-Ling Mao, Shuai Zhu, Minghui Yang, and Zujie Wen. 2023. Towards Hierarchical Policy Learning for Conversational Recommendation with Hypergraph-based Reinforcement Learning. In *IJCAI*.
 - [43] Kun Zhou, Wayne Xin Zhao, Hui Wang, Sirui Wang, Fuzheng Zhang, Zhongyuan Wang, and Ji-Rong Wen. 2020. Leveraging historical interaction data for improving conversational recommender system. In *CIKM*.
 - [44] Kun Zhou, Yuanhang Zhou, Wayne Xin Zhao, Xiaoke Wang, and Ji-Rong Wen. 2020. Towards Topic-Guided Conversational Recommender System.. In *ICML*.
 - [45] Jie Zou, Yifan Chen, and Evangelos Kanoulas. 2020. Towards question-based recommender systems. In *SIGIR*.

progress needs to be made in areas like Sim2Real (Simulation to Reality) before we deploy this system in a production setting.

ETHICAL CONSIDERATIONS

It is believed that our proposed MSMCR scenario and the corresponding solution framework MSCAA could mitigate some challenging ethical problems in several conventional recommender systems (e.g., collaborative filtering, sequential recommendation, etc.). The MSMCR scenario, which inherits but surpasses the existing MCR setting, has the following merits:

- MCR emphasizes the dynamic and precise acquisition of user preferences during *online* conversations, which is ideally less reliant on learning from offline interactions that might be noisy and bring about unintentional biases (e.g., popularity bias). Furthermore, our MSMCR additionally considers the *previous subsessions* of one session for users, comprehensively excavating online information except for feedback in the current subsession of the session. As such, MSMCR could further mitigate the bias of offline interaction data.
- Our MSMCR also takes more realistic characteristics of users into consideration. Specifically, users might continue to engage in the conversation session aimlessly until a raised topic triggers user interests. Therefore, the *activation attributes* of users are designed in the subsequent subsession of one session, so as to resemble real-world situations for satisfying user demands.

As for our solution, MSCAA framework is established on reinforcement learning, therein the *attribute selection* and *conversation policy* is learned by policy learning. Different from traditional predictive and matching methods, policy learning aims to maximize cumulative gains for all conversations, enabling to cater to long-term user satisfaction (i.e., the successful conversation). This goal is more user-friendly to enhance the overall user experience.

However, due to the difficulty of interacting with real users, collecting online user feedback is too ideal to train MSCAA for alleviating the unintentional biases in the offline interaction data. Hence, following the previous related works, we still adopt a adapted user simulator based on user interaction data, including the construction of one session, the generation of activation attributes and the simulation of user feedback. Indeed, it is undeniable that using such a simulated environment introduces potential discrepancies and biases compared to the real-world distribution. Significant further

A MAIN NOTATIONS

Table 4: Main notations in the paper.

Notations	Definitions and Descriptions
u	A user
\mathcal{U}	The set of all users
v	An item
\mathcal{V}	The set of all items
a	An attribute
\mathcal{A}	The set of all attributes
\mathcal{A}_v	The set of attributes that belong to an item v
s	A subsession
S_u	A session interacted by user u
P_{S_u}	The previous subsessions of session S_u
$P_{S_u}^{\mathcal{V}}$	The sequence of the target items of previous subsessions in the session S_u
$P_{S_u}^{\mathcal{A}}$	The sequence of the target attribute sets of previous subsessions in the session S_u
s_u^n	The current (i.e., n -th) subsession of session S_u , where n is the length of the session
v^n	The target item of subsession s_u^n (also session S_u)
$\mathcal{A}_{v^n}^*$	The activation attributes of subsession s_u^n (also session S_u)
\mathcal{A}_{acc}	The set of attributes accepted by a user in the current subsession
\mathcal{A}_{rej}	The set of attributes rejected by a user in the current subsession
\mathcal{A}_{cand}	The set of candidate attributes in the current subsession
\mathcal{V}_{cand}	The set of candidate items in the current sub-session
\mathcal{V}_{rej}	The set of items rejected by a user in the current sub-session
s_a	The state of attribute selection policy
s_c	The state of conversation policy
π_c	The conversation policy
π_a	The attribute selection policy

B FURTHER DISCUSSION

B.1 Differences between MSMCR and Existing Scenarios

Our MSMCR scenario aims to activate user interests in the current subsession for the successful recommendation by considering the correlation between previous subsessions within a session.

We highlight that it differs from two related recommendation fields: (1) sequential/session recommendation [7, 25], which lacks interactive dialogue between users and the system and only provides item recommendation once; (2) MCR [8, 17, 18], which disregards the previous subsessions and activation attributes.

B.2 The Strength and Weakness on the Simulator

Regarding the strengths, it is commendable that our proposed user simulator can handle more general behavior and preference of users (e.g., the recent subsessions imply short-term user interest), meanwhile providing more realistic feedback (e.g., “unknown” when vague interest) than previous CRS simulators [8, 17, 18].

As for the weaknesses, it is understandable that the user simulator’s reliance on activation attribute extraction may lead to some deviation in the distribution of simulated user feedback from reality. Specifically, (1) some attributes that users really like may not have been generated, resulting in the user simulator lacking more explicit responses; (2) some noisy attributes may have been mixed in with certain users’ preference attributes, causing the user simulator to conduct wrong dialogues for activating current user interest. Additionally, the construction of a session is based on the chronological user-item interaction sequence, not consecutive online dialogues between users and the system, which may also bring deviations that are somewhat inconsistent with reality. As for the setting of hyper-parameters on conversations (cf. Section 5.1.5), although we follow the existing CRS works [5, 10] to set the maximum turn T as a fixed value 10, trying to be as realistic as possible [35], this is too ideal to introduce deviations inevitably. The same is for the maximum/minimum session length.