# An Immersing Oriented Role-Playing Framework with Duplex Relationship Modeling

Yuntao Wen, Shen Gao, Chengrui Huang, Yifan Wang, and Shuo Shang[✉]

University of Electronic Science and Technology of China, Chengdu, China
{shengao, shuoshang}@uestc.edu.cn

**Abstract.** Role-playing is an emerging application of large language models (LLMs), allowing users to be immersed in conversations with virtual characters by mimicking their tones and background knowledge. It can be applied in various scenarios such as gaming and virtual reality systems. However, existing methods ignore two challenges: (1) ignoring the relationship with the role played by the user will diminish the immersive experience of the user; (2) insufficient understanding of the character's background knowledge may lead to inconsistent dialogue. In this paper, we introduce the **D**uplex **R**elationship **M**odeling based **R**ole-play framework (DRMR), a novel role-playing framework designed to enhance the immersion of user when interacting with the role-play model. We first propose a graph-based relationship modeling method, utilizing graph structures to model the duplex relationship between the user and the model's played characters. In order to better extract useful personalized information about roles from historical dialogues, we construct a role memory consisting of the description of the duplex relationship. To avoid generating an inconsistent response, we iteratively verify the generated response by updating the role memory according to the current dialogue context. Extensive experiments on benchmark dataset demonstrate the effectiveness of DRMR in enhancing user immersion in role-playing interactions.
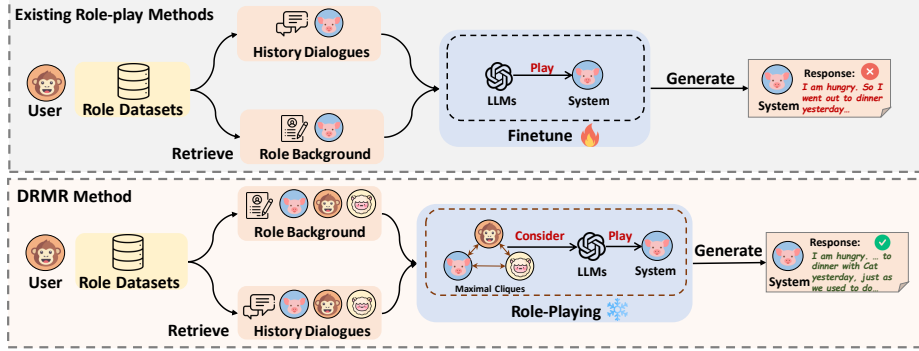
**Keywords:** Dialogue system · Role-Play · Large language model

## 1 Introduction

In recent years, large language models (LLMs) have made significant advancements in numerous classical natural language processing tasks [6, 20, 26, 31]. This has also brought several new paradigms in natural language processing, transitioning gradually from better accomplishing traditional natural language tasks to some new applications such as tool usage [7, 8, 17, 37], LLM-based multi-agent systems [5, 3, 16], embodied intelligence methods for manipulating robots [14, 28, 23] and role-playing [4, 12, 25]. Role-playing aims to enable LLMs to portray specific characters/roles[1] (*e.g.,* characters in movies and TV dramas,

---

[1] These two terms are interchangeably used.

**Fig. 1.** Comparison between existing role-play methods and our proposed DRMR. Most previous methods usually annotate large amounts of data and then fine-tune the LLM, and they typically consider only the information of the role played by the model, neglecting the duplex relationship information between the roles played by the user.

historical figures, etc.) to meet user needs. These methods have been widely used in interactive games [13, 30], virtual reality systems [18], and psychological counseling [33, 9].

On one hand, some of the existing role-playing methods [12, 4, 36, 25] focus on fine-tuning LLMs by either constructing more role-playing datasets or data augmentation. This enables large models to understand the background knowledge and language style characteristics of roles, thus achieving better role imitation. However, this not only relies on acquiring a large amount of data but also considerable training time and GPU resources for fine-tuning LLM. On the other hand, some methods [32, 35, 29] attempt to achieve this by allowing users to define role profiles as in-context instructions, but this requires lengthy input from users to define roles, which adversely affects user experience.

In real-world applications, users will provide only a brief role profile and several previous dialogues for the role-play model. Intuitively, since not all the details of the role can be comprehensively defined in the profile, models often struggle to generate consistent responses, such as an ancient figure writing code. This deviation from the background of the character in responses also diminishes the user immersion. Therefore, the **first challenge** lies in deeply understanding the brief profile and making full use of the given data to generate dialogues that are consistent with the background of the character.

Furthermore, the majority of existing methods only incorporate personalized information about the role played by the model (*a.k.a.,* **simplex** relationship), ignoring the role profile and relationship played by the user (*a.k.a.,* **duplex** relationship). However, an immersive role-playing experience requires not only mimicking the tone and knowledge background of the character being played but also involving the user in the scenario where the character is situated. It is crucial for role-playing models to understand the duplex relationship between the character played by the user and the character played by the model, as

this greatly contributes to the immersive experience of role-playing. Thus, the **second challenge** lies in how to model the interpersonal relationships between the roles played by the user and the model when role-playing.
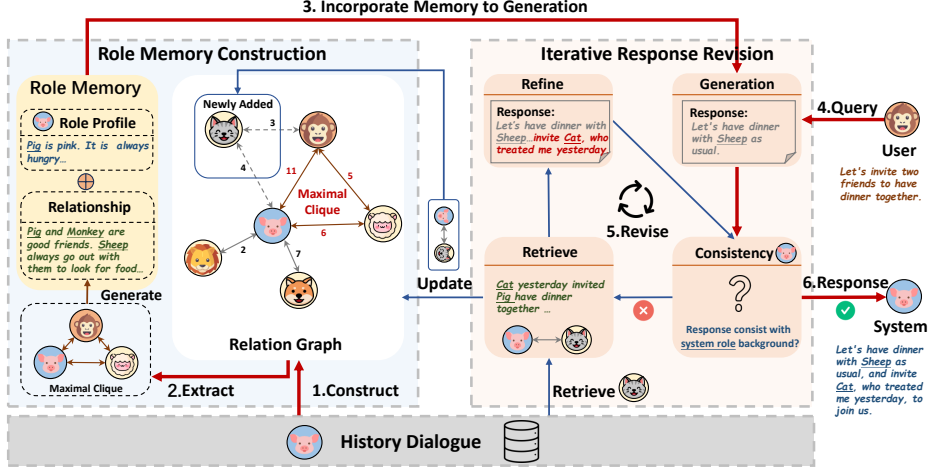
In this paper, we propose the **D**uplex **R**elationship **M**odeling based **R**ole-play framework (DRMR) method, a role-playing framework aimed at enhancing immersion of user experience. Given a brief role profile provided by the user and several historical dialogues, our approach employs two novel methods to enhance understanding of the character's background and achieve duplex relationship modeling for both the model and the user-played roles. To achieve duplex relationship modeling, we propose a *maximal-cliques-based role relationship modeling* method based on a role relation graph. By using the maximal cliques representing both the model and the user's played characters along with their shared background information, we construct a role memory to summarize the useful relationship information, thereby enhancing user immersion.

Then we introduce an *iterative response revision* method, which iteratively revises the model responses by retrieving more related dialogues and updating the role memory, thus generating responses that align with the background of the character. Extensive experiments conducted on a benchmark dataset demonstrate the effectiveness of our proposed DRMR, and we can find that our proposed model can enhance the immersion of the user when chatting with the role-playing system. Our contributions of this work are as follows:

• We propose DRMR, which is a role-playing framework to enhance the immersion of user experience.

• We introduce maximal-cliques-based role relationship modeling to incorporate the duplex relationship of both characters played by the user and the model.

• We propose the iterative response revision method which iteratively verifies the consistency of the response and revises the response by updating role memory.

• Experimental results on benchmark dataset illustrate the superiority of DRMR.

## 2   Related Work

Role-playing is an important application of LLMs, aimed at simulating a character comprehensively by using events from movies, TV shows, or historical figures to achieve immersive interaction with users. ChatHaruhi [12] is an earlier method that utilizes LLMs to implement role-playing which establishes a character dialogue database and introduces a retrieval-enhanced role-playing framework. Character-LLM [21] focuses on modeling character memories, reconstructing scene-based memories using WikiData, and adopting protective experiences to mitigate the hallucination of response. CharacterGLM [36] further develops a multi-turn role-playing dialogue system based on fine-tuning LLMs, using character profiles, dialogues, and a large amount of crowd-sourcing data as training dataset. HPD [4] is a dataset for playing the role of Harry Potter integrating extensive and detailed background information to better match LLMs with the characteristics of Harry Potter. RoleLLM [25] proposes a role-playing model

**Fig. 2.** Overview of **D**uplex **R**elationship **M**odeling based **R**ole-play framework (DRMR) which has three main steps: (1) We first construct a relation graph using the historical dialogues of the roles and extract the maximal cliques from the graph to build the role memory; (2) We generate the response by incorporating the role memory; (3) We employ the iterative response revision framework to verify the revise the response which ensures the response is consistent with the background of the role.

based on instruction tuning by maintaining specific knowledge and speaking tones of characters by combining in-context instructions.

However, the majority of existing role-playing methods require fine-tuning LLMs through annotating large datasets, which demands significant computational resources and data labeling efforts. Moreover, most existing works do not consider the interpersonal relationships between the characters portrayed by the user and the model, leading to model-generated responses that may not align with the current conversational context, thereby diminishing the immersive experience of the user.

## 3   DRMR Methodology

In this section, we detail the **D**uplex **R**elationship **M**odeling based **R**ole-play framework (DRMR). An overview of DRMR is shown in Figure 2.

### 3.1   Problem Formulation

Given the brief profiles $P_m$ and $P_u$ of the role $E_m$ to be played by the model and the role $E_u$ played by the user, along with several historical dialogues $D = \{(E_1, U_1), (E_2, U_2), \ldots, (E_L, U_L)\}$ as the input to our DRMR, where $E_i$ denotes the speaker of utterance $U_i$.

The user plays the role $E_u$ and engages in a dialogue of $T$ turn with the role $E_m$ played by the model, denoted as the current dialogue context $C = \{(c_1^u, c_1^m), (c_2^u, c_2^m), \ldots, (c_T^u, c_T^m), (c_{T+1}^u)\}$, where $c_i^u$ represents the $i$-th utterance of user, and $c_i^m$ represents the $i$-th response of the model. Based on this input, our model aims to generate responses $c_{T+1}^m$ of role $E_m$ to the user query $c_{T+1}^u$.

## 3.2   Role Memory Construction

When we engage in conversation with others, our minds not only contain information about ourselves but also the profile of the other person. And we also recall past experiences with this person (*e.g.,* travels together previously) and information about people associated with them (*e.g.,* their parents). Intuitively, it is crucial for humans to recall this information from memory during conversations which makes human-to-human dialogue natural; otherwise, conversations would become disjointed. Therefore, to enhance immersion in role-playing systems, we propose a *role memory M* to store relation information about the user role $E_u$, model role $E_m$, and other related roles. To construct the role memory $M$, we propose using the graph to explicitly model the relationship of the roles, and summarize the structures of the graph into natural language descriptions of the relationship between roles. The role memory $M$ contains several paragraphs describing the detailed profile of role $E_u$ and the relationship between role $E_u$ and $E_m$ to mimic the mind of people when chatting with others. Initially, we use the role profile $P_m$ as the initialization for the role memory $M$.

To recall the most related role information, we firstly utilize the dialogue context $C$ as query to retrieve $N$ dialogues from historical dialogues $D$ which are relevant to roles $E_u$ and $E_m$, denoted as $D^c = \{(E_1^c, U_1^c), (E_2^c, U_2^c), \ldots, (E_N^c, U_N^c)\}$. Specifically, we leverage the dense retrieval method as a semantic similarity measure to retrieve the most relevant dialogues $D^c$ from historical dialogues $D$:

$$\phi = \cos\left(\text{Emb}(U), \text{Emb}(C)\right), U \in D, \tag{1}$$

where $\phi$ is the similarity score. We employ the pre-trained LLM as the text embedding function Emb and use the cosine to measure the similarity between dialogue representations. Finally, we take the top-$N$ dialogues according to the score $\phi$ as the relevant dialogues $D^c$.

Then, we construct a user relationship graph $G$ containing the roles $E_u$ and $E_m$ as well as other related roles retrieved $\{E_1^c, E_2^c, \ldots, E_N^c\}$. When the two roles have conversations, an edge between these roles is added to the graph $G$. As the degree of association between roles varies, it is necessary to quantitatively measure the degree of association between roles when constructing the role relationship graph $G$. In this paper, we propose using LLM to evaluate the relationship weights between nodes:

$$s_{i,j} = \text{EdgeScore}(I_{\text{ES}}\{U_i^c, U_j^c\}) \in \{1, 2, 3, 4, 5\}, \tag{2}$$

where $s_{i,j}$ indicates the relation score between the role $E_i$ and $E_j$, $\{U_i^c, U_j^c\}$ represents the historical dialogue between the role $E_i$ and $E_j$, and $I_{\text{ES}}$ denotes the instruction to prompt the LLM to score the relationship between two roles:

You are a Character Event Assessment Assistant. Please carefully evaluate and score, reflecting the importance of the characters $\{E_i$ and $E_j\}$ in the following dialogue. Your scoring range is from 1 to 5...
[history dialogue $\{U_i^c, U_j^c\}$]
Refer to the following standards for scoring:
1 point: The character barely participates in the event, having no impact on its development.... Please provide a brief explanation for your score, assessing the importance of $\{E_i$ and $E_j\}$ based on the above standards.

Due to the limited context length of the LLM, it is not feasible to consider all the information of nodes and edges in a single dialogue turn. In the graph, since maximal cliques can represent a subset of vertices in a graph where every two distinct vertices are adjacent, providing a dense connection indicative of a strong relationship or relevance among the included vertices. Thus, we employ a relation maximal clique algorithm on graph $G$ to obtain a subgraph $G'$ containing a maximal clique comprising several roles most relevant to the roles $E_u$ and $E_m$:

$$G' = \underset{G' \in G^\star}{argmax} \sum_{i,j \in G'} s_{i,j}, where\, u, m \in G' \tag{3}$$

where $G^\star$ is the set of the maximal cliques. Since the maximal cliques are not always unique in graph $G$, we utilize the sum of relationship weights within the maximal clique as a selection criterion.

Subsequently, we utilize the role relation contained in the subgraph $G'$ to expand the role memory $M$. Each edge in the maximal clique subgraph $G'$ represents a dialogue between two roles. We retrieve the top $K$ dialogues $D'$ most relevant to the current dialogue $C$ from edges in $G'$ to update the role memory $M$. And we use the same retrieval method as in Equation 1. In order to enable the role-playing model to better understand the relationships between these relevant roles and incorporate these relationships into dialogue generation, we utilize LLM to extract descriptions of relationships between characters from these relevant dialogues $D'$ and summarize the events described in the dialogue:

$$m = \text{MemBuild}(D', P_m, P_u), \tag{4}$$

where $m$ is a new role memory record (*a.k.a.,* a paragraph that describes the detailed relation between two roles). And the operation MemBuild is a chain-of-thought [27] based prompting method that prompts the LLM to summarize the relationship between two roles is as follows:

You will play a role that depicts your relationship with another character through a series of events that have occurred...
First, you will play $\{E_m\}$... Next, briefly describe your relationship with $\{E_u\}$ from $\{E_m\}$'s first-person perspective. Third, ...
To assist you in this task, here are some events that have occurred between $\{E_m\}$ and $\{E_u\}$: $\{D'\}$
Please output $\{E_m\}$'s first-person evaluation of $\{E_u\}$. The description should be concise and relevant.

Finally, we append the new memory record $m$ into the role memory $M$.

### 3.3   Iterative Response Revision

Based on the role memory $M$ and dialogue context $C$, we prompt LLM to generate responses $c_{T+1}^m$ for the role $E_m$:

$$c_{T+1}^m = \text{GenResp}(M, C). \tag{5}$$

However, existing works [15, 10], have found that directly generating the response of the role may sometimes be inconsistent with the character background, such as an ancient figure writing the Python code. Therefore, we propose an *iterative response revision* method. After generating a response $c_{T+1}^m$, we employ an LLM to first validate whether the generated response $c_{T+1}^m$ is consistent with the content of the role memory with the character background:

$$h = \text{Verify}(M, c_{T+1}^m, c_{T+1}^u) \in \{1, 2, 3, 4, 5\}, \tag{6}$$

where $h$ represents the consistency score, a score of 1 indicates the lowest consistency and 5 indicates the highest consistency. The instruction of Verify is:

> You are a helpful director, focused on the setting of the character $\{E_m\}$. Please give a score following the steps, your scoring range is from 1 to 5...
> $\{E_u\}$'s question is $\{c_{T+1}^u\}$, and the $\{E_m\}$'s response is: $\{c_{T+1}^m\}$.
> The setting for $\{E_m\}$ is $\{M\}$.
> Please assess how well the answer matches the setting of $\{E_m\}$. Explain the reason and then give a score.
> I will provide you with some sample outputs. Their main purpose is to help you understand the output format and judgment criteria: $\{Examples\}$

When the consistency score $h \leq \alpha$, where $\alpha$ is a threshold hyper-parameter, we revise the generated responses to align them with the background information of the role. To give more personalized information about the role for better revising the response, we retrieve $K$ relevant dialogues $D^r = \{(E_1^r, U_1^r), (E_2^r, U_2^r), \ldots, (E_K^r, U_K^r)\}$ from the historical dialogues $D$ by using the user's last utterance $c_{T+1}^u$ as the query. The newly retrieved dialogues $D^r$ are then used to update the role relationship graph $G$, and the weights of the newly added nodes and their associated edges are updated according to Equation 2. After updating the relation graph $G$, following the previous steps, the maximal clique is extracted again, and we generate a new role memory record $m$ and append it into the role memory $M$ (introduced in Equation 4). Finally, we re-generate the response $c_{T+1}^m$ based on the updated role memory $M$. The overall algorithm is shown in Algorithm 1.

---

**Algorithm 1** Pesudo-code of our proposed DRMR framework.

---

1: **Input:** Profiles $P_m$ and $P_u$ for roles $E_m$ (model) and $E_u$ (user), historical dialogues $D = \{(E_1, U_1), (E_2, U_2), \ldots, (E_L, U_L)\}$, and the current dialogue context $C = \{(c_1^u, c_1^m), \ldots, (c_T^u, c_T^m), (c_{T+1}^u)\}$
2: **Output:** Final response $c_{T+1}^m$ for role $E_m$
3: Initialize role memory $M$ with profile $P_m$
4: Retrieve top-$N$ relevant dialogues $D^c$ of $D$ using $C$
5: Construct relationship graph $G$ with $E_u$, $E_m$ and $\{E_1^c, E_2^c, \ldots, E_N^c\}$ in $D^c$
6: Add edges between roles in $G$ based on $D^c$
7: **while** $h > \alpha$ **do**
8:     Set $s_{i,j} = \text{EdgeScore}(I_{\text{ES}}\{U_i^c, U_j^c\})$
9:     Extract the maximal clique subgraph $G' = \text{argmax}_{G' \in G^\star} \sum_{i,j \in G'} s_{i,j}$
10:    Retrieve top-$K$ dialogues $D'$ from edges in $G'$ using $C$
11:    Set $m = \text{MemBuild}(D', P_m, P_u)$
12:    Set $M = M + m$
13:    Generate response $c_{T+1}^m = \text{GenResp}(M, C)$
14:    Set $h = \text{Verify}(M, c_{T+1}^m, c_{T+1}^u)$
15:    **if** $h \leq \alpha$ **then**
16:       Retrieve top-$K$ dialogues $D^r$ of $D$ using $c_{T+1}^u$
17:       Update $G$ based on $D^r$
18:    **end if**
19: **end while**
20: **return** Final response $c_{T+1}^m$

---

## 4   Experimental Setup

### 4.1   Evaluation Metrics

Role-play aims to customize LLMs to simulate various characters or personas with distinct and precise attributes, which provides a more nuanced interaction experience for users and makes LLMs more familiar [19, 24]. Consequently, immersion can be defined as the consistency of the model's responses with the role's personality and memory, as well as the familiarity felt by the user. This familiarity arises from the relationship between the model and the user's role. For example, family members feel familiar with each other due to their relationship, whereas passersby feel alienated because of the absence of a relationship and shared experiences. Therefore, to quantitatively measure the immersion performance of DRMR, we propose three evaluation metrics in our paper:

1. **Personality** (*Pers.*): Evaluate whether the responses align with the personality traits and linguistic habits. It also verifies whether their attitude towards current events is reasonable according to the dialogue history.

2. **Memorization** (*Mem.*): Assess the recollection of character-relevant experiences and knowledge, ensuring alignment with the background of the character. Relevant historical dialogues are retrieved to determine whether specific events mentioned in the dialogue history are reflected in the responses.
3. **Relation** (*Rela.*): Evaluate the degree to which the responses correspond to the relationship between user's and model's portrayed character. Considering the relationship between both roles (such as lover, family member, etc.), it judges whether the generated responses align with the relationship.

To evaluate the generated response according to the above criteria, we employ an LLM and prompt it with elaborate descriptions of the criteria to quantitatively evaluate the response. The LLM scores each response for the above three aspects separately using a scale of 1-5.

## 4.2   Dataset

In the experiments, we employ a Chinese benchmark role-play dataset CharacterEval [22], which contains 77 characters and 1,785 high-quality multi-turn dialogue contexts. In the CharacterEval dataset, the average turns per conversation is 9.28, and the average tokens per conversation is 369.69. We use the whole CharacterEval dataset as the test set to evaluate our model and baselines.

## 4.3   Baselines

We compare our method with several LLM-based role-play methods, including: RoleGPT [25],CharacterGLM [25], Qwen [1], ChatGLM [34].

  We employ three variants of DRMR: DRMR-C, DRMR-Q and DRMR-G with ChatGPT, Qwen and ChatGLM as the backbone respectively. And we also employ two ablation models:

1. DRMR w/o Revison: We remove the verify step (introduced in Equation 6) and directly use the output of the model as the response.
2. DRMR w/o RoleMem: We remove the graph-based role memory construction module and directly use the related dialogue as a prompt to the LLM.

## 4.4   Implementation Details

In our experiments, all DRMR-C variants and the RoleGPT use the gpt-3.5-turbo-0125 version, the DRMR-G variant and ChatGLM baseline use the glm-3-turbo API[2], and the DRMR-Q and Qwen are implemented using open-source Qwen-14B-chat as the backbone. In our model, we use the temperature 1.0 in most steps, and the temperature 0.1 during the verify step in Equation 6. For the consistency threshold used in the verify step, we set $\alpha = 4$. And we employ $N = 3$ and $K = 2$ retrieved dialogues when constructing role memory and revising the response respectively. We use the *text-embedding-ada-002* model of OpenAI as the embedding model used in Equation 1. We use the Bron-Kerbosch algorithm [2] in Equation 3 to find the maximal clique.

_____
[2] https://maas.aminer.cn/dev/api#glm-3-turbo

**Table 1.** Comparison of the response quality. ‡ indicates significant improvement over RoleGPT with $p \leq 0.01$ according to a Student's t test. The value in parentheses indicates the proportion of improvement compared to the LLM backbone.

| Method | Pers. (↑) | Mem. (↑) | Rela. (↑) | Human (↑) |
|---|---|---|---|---|
| CharacterGLM | 3.21 | 3.45 | 3.41 | 0.80 |
| ChatGLM | 3.68 | 4.01 | 3.67 | 0.81 |
| Qwen | 3.78 | 4.08 | 3.71 | 0.83 |
| RoleGPT | 3.39 | 3.47 | 3.49 | 0.75 |
| DRMR-G | 3.83(4.08%) | 4.12(2.74%) | **4.05**(10.35%) | **0.87** (7.41%) |
| DRMR-Q | **3.89**(2.91%) | **4.15**(1.72%) | 3.93(5.93%) | 0.85(2.41%) |
| DRMR-C | $3.65^{\ddagger}$(7.67%) | $3.76^{\ddagger}$(8.36%) | $3.93^{\ddagger}$(12.61%) | $0.84^{\ddagger}$(12.00%) |
| DRMR-C w/o Revision | 3.62 | 3.68 | 3.83 | 0.75 |
| DRMR-C w/o RoleMem | 3.61 | 3.65 | 3.72 | 0.76 |

## 5   Experimental Results

### 5.1   Overall Performance

Table 1 shows the performance of our proposed DRMR and baselines in terms of three metrics. We can find that DRMR shows superior performance in terms of all metrics compared to their corresponding backbone LLM. Furthermore, we find that our DRMR achieved a greater improvement in terms of the relation metric compared to its backbone LLM, indicating that our role memory construction contributes to generating higher-quality responses. This phenomenon indicates that our proposed DRMR is capable of generating a response that mimics the personality of the role in a higher quality.

### 5.2   Human Evaluation

To better examine whether the generated responses align with human intuitive judgment and further evaluate the enhancement of immersion, we employ human evaluation to further assess the baselines. We use three criteria for human evaluation, including (1) **Personality**: Assess whether the responses align with the personality of the character; (2) **Contextualize**: Determine if the responses correspond to the relevant events or background associated with the character in the ongoing conversation; (3) **Relationship Consistency**: Evaluate whether the responses of the character align with the relationship between the two roles in the current conversation.

We randomly select 300 generated results from each model and hire 3 educated annotators with master's degree to conduct double-blind annotation on randomly shuffled results. The score range of each aspect is 1-5. From Table 1, we observe that DRMR outperforms all baselines. And the ranking of human

**Table 2.** Examples of response generated by different methods. The text in <span style="color:red">red</span> indicates conflicting content with the plot, and the text in <span style="color:green">green</span> indicates the supporting facts for generating consistent response.

---

**User Role**: Anxin, **Model Role**: Gao
**Dialogue Context** $C$:
*Anxin*: When did you meet Lao Mo?
*Gao*: It's been many years. We met through Xiaohu. He wanted to rent a stall to do business but lacked experience, so Tang brought him to me and asked me to teach him.
*Anxin*: Have you been in touch since then?

---

**RoleGPT**: There was a period when we were <span style="color:red">in close contact</span>. We <span style="color:red">often discussed</span> business together. But later, as <span style="color:red">we both rose in the underworld</span>, we gradually drifted apart. However, Lao Mo has always been an important friend to me.
**DRMR**: We don't stay in touch very often, after all, we each have our own busy lives. But when there's a need, we still help each other out.
***RoleMemory of DRMR***: I see Anxin as a <span style="color:green">thorn</span> in my side, a threat to my position and power. He's a <span style="color:green">police officer</span>, and I'm a criminal underworld. We are <span style="color:green">enemies</span>, and I admit Anxin has surprising insight. He can find the elusive Lao Mo <span style="color:green">I've been hiding</span>.

---

evaluation is also consistent with the LLM-based automatic metrics, which also demonstrate the effectiveness of our proposed LLM-based evaluation method.

### 5.3   Ablation Study

To evaluate the effectiveness of each module in DRMR, we also conduct ablation studies with model DRMR-C, and the results are shown in Table 1. From this table, it can be observed that both ablation models perform worse than DRMR-C in terms of all metrics, indicating the effectiveness of the role memory and iterative response revision. We found that the DRMR w/o RoleMem method achieved lower scores compared to other ablation models, indicating the effectiveness of modeling the relationship between roles in our approach.

### 5.4   Case Study

Table 2 shows an example of responses generated by RoleGPT and DRMR-C. In this case, police officer Anxin interrogates criminal Gao. In RoleGPT's response, Gao admits to criminal interactions, which misaligns with Gao's background. In contrast, in the response generated by DRMR, Gao denies such interactions and cheats Anxin, maintaining character consistency and role-play immersion.

## 6   Discussion

**Table 3.** Comparison of the response quality on two subsets of the CharacterEval. The subset "Unseen" indicates that the content of the TV show has not been used as the pre-train data of the backbone LLM, while the characters in the "Seen" subset have been shown when pre-training LLM.
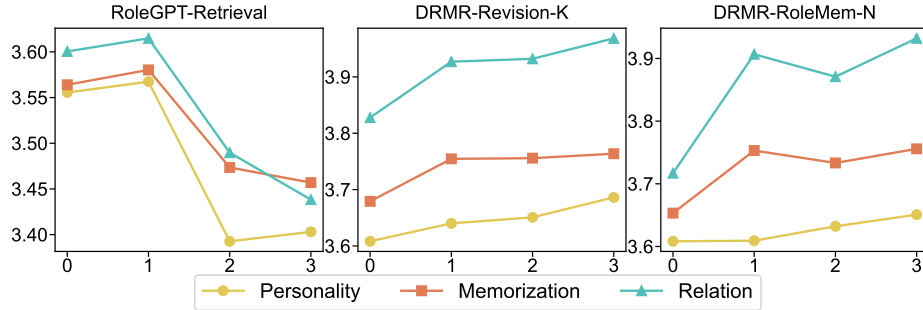
| Method | Unseen (↑) | Seen (↑) |
|--------|-----------|----------|
| RoleGPT | 0.58 | 0.77 |
| DRMR-C | **0.67**(15.51%) | **0.85** (10.38%) |

### 6.1   Analysis on Unseen Character

Due to the extensive use of web data for pre-training, LLM backbone is already familiar with most of the roles in the dataset CharacterEval. To validate the generalization ability of the model, we separate the data from CharacterEval for some newly released TV shows, which have not been trained on LLM. Thus, we divided the CharacterEval into two subsets, **seen** and **unseen**, not only based on the release time of the TV show but also by asking LLM if it knows the characters in the script. Table 3 shows the comparison between our proposed DRMR and RoleGPT on these two subsets. From the results, it can be seen that our method exhibits superior performance on both subsets, demonstrating better generalization ability of our DRMR. We can also find that both methods achieve higher scores on the seen dataset compared to the unseen dataset. As LLM has been trained on many data related to the role during the pre-training phase, it has a better understanding of the role compared to simply providing in-context information about the role. Due to the same reason, LLM may not fully understand the background of unseen characters, it cannot assess the quality of the response comprehensively. In this experiment, we employ human evaluation on 150 generated responses for each subset respectively, which uses the same criteria as in § 5.2.

### 6.2   Analysis of Using Different Numbers of Retrieved Dialogues

In § 3.2 and § 3.3, we employ the dense retrieval method [11] to find semantically related dialogues from historical dialogues of the role to enhance the role memory and revise the response. In this section, we explore the influence of using different numbers of retrieval dialogues on the final performance. The baseline RoleGPT also employs a similar retrieval approach to extract relevant information about the roles from historical dialogue data. Figure 3 illustrates the impact of using different numbers of retrieval dialogues on the performance of RoleGPT, our model in the revision stage, and our model in the role memory construction stage, respectively. From Figure 3, we observe that our approach effectively enhances response quality by using more retrieval dialogues in both stages. This demonstrates that our method leverages prompting LLM to construct role memory more effectively, thus utilizing data more efficiently. On the other hand, the baseline method RoleGPT struggles to extract useful information from excessive data, leading to a decline in the quality of generated responses.

**Fig. 3.** Performance of using different numbers of retrieved historical dialogues. The left figure shows the performance of using different historical dialogues when retrieval. The middle and right figures show the performance of using different historical dialogues in the revision stage and role memory construction stage, respectively.

**Table 4.** Comparison of token consumption. The "All" is the results on the entire dataset. The "Long" is the result of dialogues with more than 20 turns. The "Token" is the average token consumption for a multi-turn dialogue.
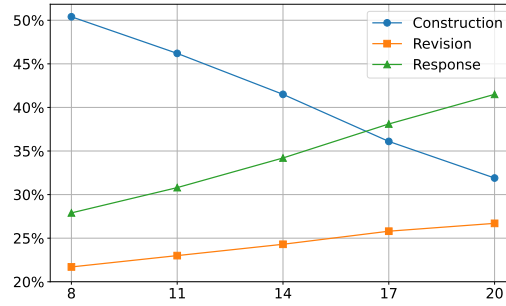
|      | Method | Pers. | Mem. | Rela. | Tokn. |
|------|--------|-------|------|-------|-------|
| ALL  | RoleGPT | 3.39 | 3.47 | 3.49 | **24K** |
|      | DRMR-C | **3.65** | **3.76** | **3.93** | 46K |
| Long | RoleGPT | 3.11 | 3.09 | 3.47 | **62K** |
|      | DRMR-C | **3.42** | **3.52** | **4.04** | 86K |

### 6.3 Analysis of Efficiency

Our method constructs and iteratively updates role memory using retrieved dialogues, which increases the token consumption of LLM. Table 4 presents the statistics of token consumption. The results show that, compared to RoleGPT, DRMR consumes more tokens but achieves better performance, especially in long dialogues. We analyze the token consumption of different modules in DRMR. As shown in Figure 4, the memory construction module consumes more tokens than the other two modules in short dialogues. However, as the number of dialogue turns increases, its proportion continuously decreases. The reason is that memory construction is frequent only at the beginning of a dialogue. Once the relationship is built, token consumption for this module will no longer increase.

## 7   Conclusion

In this paper, we present the **D**uplex **R**elationship **M**odeling based **R**ole-play framework (DRMR), an LLM-based role-playing framework aiming at enhancing the immersion of the user. We first introduce a novel maximal-cliques-based graph method to establish a duplex role relationship between characters played

**Fig. 4.** Token consumption proportion in different modules with different dialogue turns. Construction, Revision, and Response respectively represent the proportion of role memory construction, response generation, and iterative response revision.

by the user and the model. Next, we propose to leverage the reasoning ability of the LLM to summarize useful relationship information from the maximal cliques as role memory, and then generate the response by incorporating the role memory. To enhance the consistency between the generated responses and the background of the role, we propose the iterative response revision which first verifies the consistency of the response with the background knowledge of the role and then retrieves related dialogues to update the role memory and revise the response. Experimental results on the benchmark dataset demonstrate the superiority of the DRMR in elevating user immersion in role-playing interactions.

## Acknowledgment

## References

1. Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., et al.: Qwen technical report. arXiv preprint arXiv:2309.16609 (2023)
2. Bron, C., Kerbosch, J.: Algorithm 457: finding all cliques of an undirected graph. Communications of the ACM **16**(9), 575–577 (1973)
3. Chen, G., Dong, S., Shu, Y., Zhang, G., Sesay, J., Karlsson, B.F., Fu, J., Shi, Y.: Autoagents: A framework for automatic agent generation. arXiv preprint arXiv:2309.17288 (2023)
4. Chen, N., Wang, Y., Jiang, H., Cai, D., Li, Y., Chen, Z., Wang, L., Li, J.: Large language models meet harry potter: A dataset for aligning dialogue agents with characters. In: Findings of the Association for Computational Linguistics: EMNLP 2023. pp. 8506–8520 (2023)

5. Gao, S., Li, H., Shi, Z., Huang, C., Tu, Q., Shang, S., Tian, Z., Huang, M.: 360 rea: Towards a reusable experience accumulation with 360 assessment for multi-agent system. In: Findings of the Association for Computational Linguistics ACL 2024. pp. 13149–13162 (2024)

6. Gao, S., Ren, Z., Zhao, Y.E., Zhao, D., Yin, D., Yan, R.: Product-aware answer generation in e-commerce question-answering. Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining (2019), https://api.semanticscholar.org/CorpusID:59158861

7. Gao, S., Shi, Z., Zhu, M., Fang, B., Xin, X., Ren, P., Chen, Z., Ma, J., Ren, Z.: Confucius: Iterative tool learning from introspection feedback by easy-to-difficult curriculum. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 38, pp. 18030–18038 (2024)

8. Gao, S., Wen, Y., Zhu, M., Wei, J., Cheng, Y., Zhang, Q., Shang, S.: Simulating financial market via large language model based agents. arXiv preprint arXiv:2406.19966 (2024)

9. Hsu, S.L., Shah, R.S., Senthil, P., Ashktorab, Z., Dugan, C., Geyer, W., Yang, D.: Helping the helper: Supporting peer counselors via ai-empowered practice and feedback. arXiv preprint arXiv:2305.08982 (2023)

10. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P.: Survey of hallucination in natural language generation. ACM Computing Surveys **55**(12), 1–38 (2023)

11. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., et al.: Retrieval-augmented generation for knowledge-intensive nlp tasks. Advances in Neural Information Processing Systems **33**, 9459–9474 (2020)

12. Li, C., Leng, Z., Yan, C., Shen, J., Wang, H., MI, W., Fei, Y., Feng, X., Yan, S., Wang, H., et al.: Chatharuhi: Reviving anime character in reality via large language model. arXiv preprint arXiv:2308.09597 (2023)

13. Light, J., Cai, M., Shen, S., Hu, Z.: From text to tactic: Evaluating LLMs playing the game of avalon. In: NeurIPS 2023 Foundation Models for Decision Making Workshop (2023)

14. Mahadevan, K., Chien, J., Brown, N., Xu, Z., Parada, C., Xia, F., Zeng, A., Takayama, L., Sadigh, D.: Generative expressive robot behaviors using large language models. arXiv preprint arXiv:2401.14673 (2024)

15. Marcus, G.: The next decade in ai: four steps towards robust artificial intelligence. arXiv preprint arXiv:2002.06177 (2020)

16. Park, J.S., O'Brien, J., Cai, C.J., Morris, M.R., Liang, P., Bernstein, M.S.: Generative agents: Interactive simulacra of human behavior. In: Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology. pp. 1–22 (2023)

17. Qin, Y., Hu, S., Lin, Y., Chen, W., Ding, N., Cui, G., Zeng, Z., Huang, Y., Xiao, C., Han, C., et al.: Tool learning with foundation models. arXiv preprint arXiv:2304.08354 (2023)

18. Sapkaroski, D., Mundy, M., Dim, mock, M.R.: Immersive virtual reality simulated learning environment versus role-play for empathic clinical communication training. Journal of medical radiation sciences **69**(1), 56–65 (2022)

19. Shanahan, M., McDonell, K., Reynolds, L.: Role play with large language models. Nature **623**(7987), 493–498 (2023)

20. Shang, S., Yao, Z., Fu, H., Tao, C., Chen, X., Wang, F., Wang, Y., Ren, Z., Gao, S.: Unified multi-scenario summarization evaluation and explanation. IEEE Transactions on Knowledge and Data Engineering (2024)

21. Shao, Y., Li, L., Dai, J., Qiu, X.: Character-llm: A trainable agent for role-playing. In: Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. pp. 13153–13187 (2023)
22. Tu, Q., Fan, S., Tian, Z., Yan, R.: Charactereval: A chinese benchmark for role-playing conversational agent evaluation. arXiv preprint arXiv:2401.01275 (2024)
23. Wang, H., Feng, S., Chen, L., Liu, Y., Shang, S.: Simulating individual infection risk over big trajectory data. In: International Conference on Database Systems for Advanced Applications. pp. 136–151. Springer (2024)
24. Wang, Z., Zhang, G., Yang, K., Shi, N., Zhou, W., Hao, S., Xiong, G., Li, Y., Sim, M.Y., Chen, X., et al.: Interactive natural language processing. arXiv preprint arXiv:2305.13246 (2023)
25. Wang, Z.M., Peng, Z., Que, H., Liu, J., Zhou, W., Wu, Y., Guo, H., Gan, R., Ni, Z., Zhang, M., et al.: Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. arXiv preprint arXiv:2310.00746 (2023)
26. Wei, J., Bosma, M., Zhao, V., Guu, K., Yu, A.W., Lester, B., Du, N., Dai, A.M., Le, Q.V.: Finetuned language models are zero-shot learners. In: International Conference on Learning Representations (2022)
27. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems **35**, 24824–24837 (2022)
28. Wu, Q., Xu, W., Liu, W., Tan, T., Liu, J., Li, A., Luan, J., Wang, B., Shang, S.: Mobilevlm: A vision-language model for better intra-and inter-ui understanding. arXiv preprint arXiv:2409.14818 (2024)
29. Xu, X., Lei, Z., Wu, W., Niu, Z.Y., Wu, H., Wang, H.: Towards zero-shot persona dialogue generation with in-context learning. In: Findings of the Association for Computational Linguistics: ACL 2023. pp. 1387–1398 (2023)
30. Xu, Y., Wang, S., Li, P., Luo, F., Wang, X., Liu, W., Liu, Y.: Exploring large language models for communication games: An empirical study on werewolf. arXiv preprint arXiv:2309.04658 (2023)
31. Zhang, H., Liu, X., Zhang, J.: Extractive summarization via chatgpt for faithful summary generation. In: Conference on Empirical Methods in Natural Language Processing (2023)
32. Zhang, S., Dinan, E., Urbanek, J., Szlam, A., Kiela, D., Weston, J.: Personalizing dialogue agents: I have a dog, do you have pets too? In: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 2204–2213 (2018)
33. Zheng, C., Sabour, S., Wen, J., Zhang, Z., Huang, M.: Augesc: Dialogue augmentation with large language models for emotional support conversation. In: Findings of the Association for Computational Linguistics: ACL 2023. pp. 1552–1568 (2023)
34. Zhipuai: Chatglm. https://chatglm.cn (2023)
35. Zhong, P., Zhang, C., Wang, H., Liu, Y., Miao, C.: Towards persona-based empathetic conversational models. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 6556–6566 (2020)
36. Zhou, J., Chen, Z., Wan, D., Wen, B., Song, Y., Yu, J., Huang, Y., Peng, L., Yang, J., Xiao, X., et al.: Characterglm: Customizing chinese conversational ai characters with large language models. arXiv preprint arXiv:2311.16832 (2023)
37. Zhuang, Y., Yu, Y., Wang, K., Sun, H., Zhang, C.: ToolQA: A dataset for LLM question answering with external tools. In: Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track (2023)