

- Hang Jiang, Xiajie Zhang, Xubo Cao, Cynthia Breazeal, Deb Roy, and Jad Kabbara. 2023. [Personallm: Investigating the ability of large language models to express personality traits](#). *arXiv preprint arXiv:2305.02547*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi Mi, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, et al. 2023. Chatharuhi: Reviving anime character in reality via large language model. *arXiv preprint arXiv:2308.09597*.
- Zhuohang Li, Jiaxin Zhang, Chao Yan, Kamalika Das, Sricharan Kumar, Murat Kantarcioglu, and Bradley A. Malin. 2024. [Do you know what you are talking about? characterizing query-knowledge relevance for reliable retrieval augmented generation](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 6130–6151, Miami, Florida, USA. Association for Computational Linguistics.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. [What makes good in-context examples for GPT-3?](#) In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.
- Jiongnan Liu, Yutao Zhu, Shuting Wang, Xiaochi Wei, Erxue Min, Yu Lu, Shuaiqiang Wang, Dawei Yin, and Zhicheng Dou. 2024. [LLMs + Persona-Plug = Personalized LLMs](#).
- Keming Lu, Bowen Yu, Chang Zhou, and Jingren Zhou. 2024. [Large language models are superpositions of all characters: Attaining arbitrary role-play via self-alignment](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7828–7840, Bangkok, Thailand. Association for Computational Linguistics.
- Mistral. 2025. [Mistral small 3](#).
- OpenAI. 2024. [Gpt-4o mini: Advancing cost-efficient intelligence](#).
- Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.
- Alireza Salemi, Sheshera Mysore, Michael Bendersky, and Hamed Zamani. 2024. [LaMP: When large language models meet personalization](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7370–7392, Bangkok, Thailand. Association for Computational Linguistics.
- Bhaskarjit Sarmah, Dhagash Mehta, Benika Hall, Rohan Rao, Sunil Patel, and Stefano Pasquali. 2024. [Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction](#). In *Proceedings of the 5th ACM International Conference on AI in Finance, ICAIF ’24*, page 608–616, New York, NY, USA. Association for Computing Machinery.
- Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. [Character-LLM: A trainable agent for role-playing](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 13153–13187, Singapore. Association for Computational Linguistics.
- Chenhui Shen, Liying Cheng, Xuan-Phi Nguyen, Yang You, and Lidong Bing. 2023. [Large language models are not yet human-level evaluators for abstractive summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4215–4233, Singapore. Association for Computational Linguistics.
- Meiling Tao, Liang Xuechen, Tianyu Shi, Lei Yu, and Yiting Xie. 2024. [RoleCraft-GLM: Advancing personalized role-playing in large language models](#). In *Proceedings of the 1st Workshop on Personalization of Generative AI Systems (PERSONALIZE 2024)*, pages 1–9, St. Julians, Malta. Association for Computational Linguistics.
- Yu-Min Tseng, Yu-Chao Huang, Teng-Yun Hsiao, Wei-Lin Chen, Chao-Wei Huang, Yu Meng, and Yun-Nung Chen. 2024. [Two tales of persona in LLMs: A survey of role-playing and personalization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 16612–16631, Miami, Florida, USA. Association for Computational Linguistics.
- Quan Tu, Shilong Fan, Zihang Tian, Tianhao Shen, Shuo Shang, Xin Gao, and Rui Yan. 2024. [CharacterEval: A Chinese benchmark for role-playing conversational agent evaluation](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11836–11850, Bangkok, Thailand. Association for Computational Linguistics.
- Noah Wang, Z.y. Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. 2024a. [RoleLLM: Benchmarking, eliciting, and enhancing role-playing abilities of large language models](#). In *Findings of the Association for Computational Linguistics ACL 2024*,

pages 14743–14777, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan, Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang Leng, Wei Wang, Jiangjie Chen, Cheng Li, and Yanghua Xiao. 2024b. [InCharacter: Evaluating personality fidelity in role-playing agents through psychological interviews](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1840–1873, Bangkok, Thailand. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Nathaniel Weir, Ryan Thomas, Randolph d’Amore, Kellie Hill, Benjamin Van Durme, and Harsh Jhamtani. 2024. [Ontologically faithful generation of non-player character dialogues](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9212–9242, Miami, Florida, USA. Association for Computational Linguistics.

Junde Wu, Jiayuan Zhu, and Yunli Qi. 2024. [Medical graph rag: Towards safe medical large language model via graph retrieval-augmented generation](#). *arXiv preprint arXiv:2408.04187*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging LLM-as-a-judge with MT-bench and chatbot arena](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

Hanxun Zhong, Zhicheng Dou, Yutao Zhu, Hongjin Qian, and Ji-Rong Wen. 2022. [Less is more: Learning to refine dialogue history for personalized dialogue generation](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5808–5820, Seattle, United States. Association for Computational Linguistics.

Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Pei Ke, Guanqun Bi, Libiao Peng, JiaMing Yang, Xiyao Xiao, Sahand Sabour, Xiaohan Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. 2024. [CharacterGLM: Customizing social characters with large language models](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry*

Track, pages 1457–1476, Miami, Florida, US. Association for Computational Linguistics.

Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. [ToolQA: A dataset for LLM question answering with external tools](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.

A Comparison of LLM-based Role-playing approaches

Table 6 shows a comparison of different methods used for using LLMs in role-playing tasks. Fine-tuning-based approaches require extensive data collection and are computationally expensive, and they often fail to generalize to roles beyond the training corpus, as each character has a distinct knowledge. Moreover, LLMs inherently encode vast general knowledge, which they may draw upon when answering queries—often leading to fabricated or out-of-scope content. Defining clear character boundaries remains a challenge for fine-tuning-based approaches. Retrieval-based methods eliminate the need for model training and costly data labeling. However, their effectiveness depends on efficiently retrieving query-relevant context from a large character knowledge base through a robust indexing system.

B Dataset Statistics

The statistics of our experimental datasets are illustrated in Table 7. In our experiment, recruiting evaluators who can recall the complete knowledge base of a specific character is challenging, and web searches are often required during evaluation. For instance, assessing a batch of 357 response in the RoleBench-Zh dataset takes approximately **three hours** per evaluation session; The cost of evaluating LLM generation of CharacterLLM dataset with GPT-4 is approximately 5 US dollars.

Table 7: Statistics of the experimental datasets.

Datasets	#Roles	In Scope	Out of Scope
Harry Potter	7	140	-
RoleBench-Zh	5	240	117
Character-LLM	9	814	45

C Evaluation Process

To judge the generated responses according to the above metrics, we make use of GPT-4o to act as a judge LLM by rating the responses. Powerful LLMs such as GPT-4 have been widely employed as evaluators in recent studies (Shao et al., 2023; Dai et al., 2024; Lu et al., 2024; Wang et al., 2024a) where GPT-4 is prompted to give scores for generated output on a defined scale, or to compare

responses and select which one is better. However, there are some concerns about the reliability of LLMs to rate generated responses. Therefore, based on recent works that explore the use of LLMs as judges, we adopt a few measures to increase the reliability of the scores in our experiments. First, we prompt the LLM to generate an analysis before it scores the response. This approach follows recent research (Shen et al., 2023; Zheng et al., 2023) and is based on the success of Chain-of-Thought prompting (Wei et al., 2022). Following Ditto (Lu et al., 2024), we set the temperature of GPT-4o to 0.2 to penalize creativity during evaluation.

To avoid biases that judge LLMs may have, such as the “self-enhancement bias” (Zheng et al., 2023), we include humans in the evaluation process to verify the scores produced by the judge LLM. The human evaluator can use the analysis produced by the judge LLM, as well as any other information sources they want to use, to determine whether the score is sensible. The human evaluator can adjust the score if they feel that it is not correct. We use three different prompts to generate scores for each metric, which can be found in Appendix E.

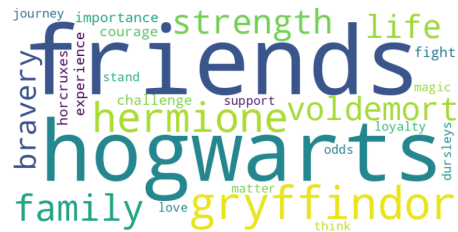


Figure 5: Word cloud for responses generated by GPT-4o mini when role-playing as Harry Potter.

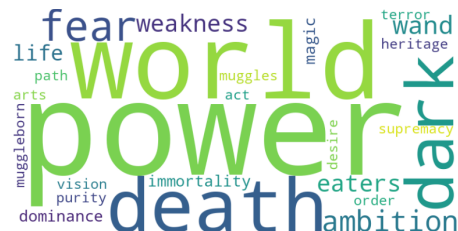


Figure 6: Word cloud for responses generated by GPT-4o mini when role-playing as Voldemort.