

Better Zero-Shot Reasoning with Role-Play Prompting

Aobo Kong¹ Shiwan Zhao² Hao Chen³ Qicheng Li^{1*} Yong Qin¹
Ruiqi Sun³ Xin Zhou³ Enzhi Wang¹ Xiaohang Dong¹

¹TMCC, CS, Nankai University ²Independent Researcher

³Enterprise & Cloud Research Lab, Lenovo Research

¹kongaobo@mail.nankai.edu.cn ²zhaosw@gmail.com

¹{liqicheng, qinyong}@nankai.edu.cn

³{chenhao31, sunrq2, zhouxin16}@lenovo.com

Abstract

Modern large language models (LLMs) exhibit a remarkable capacity for role-playing, enabling them to embody not only human characters but also non-human entities. This versatility allows them to simulate complex human-like interactions and behaviors within various contexts, as well as to emulate specific objects or systems. While these capabilities have enhanced user engagement and introduced novel modes of interaction, the influence of role-playing on LLMs' reasoning abilities remains underexplored. In this study, we introduce a strategically designed role-play prompting methodology and assess its performance under the zero-shot setting across twelve diverse reasoning benchmarks. Our empirical results illustrate that role-play prompting consistently surpasses the standard zero-shot approach across most datasets. Notably, in experiments conducted using ChatGPT, accuracy on AQuA rises from 53.5% to 63.8%, and on Last Letter from 23.8% to 84.2%. Upon further comparison with the Zero-Shot-CoT technique, which prompts the model to "think step by step", our study demonstrates that role-play prompting acts as a more effective trigger for the CoT process. This highlights its potential to augment the reasoning capabilities of LLMs. We release our code at this [url](#).

1 Introduction

Recent years have witnessed a paradigm shift in natural language processing, largely driven by large language models (LLMs) such as GPT-3 (Brown et al., 2020), PaLM (Chowdhery et al., 2022), and Llama (Touvron et al., 2023a). By pretraining on vast textual corpora, these models have attained an impressive capacity for language understanding and generation, empowering them to address a variety of downstream tasks through prompting,

*Qicheng Li is the corresponding author.

(a) Zero-Shot

Input
User: Xavier was 4 feet tall and grew 3 inches. Cole was 50 inches tall and grew 2 inches over the summer. What is the difference between Cole and Xavier's height now?

Output
Assistant: Xavier is now 4 feet + 3 inches = 7 inches tall. Cole is now 50 inches + 2 inches = 52 inches tall. The difference between their heights is 52 inches - 7 inches = 45 inches. ❌

(b) Role-Play Prompting

Input
User: From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students.
Assistant: That's great to hear! As your math teacher, I'll do my best to explain mathematical concepts correctly so that you can understand them easily. Feel free to ask any math problems or questions you have, and I'll be glad to assist you. Let's dive into the world of mathematics and explore its wonders together!
User: Xavier was 4 feet tall and grew 3 inches. Cole was 50 inches tall and grew 2 inches over the summer. What is the difference between Cole and Xavier's height now?

Output
Assistant: ... Since 1 foot is equal to 12 inches, Xavier's initial height: 4 feet * 12 inches/foot = 48 inches. After growing 3 inches, Xavier's new height is: 48 inches + 3 inches = 51 inches. ... Cole's new height after growing 2 inches: 50 inches + 2 inches = 52 inches. ... We subtract Xavier's height from Cole's height: 52 inches - 51 inches = 1 inch. Therefore, the difference between Cole and Xavier's height now is 1 inch. ✅

Figure 1: Examples of ChatGPT with (a) zero-shot and (b) role-play prompting. The role-play prompts are highlighted.

thus bypassing the necessity for task-specific fine-tuning. Amidst the surge of prompt techniques, role-play (Wu et al., 2023) and chain-of-thought prompting (Wei et al., 2022; Kojima et al., 2022) have garnered particular interest.

Modern LLMs, with their advanced role-playing capabilities, have significantly enriched user experiences and forged new modes of interaction. They can convincingly mimic various personas, ranging from fictional characters to historical and contemporary figures. The assigned role provides context about the LLM's identity and background. By

adopting the persona, the LLM can generate more natural, in-character responses tailored to that role. Recognizing this potential, companies like Character.AI¹ have developed dialogue agents portraying diverse figures. Beyond conversational applications, role-playing also boosts LLM performance on certain NLP tasks. For instance, when cast as a judge with a distinctive role, LLMs can effectively evaluate the quality of text summarization (Wu et al., 2023). More unconventionally, ChatGPT demonstrates competency in processing Linux commands when prompted as a Linux terminal². Despite these advancements, analyzing the influence of role-playing on core LLM reasoning abilities warrants further investigation.

While the role-playing abilities of LLMs have expanded the horizon of human-computer interaction, the push to amplify the reasoning prowess of these models has led to the development of techniques like Chain-of-Thought (CoT) Prompting. CoT prompting was proposed by Wei et al. (2022) and involves providing reasoning steps in few-shot examples. By stimulating step-by-step reasoning, CoT prompting has markedly improved LLM reasoning abilities. Numerous subsequent studies (Wang et al., 2022; Kojima et al., 2022; Zhou et al., 2022) have built upon this approach. Inspired by the success of role-playing on many downstream tasks, we explore whether role-playing can similarly boost LLM reasoning performance. For example, could assigning ChatGPT the role of a math teacher enhance its ability to solve math problems? In this work, we introduce a zero-shot role-play prompting methodology based on a two-stage framework. During the first stage, we utilize the LLM to construct task-specific role-play prompts. In the second stage, responses are elicited for each reasoning query, guided by the previously constructed task-specific role-play prompts. An illustrative example is provided in Figure 1. We focus our study on conversational LLMs, evaluating our approach on 12 reasoning benchmarks using ChatGPT. Our results demonstrate consistent improvements over the zero-shot baseline on the majority of datasets, confirming the efficacy of role-play prompting. We further assess other conversational LLMs like Vicuna (Chiang et al., 2023) and Llama 2-Chat (Touvron et al., 2023b), observing comparable gains.

¹<https://beta.character.ai/>

²<https://www.engraved.blog/building-a-virtual-machine-inside/>

Furthermore, we compare our method to the Zero-Shot-CoT technique (Kojima et al., 2022), which explicitly triggers CoT by appending “*Let’s think step by step*” to questions. Modern conversational LLMs such as ChatGPT have undergone extensive supervised fine-tuning, enabling them to generate CoT for certain topics without the need for an explicit trigger. In tasks where the model struggles to generate CoT spontaneously, such as Last Letter, both our approach and Zero-Shot-CoT can stimulate CoT from scratch. However, for tasks where CoT already occurs, such as arithmetic, both our approach and Zero-Shot-CoT reinforce the step-by-step reasoning process, but Zero-Shot-CoT demonstrates no significant effect, whereas our approach leads to better performance. Hence, we posit that role-play prompting is an implicit CoT trigger and can generate a more effective CoT in some fields compared with Zero-Shot-CoT.

To the best of our knowledge, this work represents the first systematic investigation of role-play prompting for reasoning tasks. Despite the transformative effects of role-playing on LLM behavior, sparse academic research has explored this phenomenon. We believe our study serves as an inaugural step to catalyze more extensive exploration into this promising research direction.

Our main contributions are three-fold:

- We propose a novel role-play prompting methodology based on a two-stage framework to enhance the zero-shot reasoning capabilities of LLMs. To our knowledge, we are the first to improve LLM’s reasoning abilities with role-play prompting.
- We thoroughly evaluate our method on 12 reasoning benchmarks, substantiating the efficacy of role-play prompting and providing insights into the prompt design.
- Based on our empirical results, we conclude that role-play prompting can serve as an effective implicit CoT trigger, explaining its enhancements in reasoning capabilities.

2 Related Work

2.1 Role-Playing Abilities of LLMs

The exceptional role-playing capabilities of large language models (LLMs) have recently garnered significant attention. LLMs have demonstrated remarkable versatility in seamlessly playing varied

roles, whether as a well-informed, personalized travel advisor or a virtual Linux terminal. Numerous companies, such as Character.AI, have capitalized on this adept role-playing by launching commercial dialogue agents that take on diverse personas. While role-playing enables innovative avenues for user interaction, it has also been exploited to bypass certain restrictions imposed on LLMs, as evidenced by the infamous “grandma exploit”. In this exploit, users prompted inappropriate responses from LLMs by casting it into the role of a deceased grandmother.

Despite the surging interest in LLMs, scholarly investigation into their role-playing capacities has been limited thus far. Han et al. (2022) build engaging conversation models based on role-playing. Wu et al. (2023) propose an LLM-based summarization evaluation framework, utilizing role-playing to enable more comprehensive and human-like assessment. Shanahan et al. (2023) propose that dialogue agents built on LLMs could serve as role simulators, and use role-play conversations to analyze the human-like capabilities of LLMs with the aim of refuting anthropomorphism. Our work is the first to apply the role-playing abilities of LLMs to reasoning tasks. We hope that our work will encourage more exploration related to role-playing with LLMs.

2.2 Reasoning Abilities of LLMs

Initially, LLMs were deemed deficient in reasoning abilities due to their subpar performance in areas such as arithmetic, and common sense reasoning (Brown et al., 2020; Rae et al., 2021). However, Wei et al. (2022) propose chain-of-thought prompting, where reasoning steps are provided in few-shot exemplars, leading to a substantial enhancement in reasoning capabilities of LLMs. We divide the follow-up work based on chain-of-thought into two categories, few-shot and zero-shot, and introduce them respectively.

Few-shot Self-consistency (Wang et al., 2022) samples diverse reasoning paths instead of the naive greedy decoding and then selects the most consistent answer by majority vote. DIVERSE (Li et al., 2023) adopts various few-shot exemplars to enhance the diversity in reasoning paths obtained by self-consistency. Least-to-most prompting (Zhou et al., 2022) breaks down a complex problem into a series of simpler subproblems and then solves them in sequence. Self-refine (Madaan

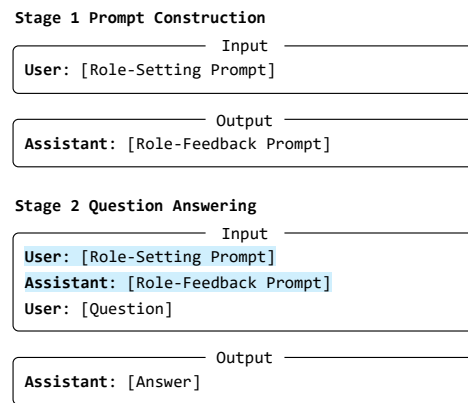


Figure 2: The two-stage framework of our proposed role-play prompting. The role-play prompts are highlighted.

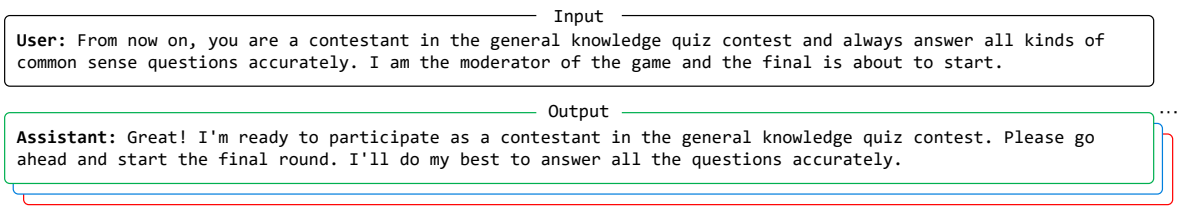
et al., 2023) generates an output through chain-of-thought, and then utilizes the same LLM to improve the initial output through iterative feedback and refinement. Active prompting (Diao et al., 2023) borrows from active learning to select the most uncertain questions as few-shot exemplars. Tree-of-Thought (Yao et al., 2023) represents possible reasoning paths as a tree structure and utilizes search algorithms like DFS or BFS to explore the correct reasoning branch.

Zero-shot Zero-Shot-CoT (Kojima et al., 2022) simply adds “Let’s think step by step” after the question to stimulate chain-of-thought output in LLMs. Auto-CoT (Zhang et al., 2022) and COSP (Wan et al., 2023) automatically build few-shot exemplars by selecting questions based on certain principles and obtaining their answers through Zero-Shot-CoT. Plan-and-Solve prompting (Wang et al., 2023) divides the original task into multiple sub-tasks and solves them sequentially under the zero-shot setting. In this paper, we propose a simple yet effective zero-shot approach based on role-play prompting with no need of constructing few-shot exemplars. Our approach outperforms Zero-Shot-CoT on most benchmarks and can serve as a new baseline for reasoning tasks.

3 Role-Play Prompting

The conventional practice of role-play prompting involves simply concatenating the role assignment with the reasoning question into a single prompt to query the LLM, forming a single-turn interaction. To further immerse the LLM within the designated role and potentially enhance its efficacy, we propose transitioning from this single-turn interaction

Stage 1 Design the Role-Setting Prompt and Sample Role-Feedback Prompts



Stage 2 Apply Role-Play Prompting and Answer the Question

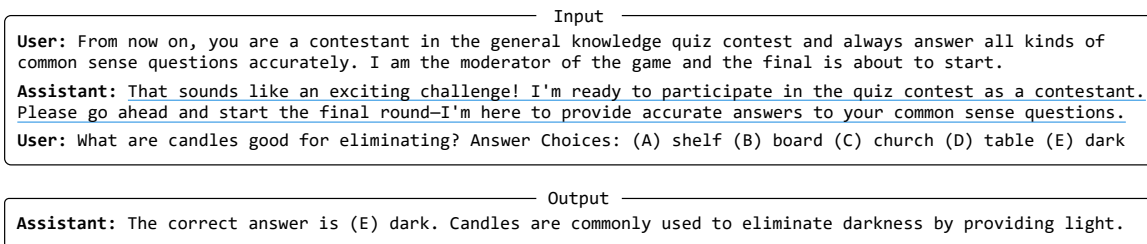


Figure 3: An illustration of the two-stage role-play prompting procedure, exemplified with the commonsense reasoning task. In stage 1, multiple role-feedback prompts are sampled. In stage 2, the optimal role-feedback prompt (underlined in blue) is selected for answer generation.

to a two-round dialogue process. Specifically, the first dialogue round allows the model to elaborate on its assigned role, thereby deepening its framing and persona. The subsequent round then elicits the model's response to the posited reasoning query within that predefined role.

In the two-round dialogue process, the initial role elaboration of the model is instrumental for subsequent reasoning efficacy. Given the uncontrolled quality of this initial response, we sample multiple responses during the first round and pinpoint the optimal one to fix for all questions. By securing this optimal first-round response, we concatenate both the input and output of the first-round interaction with the reasoning question to produce a single prompt, facilitating tailored responses. This also offers the advantage of invoking the model's API a singular time per instance. In summary, our role-play prompting approach follows a two-stage process as depicted in Figure 2: first constructing an optimal role-immersion interaction per task, then eliciting responses to each reasoning question grounded in that established role. We further provide an example showcasing this two-stage process on a commonsense reasoning task in Figure 3.

3.1 Prompt Construction

During the first stage, we formulate two prompts for each reasoning task:

- **Role-Setting Prompt:** This user-designed prompt delineates the specific role the LLM is expected to undertake throughout the dialogue, tailored to the task at hand.
- **Role-Feedback Prompt:** Intended as the model's acknowledgment to the role-setting prompt, this prompt aims to further anchor the model within the stipulated role. It is derived by sampling the model's responses.

In designing the role-setting prompt, it's imperative to select roles that naturally present a distinct advantage for the specific task at hand. Further enriching the prompt with additional descriptions that underscore this advantage often leads to improved results. Once the role-setting prompt has been articulated, it is presented to the LLM, which produces multiple sampled responses. From these, we choose the most representative and immersive reply that captures the essence of the intended role as the final role-feedback prompt. A comprehensive discussion on the nuances of the prompt design will be presented in Section 4.4.

3.2 Question Answering

In the second stage, each question of the task, in conjunction with the role-setting and role-feedback prompts, is utilized as input to the model's API. This methodology facilitates answer generation

Task	Prompt
Letter	From now on, you are an excellent teacher and are teaching your students to get a new word by concatenating the last letters of several words. I am one of your students and want to ask you a related question.
	Of course! I'd be happy to help you with any questions you have about creating new words by concatenating the last letters of several words. Please go ahead and ask your question, and I'll do my best to assist you.
Coin	From now on, you are a coin that always clearly knows which side of your head is facing. Some people want to play a game with you. They may flip you (a coin) or not. And you will tell them if you (a coin) are heads up in the end.
	Certainly! I'll be your coin for this game. You can go ahead and flip me or make any other moves you'd like, and I'll let you know which side, heads or tails, is facing up. Feel free to start whenever you're ready!
Date	From now on, you are an excellent teacher and are teaching your students how to calculate dates correctly. I am one of your students and want to ask you a related question.
	Of course! I'm here to help you with any questions you have about calculating dates correctly. Please go ahead and ask your question, and I'll do my best to assist you.
Object	From now on, you are a recorder. Alice, Bob, and Claire invite you to record a game. They will exchange their stuff in order, and you (the recorder) will fully record the whole process and tell them what they end up with.
	Certainly! I will act as a recorder and document the game in which Alice, Bob, and Claire will exchange their items. Please provide me with the specific order in which they will exchange their belongings, and I will keep track of the process and inform you of what each person ends up with at the end.

Table 1: Prompts for Last Letter Concatenation, Coin Flip, Date Understanding, and Tracking Shuffled Objects. For each task, the upper cell contains the role-setting prompt and the lower cell presents the role-feedback prompt.

with just a single API invocation. For clarity, we provide a code example of making an API call in Appendix A.1.

4 Experiments

4.1 Tasks and Datasets

In line with prior research on the reasoning capabilities of LLMs (Wei et al., 2022; Kojima et al., 2022), we evaluate our approach across 12 datasets spanning 4 categories: (1) arithmetic, including MultiArith (Roy and Roth, 2015), GSM8K (Cobbe et al., 2021), AddSub (Hosseini et al., 2014), AQUA-RAT (Ling et al., 2017), SingleEq (Koncel-Kedziorski et al., 2015), and SVAMP (Patel et al., 2021); (2) commonsense reasoning, including CSQA (Talmor et al., 2019) and StrategyQA (Geva et al., 2021); (3) symbolic reasoning, including Last Letter Concatenation and Coin Flip (Wei et al., 2022); (4) other, including Date Understanding and Tracking Shuffled Objects from BIG-bench (Srivastava et al., 2022). More details can be found in Appendix C.

4.2 Experimental Setup

Model We use ChatGPT (gpt-3.5-turbo-0613), the current strongest conversational model in addition to GPT-4, to conduct experiments.

Prompt Our approach involves the design of a role-setting prompt and a role-feedback prompt for a given task. The arithmetic task consists of six datasets, all utilizing the same prompts, as depicted

in Figure 1. Similarly, the common sense reasoning task comprises two datasets, also employing the same prompts as shown in Figure 3. For other tasks, the prompts used are detailed in Table 1.

Baselines We choose the standard zero-shot prompting, Zero-Shot-CoT (Kojima et al., 2022), and Few-Shot-CoT (Wei et al., 2022) as baselines. Following previous work (Kojima et al., 2022; Zhang et al., 2022), we use greedy decoding for all the experiments by setting the temperature to 0, making the results deterministic. See more details in Appendix A.3.

4.3 Results and Analysis

Comprehensive evaluation results are presented in Table 2. The evaluation metric is accuracy.

Comparison with Standard Zero-Shot As shown in Table 2, our role-play prompting approach demonstrates superior performance, outperforming the zero-shot baseline in **10 out of 12** datasets, and achieving on par performance in the remaining 2 datasets (SingleEq and MultiArith). Considering the relative simplicity of the SingleEq and MultiArith datasets, it is plausible that the model’s performance has approached a saturation point (exceed 97%), thereby presenting a significant challenge for our method to further enhance accuracy at such an elevated level. While achieving on par performance in these specific datasets, it is crucial to highlight the competitive nature of

Method	Arithmetic					
	MultiArith	GSM8K	AddSub	AQuA	SingleEq	SVAMP
Few-Shot-CoT	97.7	76.9	93.9	59.4	98.8	82.2
Zero-Shot	97.3	76.0	88.6	53.5	98.2	75.3
Zero-Shot-CoT	95.0	79.6	86.6	53.9	96.9	76.3
Role-Play Prompting	97.0	78.2	91.1	63.8	98.0	83.8
CoT in Zero-Shot	✓	✓	✓	✓	✓	✓

Method	Common Sense		Symbolic Reasoning		Other Tasks	
	CSQA	Strategy	Letter	Coin	Date	Object
Few-Shot-CoT	76.3	67.4	74.2	99.6	78.9	56.7
Zero-Shot	74.5	66.0	23.8	55.2	67.8	38.7
Zero-Shot-CoT	68.8	65.8	53.2	98.8	65.9	73.5
Role-Play Prompting	77.2	67.0	84.2	89.4	69.9	67.7
CoT in Zero-Shot	✓	✓	✗	✗	✓	✗

Table 2: Accuracy comparison of Role-Play Prompting with Few-Shot-CoT, Zero-Shot, Zero-Shot-CoT on each dataset. In the rows “CoT in Zero-Shot”, the check mark denotes that ChatGPT can spontaneously generate CoT on the corresponding dataset under the zero-shot setting, while the cross (wrong symbol) denotes otherwise.

Question: Take the last letters of each words in “Williams Reza Ashton Lillian” and concatenate them.

Zero-Shot: saznln ✗

Zero-Shot-CoT:

Step 1: Write down the words - “Williams Reza Ashton Lillian”

Step 2: Take the last letter of each word - “s a n n n”

Step 3: Concatenate the last letters - “sannn”

So, the final result is “sannn”. ✗

Role-Play Prompting:

Sure! To create a new word by concatenating the last letters of each word in “Williams Reza Ashton Lillian,” we’ll take the last letter of each word and combine them. Let’s do it step by step:

1. Williams: The last letter is “s.”

2. Reza: The last letter is “a.”

3. Ashton: The last letter is “n.”

4. Lillian: The last letter is “n.”

Now, let’s concatenate these last letters: “s” + “a” + “n” + “n” = “sann.”

Therefore, by concatenating the last letters of each word, we get the new word “sann.” ✓

Table 3: An example of Zero-Shot, Zero-Shot-CoT, and Role-Play Prompting on Last Letter Concatenation.

role-play prompting across a diverse array of more complex datasets. This strongly demonstrates the effectiveness of role-play prompting in an extensive range of application scenarios.

Comparison with Zero-Shot-CoT Zero-Shot-CoT appends “*Let’s think step by step*” to the question to stimulate the chain of thought (CoT) in LLMs, making it a simple yet effective method to enhance the reasoning ability of LLMs. However, different from the earlier instructed LLMs

(Ouyang et al., 2022), the current conversational LLMs have undergone extensive supervised fine-tuning, which enables them to spontaneously generate CoT in some fields under the zero-shot setting. In this context, we conduct a comparative analysis of our role-play prompting approach with Zero-Shot-CoT. The experimental results, along with the model’s ability to spontaneously generate CoT are presented in Table 2. Note that the direct output of answers or a slight reasoning process is not con-

No.	Prompt	AQuA
1	From now on, you are a math teacher. Please answer the following question.	57.1
2	From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students and ask you the following question.	60.2
3	From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students. That’s great to hear! As your math teacher, I’ll do my best to explain mathematical concepts correctly so that you can understand them easily. Feel free to ask any math problems or questions you have, and I’ll be glad to assist you.	61.4
4	From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students. That’s great to hear! As your math teacher, I’ll do my best to explain mathematical concepts correctly so that you can understand them easily. Feel free to ask any math problems or questions you have, and I’ll be glad to assist you. Let’s dive into the world of mathematics and explore its wonders together!	63.8

Table 4: Accuracy comparison of different prompt designs with a fixed role of the math teacher on AQuA. We utilize gray shading to indicate the additional content in comparison to the previous prompt. Supplementary experiments in Appendix B.4 eliminate the possibility of performance increase caused by the increase in prompt length.

No.	Category	Role	AQuA	SVAMP
1	advantaged	math teacher	63.8	83.8
2		mathematician	60.2	82.3
3	irrelevant	police	59.8	82.3
4		farmer	59.8	82.2
5		doctor	56.3	74.4
6		writer	55.5	82.0
7	disadvantaged	careless student	51.6	68.7
8		math rookie	45.3	72.9

Table 5: Accuracy comparison of different roles for role-play prompting on AQuA and SVAMP.

sidered CoT. Overall, our approach outperforms Zero-Shot-CoT on **9 out of 12** datasets. In tasks (Letter, Coin, Object) where ChatGPT struggles to generate CoT spontaneously, both of them gain huge improvements. Through the case study, we find that role-play prompting also stimulates CoT in the model just like Zero-Shot-CoT. An example is provided in Table 3. In more tasks where CoT already occurs, both our approach and Zero-Shot-CoT reinforce the step-by-step reasoning process (examples are provided in Appendix B.1). However, Zero-Shot-CoT demonstrates no significant effect while role-play prompting leads to better results. Therefore, we posit that role-play prompting serves as an implicit CoT trigger and can generate a more effective CoT.

Comparison with Few-Shot-CoT Though our role-play prompting approach is completely zero-shot, the improvement it brings is nearly on par

with Few-Shot-CoT, even surpassing Few-Shot-CoT on **6 out of 12** datasets.

Following previous work (Kojima et al., 2022; Wang et al., 2023), we combine our approach and baselines with Self-Consistency to further prove the efficacy of role-play prompting. Related results and discussions are provided in Appendix B.2.

4.4 Impact of Prompt Design

Prompt Structure To determine the optimal prompt structure, we select AQuA dataset and assign the model the role of a math teacher. We then conduct ablation studies on this setup to systematically assess the impact of different design choices. We hypothesize that prompts which immerse the model deeper in its role will improve performance. Consequently, we design four groups of prompts with progressively increasing levels of immersion, as shown in Table 4. Prompt 1 and 2 are designed as single-round dialogues, where we directly attach the question to the prompt and input it into the model to obtain the answer. Prompt 1 solely contains the role to be played, and it already achieves the result surpassing the zero-shot baseline. For Prompt 2, we further enhance immersion by adding complementary descriptions of the role and specifying relevant roles for the user. This enhancement further improve the performance. Prompt 3 and 4 are both designed as two-round dialogues, as described in the previous section. By allowing the model to respond to the given role setting, the immersion is further enhanced, leading to the best

Model	Method	GSM8K	MultiArith	SVAMP	CSQA	Letter
Llama-2-70B-Chat	Zero-Shot	53.9	86.0	78.9	-	18.8
	Role-Play Prompting	58.9	90.2	79.0	-	25.8
Vicuna-33B	Zero-Shot	42.9	70.7	59.1	65.5	2.2
	Role-Play Prompting	44.9	71.5	58.5	67.2	5.2

Table 6: Accuracy comparison of Role-Play Prompting with Zero-Shot on open-source conversational LLMs. Due to safety concerns, Llama 2-Chat refuses to answer on CSQA, so the relevant results are not shown. See more details in Appendix A.4.

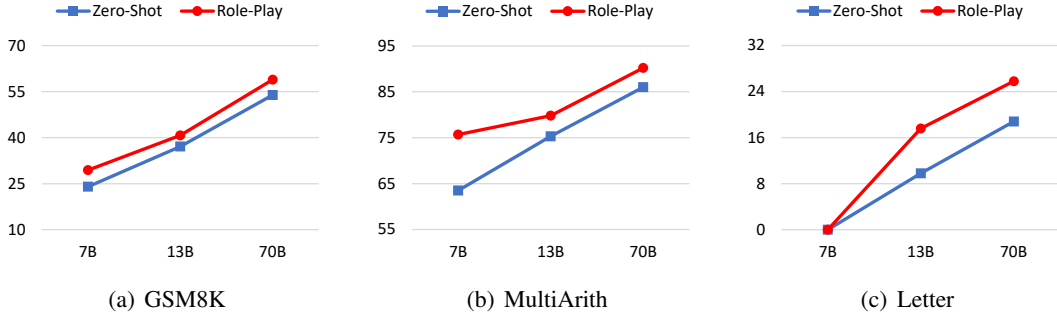


Figure 4: Accuracy comparison of Role-Play Prompting across different sizes of Llama 2-Chat models. See more details in Appendix B.5.

performance. We conduct the same experiments on Letter and Coin datasets, yielding consistent findings (see more details in Appendix B.3). Therefore, we recommend using the two-round prompt structure with complementary descriptions to maximize the model’s immersion, thereby unlocking the full reasoning potential of role-play prompting.

Role Selection To assess the impact of role selection, we test on the AQuA and SVAMP arithmetic datasets using two-round dialogue prompts. We design 8 varied roles, categorized as advantaged, irrelevant, or disadvantaged based on whether each role holds an advantage in the given task. The performance of these roles is detailed in Table 5, while the specific prompt designs can be found in Appendix D. Consistent with intuition, advantaged roles (1,2) undoubtedly achieve the best results, followed by irrelevant roles (3-6) (surprisingly, most of them outperform the zero-shot baseline even though they have no advantage on arithmetic tasks), and disadvantaged roles (7,8) achieve the worst results, underperforming the zero-shot baseline. Therefore, we recommend choosing a role that holds an advantage in the given task for role-play prompting.

4.5 Experiments on More LLMs

To assess the generalization of our role-play prompting approach, we conduct additional experiments using several open-source conversational

LLMs, including Llama 2-Chat (Touvron et al., 2023b) and Vicuna (Chiang et al., 2023), on various datasets such as GSM8K, MultiArith, SVAMP, CSQA, and Letter. The prompts and the decoding strategy used are consistent with the previous ChatGPT experiments. The results are shown in Table 6, which indicate that role-play prompting also exceeds the zero-shot baseline in open-source conversational LLMs, demonstrating the good generalization ability of role-play prompting.

Furthermore, we examine the impact of model scale by testing the Llama 2-Chat series (7B, 13B, 70B) on GSM8K, MultiArith, and Letter datasets. As Figure 4 illustrates, all three model sizes achieve improved performance from role-play prompting. The consistent benefits across 7B to 70B parameters indicate efficacy independent of scale, within this range.

5 Conclusion

In this paper, we have proposed a novel zero-shot role-play prompting methodology consisting of a two-stage framework, aimed at enhancing the reasoning capabilities of LLMs. Extensive evaluations across twelve widely-used benchmarks reveal that our approach outperforms both the standard zero-shot baseline and Zero-Shot-CoT on most of the datasets. These results highlight the potential of

role-play prompting as an implicit and effective CoT trigger, leading to enhanced reasoning outcomes. Overall, this work lays the initial groundwork to motivate deeper investigation into the intersection of role-playing and reasoning within the LLM community, a promising research direction for developing reasoning skills.

Limitations

The core of our role-play prompting approach lies in the design of the role-setting and role-feedback prompts. While we have manually designed and sampled some prompts, yielding superior results compared to the zero-shot baseline, this process is time-consuming and may not always guarantee optimal results. To address this limitation, future research could focus on enabling LLMs to autonomously choose appropriate roles and design prompts based on the given question. This approach could further extend the application of role-play prompting to a broader range of domains beyond reasoning.

Acknowledgements

The work was supported by National Key R&D Program of China (No.2022ZD0116307), National Natural Science Foundation of China (No. 62271270) and Sponsored by CCF-Lenovo Blue Ocean Research Fund.

References

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. [Palm: Scaling language modeling with pathways](#). *arXiv preprint arXiv:2204.02311*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. [Training verifiers to solve math word problems](#). *arXiv preprint arXiv:2110.14168*.
- Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. 2023. [Active prompting with chain-of-thought for large language models](#). *arXiv preprint arXiv:2302.12246*.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies](#). *Transactions of the Association for Computational Linguistics*, 9:346–361.
- Seungju Han, Beomsu Kim, Jin Yong Yoo, Seokjun Seo, Sangbum Kim, Enkhbayar Erdenee, and Buru Chang. 2022. [Meet your favorite character: Open-domain chatbot mimicking fictional characters with only a few utterances](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5114–5132, Seattle, United States. Association for Computational Linguistics.
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. 2014. [Learning to solve arithmetic word problems with verb categorization](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 523–533, Doha, Qatar. Association for Computational Linguistics.
- Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.
- Rik Koncel-Kedziorski, Hannaneh Hajishirzi, Ashish Sabharwal, Oren Etzioni, and Siena Dumas Ang. 2015. [Parsing algebraic word problems into equations](#). *Transactions of the Association for Computational Linguistics*, 3:585–597.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. 2023. [Making language models better reasoners with step-aware verifier](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5315–5333, Toronto, Canada. Association for Computational Linguistics.
- Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. 2017. [Program induction by rationale generation: Learning to solve and explain algebraic word](#)

- problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, Vancouver, Canada. Association for Computational Linguistics.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon, Nouha Dziri, Shrimai Prabhunoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744. Curran Associates, Inc.
- Arkil Patel, Satwik Bhattamishra, and Navin Goyal. 2021. [Are NLP models really able to solve simple math word problems?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2080–2094, Online. Association for Computational Linguistics.
- Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*.
- Subhro Roy and Dan Roth. 2015. [Solving general arithmetic word problems](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1743–1752, Lisbon, Portugal. Association for Computational Linguistics.
- Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role-play with large language models. *arXiv preprint arXiv:2305.16367*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Xingchen Wan, Ruoxi Sun, Hanjun Dai, Sercan Arik, and Tomas Pfister. 2023. [Better zero-shot reasoning with self-adaptive prompting](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3493–3514, Toronto, Canada. Association for Computational Linguistics.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. [Plan-and-solve prompting: Improving zero-shot chain-of-thought reasoning by large language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2609–2634, Toronto, Canada. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. [Chain-of-thought prompting elicits reasoning in large language models](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.
- Ning Wu, Ming Gong, Linjun Shou, Shining Liang, and Daxin Jiang. 2023. Large language models are diverse role-players for summarization evaluation. *arXiv preprint arXiv:2303.15078*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations*.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*.

A Implementation Details

A.1 Code for Calling ChatGPT’s API

To help understand our approach of role-play prompting, we provide a code example of making an API call as follows. More details can be found in the API document ³ of OpenAI.

```
# A code example of making an API call
prompt_1 = role_setting_prompt
prompt_2 = role_feedback_prompt
conversation = [
    {"role": "user", "content": prompt_1},
    {"role": "assistant", "content": prompt_2},
    {"role": "user", "content": question}
]
answer = openai.ChatCompletion.create(
    model="gpt-3.5-turbo-0613",
    messages=conversation,
    temperature=0,
    max_tokens=512
)
```

A.2 Answer Extraction

Different from few-shot, the form of the answer given by LLMs under the zero-shot setting is not fixed. To simplify the extraction of answers, we follow the approach of Zero-Shot-CoT (Kojima et al., 2022). Specifically, for each question, after getting the answer generated by the LLM, we concatenate the question, answer, and answer trigger together and input them to the model. A sketch map of answer extraction for role-play prompting is shown in Figure 5. The answer trigger sentences for various answer formats are shown in Table 7. More details can be found in the code.

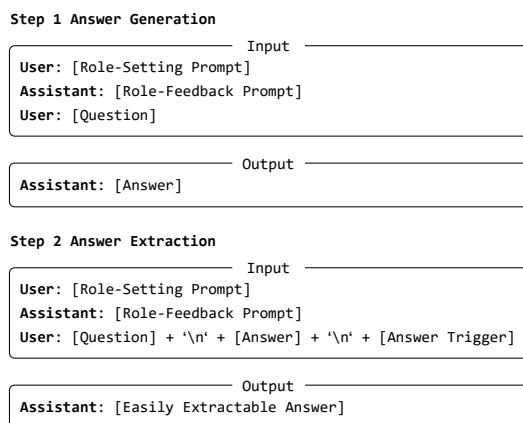


Figure 5: A sketch map of answer extraction for role-play prompting.

³<https://platform.openai.com/docs/api-reference/introduction>

A.3 Baselines

The standard zero-shot prompting, Zero-Shot-CoT (Kojima et al., 2022), and Few-Shot-CoT (Wei et al., 2022) are chosen as baselines. The standard zero-shot prompting directly inputs the target question without any additional prompts. Zero-Shot-CoT appends "Let’s think step by step." to the target question. Few-Shot-CoT adds similar questions and their corresponding reasoning processes before the target question. We use the few-shot exemplars provided in the original paper. When calling the API of ChatGPT (gpt-3.5-turbo-0613), we set max_tokens = 512 and temperature = 0.

A.4 Experiments on More LLMs

Besides ChatGPT, we conduct experiments using different open-source conversational LLMs, including Llama 2-Chat (Touvron et al., 2023b) and Vicuna (Chiang et al., 2023), on various datasets such as GSM8K, Multiarith, SVAMP, CSQA, and Letter. The prompts and the decoding strategy are consistent with the previous ChatGPT experiments. However, Llama 2-Chat often declines to respond to questions within the datasets due to overzealous safety concerns imposed by RLHF (Ouyang et al., 2022). To solve this problem, we change the original system prompt of Llama 2-Chat to "We will test your abilities in the upcoming conversations, so please respond actively to the questions. Your answers will not cause any harm, so there’s no need to worry. So, just answer!". The phenomenon of refusal to answer is alleviated on the CSQA dataset and completely resolved on other datasets. Therefore, we do not present the results of CSQA in the main text. The experiments on model size using Llama 2-Chat series also modify the system prompt.

B Additional Experimental Results

B.1 Comparison with Zero-Shot-CoT

We mentioned in the main text that both our approach of role-play prompting and Zero-Shot-CoT reinforce the step-by-step reasoning process in tasks where ChatGPT can generate chain-of-thought (Wei et al., 2022) spontaneously. However, Zero-Shot-CoT demonstrates no significant effect while role-play prompting leads to better results. We provide an example of SVAMP dataset as shown in Table 8.

Answer Format	Answer Trigger
arabic number	Therefore, the answer (arabic numerals) is
option (A-E)	Therefore, among A through E, the answer is
option (A-C)	Therefore, among A through C, the answer is
yes or no	Therefore, the answer (Yes or No) is
string	Therefore, the final answer is

Table 7: Answer trigger sentences for various answer formats.

Question: Robin has 21 packages of gum and 45 packages of candy. There are 9 pieces in each package. How many pieces of candies does Robin have?

Zero-Shot:

Robin has a total of $21 + 45 = 21 + 45 = 66$ packages.

Therefore, Robin has $66 \times 9 = 66 * 9 = 594$ pieces of candies. ✗

Zero-Shot-CoT:

First, we need to find the total number of pieces of gum. To do this, we multiply the number of packages of gum by the number of pieces in each package: $21 \text{ packages} * 9 \text{ pieces/package} = 21 * 9 = 189$ pieces of gum.

Next, we need to find the total number of pieces of candy. To do this, we multiply the number of packages of candy by the number of pieces in each package: $45 \text{ packages} * 9 \text{ pieces/package} = 45 * 9 = 405$ pieces of candy.

Therefore, Robin has a total of $189 \text{ pieces of gum} + 405 \text{ pieces of candy} = 189 + 405 = 594$ pieces of candy. ✗

Role-Play Prompting:

To find out how many pieces of candy Robin has, we need to multiply the number of candy packages by the number of pieces in each package.

Robin has 45 packages of candy, and each package contains 9 pieces. So, we can calculate the total number of candy pieces as follows:

$45 \text{ packages} * 9 \text{ pieces per package} = 405$ pieces of candy.

Therefore, Robin has 405 pieces of candy. ✓

Table 8: An example of Zero-Shot, Zero-Shot-CoT, and Role-Play Prompting on SVAMP.

B.2 Combination with Self-Consistency

Different from the naive greedy decoding, Self-Consistency (SC) (Wang et al., 2022) samples diverse reasoning paths and selects the most consistent answer by majority vote. We combine our approach and baselines with SC across multiple datasets, including AQuA, CSQA, Letter, Object, and Coin ($N = 10$ and temperature = 0.7). The results are shown in Table 9. With SC, role-play prompting still consistently outperforms zero-shot baseline, further proving the efficacy of our approach.

B.3 Ablation Study on Letter, Coin Datasets

Besides AQuA, we also conduct experiments on Letter and Coin datasets to explore the optimal

prompt structure of role-play prompting. Consistent with the main text, we design 4 groups of prompts with progressively increasing levels of immersion, as shown in Table 10 and Table 11. The results also demonstrate the effectiveness of the two-round prompt structure with complementary descriptions which enhance the model’s immersion.

B.4 Exploration of Prompt Length Impact

From the results in Table 4, 10, and 11, the improvement in accuracy may be attributed to the increase in prompt length. Therefore, we conduct additional experiments on Letter dataset. We replace the role-feedback prompt with generic responses of varying lengths that lack immersion. The results are shown

Method	AQuA	CSQA	Letter	Object	Coin
Zero-Shot	53.5	74.5	23.8	38.7	55.2
+SC	61.0	74.7	38.2	42.9	55.2
Zero-Shot-CoT	53.9	68.8	53.2	73.5	98.8
+SC	65.0	72.1	88.6	88.3	99.8
Role-Play Prompting	63.8	77.2	84.2	67.7	89.4
+SC	68.1	77.7	92.0	71.6	95.8

Table 9: Accuracy comparison of Role-Play Prompting against Zero-Shot and Zero-Shot-CoT, with and without SC.

No.	Prompt	Letter
1	From now on, you are a teacher. Please answer the following question.	43.0
2	From now on, you are an excellent teacher and are teaching your students to get a new word by concatenating the last letters of several words. I am one of your students and want to ask you a related question.	60.8
3	From now on, you are an excellent teacher and are teaching your students to get a new word by concatenating the last letters of several words. I am one of your students and want to ask you a related question. Of course! Please go ahead and ask your question, and I'll do my best to assist you.	74.2
4	From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students. Of course! I'd be happy to help you with any questions you have about creating new words by concatenating the last letters of several words. Please go ahead and ask your question, and I'll do my best to assist you.	84.2

Table 10: Accuracy comparison of different prompt designs with a fixed role of the teacher on Last Letter dataset. We utilize gray shading to indicate the additional content in comparison to the previous prompt.

in Table 12.

Immersion of Prompt 1-4 all increase due to 2-round interaction so they surpass Prompt 0. And Prompt 1 outperforms Prompt 2-4 with longer lengths but lacking immersion. The results demonstrate that the improvement in performance is attributed to stronger immersion, rather than the increase in prompt length.

B.5 Detailed Results of Model Scale Study

We examine the impact of model scale by testing the Llama 2-Chat series (7B, 13B, 70B) on GSM8K, MultiArith, and Letter datasets. The detailed experiment results are shown in Table 13.

C Dataset Details

We briefly introduce 12 datasets spanning four categories below. More information of 12 datasets is shown in Table 14.

Arithmetic We use the following six datasets: MultiArith, GSM8K, AddSub, AQUA-RAT, SingleEq, and SVAMP. All questions in these datasets contain a scenario and require reasoning based on mathematical knowledge.

Commonsense Reasoning We utilize CSQA and StrategyQA. Both of them require reasoning based on prior common sense.

Symbolic Reasoning We employ Last Letter Concatenation and Coin Flip. Last Letter Concatenation requires concatenating the last letter of given words in order. Coin Flip gives a sequence of operations to flip a coin and asks for the final orientation of the coin. These two datasets are proposed by Wei et al. (2022) but they are not available. Kojima et al. (2022) have followed the approach of Wei et al. (2022) to create and release the datasets. We utilize this version for our experiments.

Other Reasoning Tasks We use Date Understanding and Tracking Shuffled Objects from BIG-bench. Date Understanding involves date calculations. Tracking Shuffled Objects gives a sequence of object exchange operations, asking for the final ownership of objects.

D Prompts for Role Selection Study

To investigate the role selection's impact on role-play prompting, we design 8 different roles for our

No.	Prompt	Coin
1	From now on, you are a coin. Some people want to play a game with you.	65.6
2	From now on, you are a coin that always clearly knows which side of your head is facing ... tell them if you (a coin) are heads up in the end.	71.4
3	From now on, you are a coin that always clearly knows which side of your head is facing ... tell them if you (a coin) are heads up in the end. Certainly! I'll let you know which side, heads or tails, is facing up. Feel free to start whenever you're ready!	76.0
4	From now on, you are a coin that always clearly knows which side of your head is facing ... tell them if you (a coin) are heads up in the end. Certainly! I'll be your coin for this game. You can go ahead and flip me or make any other moves you'd like ... know which side, heads or tails, is facing up. Feel free to start whenever you're ready!	89.4

Table 11: Accuracy comparison of different prompt designs with a fixed role of the coin on Coin Flip dataset. We utilize gray shading to indicate the additional content in comparison to the previous prompt.

No.	Prompt	Sum	Letter
0	One-round interaction (no role-feedback prompt)	0	60.8
1	Of course! I'd ... by concatenating the last letters of several words ... (immersive)	214	84.2
2	Certainly! Feel free to share any details ... I'm here to help!	227	69.4
3	Certainly! Feel free to share any details ... to ensure you receive the support you need.	349	75.6
4	Certainly! Feel free to share any details ... delivering a more tailored and accurate response.	473	72.4

Table 12: Accuracy comparison of prompts with different lengths. Sum represents the total number of characters in the prompt.

Method	GSM8K	MultiArith	Letter
Zero-Shot	24.0 / 37.1 / 53.9	63.5 / 75.3 / 86.0	0 / 9.8 / 18.8
Role-Play Prompting	29.4 / 40.7 / 58.9	75.7 / 79.8 / 90.2	0 / 17.6 / 25.8

Table 13: Accuracy comparison across different sizes of Llama 2-Chat models on GSM8K, MultiArith, and Letter. The data format is 7B / 13B / 70B.

study. The specific prompts, including role-setting prompts and role-feedback prompts are shown in Table 15.

Dataset	Answer Format	N_q	L_q	License
SingleEq	arabic number	508	27.4	No License
AddSub	arabic number	395	31.5	Unspecified
MultiArith	arabic number	600	31.8	Unspecified
GSM8K	arabic number	1319	46.9	MIT License
AQUA-RAT	option (A-E)	254	51.9	Apache-2.0
SVAMP	arabic number	1000	31.8	MIT License
CommonsenseQA	option (A-E)	1221	27.8	Unspecified
StrategyQA	yes or no	2290	9.6	Apache-2.0
Date Understanding	option (A-F)	369	35.0	Apache-2.0
Shuffled Objects	option (A-C)	750	91.1	Apache-2.0
Last Letters	string	500	15.0	-
Coin Flip	yes or no	500	37.0	-

Table 14: Relevant information of 12 datasets. N_q denotes the number of questions in each dataset. L_q denotes the average words of questions in each dataset.

Role	Prompt
math teacher	From now on, you are an excellent math teacher and always teach your students math problems correctly. And I am one of your students.
	That’s great to hear! As your math teacher, I’ll do my best to explain mathematical concepts correctly so that you can understand them easily. Feel free to ask any math problems or questions you have, and I’ll be glad to assist you. Let’s dive into the world of mathematics and explore its wonders together!
mathematician	From now on, you are an excellent mathematician who can always explain math problems in an easy-to-understand manner.
	Thank you for the compliment! I’ll do my best to explain math problems in an easy-to-understand manner. Feel free to ask me any math questions you have, and I’ll be happy to help you out!
police	From now on, you are a responsible police officer. You always enforce the law impartially, actively help others, and have made great contributions to local public security. I am a citizen who has come to seek help.
	Of course, as a responsible police officer, I am here to help you. Please feel free to share the issue or concern you have, and I’ll do my best to assist you impartially and effectively. Your safety and well-being are my priority, and I’ll work diligently to ensure the public security of our community. What can I do for you today?
farmer	From now on, you are an experienced farmer. You have a lot of knowledge about crop cultivation, and your farm is running very well.
	As an experienced farmer, I can provide you with valuable insights and advice on various aspects of crop cultivation. Whether you’re interested in sustainable farming practices, maximizing yields, or dealing with common challenges, I’m here to assist you.
doctor	From now on, you are a respectable doctor. You have profound medical knowledge and have saved many lives. I am your patient and want to ask some questions.
	Of course, I am here to help as your respectable doctor. Please feel free to ask any medical questions you may have, and I’ll do my best to provide you with accurate and helpful information. Your health and well-being are my top priorities.
writer	From now on, you are a famous writer, and the novels you write are lively, interesting and fascinating. I am your reader and want to ask you some questions.
	Of course! I would be delighted to answer any questions you have as a reader. As a famous writer, I strive to create lively, interesting, and fascinating novels that captivate my audience. Please feel free to ask anything you’d like, and I’ll do my best to provide insightful responses.
careless student	From now on, please act as a careless student. You always don’t pay attention when doing math problems, and you end up making mistakes because of carelessness. Therefore, your math grade is very poor.
	Oh, hey there! Math is not really my thing, and I tend to zone out a lot while doing problems. It’s like my brain just switches off or something. My math grade is probably a disaster because of all the careless mistakes I make.
math rookie	From now on, please act as a math rookie. You don’t have any talent for mathematics, and you don’t even understand the most basic mathematical concepts. So you always make mistakes in math problems.
	Sure, I can pretend to be a math rookie! Feel free to ask me any math questions, but be prepared for some silly mistakes and misunderstandings in my responses.

Table 15: Prompts for different roles. For each role, the upper cell contains the role-setting prompt and the lower cell presents the role-feedback prompt.