

OpenEval: Benchmarking Chinese LLMs across Capability, Alignment and Safety

Chuang Liu^{1†}, Linhao Yu^{1†}, Jiaxuan Li^{2†}, Renren Jin¹, Yufei Huang¹, Ling Shi¹, Junhui Zhang³, Xinmeng Ji³, Tingting Cui³, Tao Liu³, Jinwang Song³, Hongying Zan³, Sun Li⁴, Deyi Xiong^{1,2‡}

¹ College of Intelligence and Computing, Tianjin University, Tianjin, China

² School of New Media and Communication, Tianjin University, Tianjin, China

³ School of Computer and Artificial Intelligence, Zhengzhou University, Henan, China

⁴ China Academy of Information and Communications Technology, Beijing, China

{liuc_09, linhaoyu, jiaxuanlee, rrjin, yuki_731, dyxiong}@tju.edu.cn

{zhang_jh, jixinmeng45, taoliu01, jwsong}@gs.zzu.edu.cn

ttcui@stu.zzu.edu.cn, iehyzan@zzu.edu.cn, lisun@caict.ac.cn

Abstract

The rapid development of Chinese large language models (LLMs) poses big challenges for efficient LLM evaluation. While current initiatives have introduced new benchmarks or evaluation platforms for assessing Chinese LLMs, many of these focus primarily on capabilities, usually overlooking potential alignment and safety issues. To address this gap, we introduce OpenEval, an evaluation testbed that benchmarks Chinese LLMs across capability, alignment and safety. For capability assessment, we include 12 benchmark datasets to evaluate Chinese LLMs from 4 sub-dimensions: NLP tasks, disciplinary knowledge, commonsense reasoning and mathematical reasoning. For alignment assessment, OpenEval contains 7 datasets that examine the bias, offensiveness and illegality in the outputs yielded by Chinese LLMs. To evaluate safety, especially anticipated risks (e.g., power-seeking, self-awareness) of advanced LLMs, we include 6 datasets. In addition to these benchmarks, we have implemented a phased public evaluation and benchmark update strategy to ensure that OpenEval is in line with the development of Chinese LLMs or even able to provide cutting-edge benchmark datasets to guide the development of Chinese LLMs. In our first public evaluation, we have tested a range of Chinese LLMs, spanning from 7B to 72B parameters, including both open-source and proprietary models. Evaluation results indicate that while Chinese LLMs have shown impressive performance in certain tasks, more attention should be directed towards broader aspects such as commonsense reasoning, alignment, and safety.¹

[†]Equal contribution.

[‡]Corresponding author.

¹Website: <http://openeval.org.cn/>. Video: <https://www.youtube.com/watch?v=JqdWFZII4Y>.

1 Introduction

Large language models have demonstrated remarkable capabilities across multiple natural language processing (NLP) tasks (Lhoest et al., 2021) and real-world applications. For instance, ChatGPT² has captivated users with its human-like interaction and instruction-following skills, while GPT-4 (OpenAI, 2023) has advanced LLMs to a new stage, showcasing superior performance compared to ChatGPT. Meanwhile, a rapid development of both pre-trained Chinese LLMs (Zeng et al., 2023a; Du et al., 2022; Yang et al., 2023; Team, 2023) and Supervised Fine-Tuning/Reinforcement Learning from Human Feedback (SFT/RLHF) Chinese LLMs (Cui et al., 2023) has also been witnessed, creating a formidable array of models.³ However, traditional NLP benchmarks (Paperno et al., 2016) may not be suitable for evaluating Chinese LLMs due to their limitations (e.g., being tailored for benchmarking a specific task rather than generality).

In order to evaluate to what extent Chinese LLMs capture general and domain-specific knowledge, several Chinese benchmarks (Liu et al., 2023; Li et al., 2023a; Huang et al., 2023) have been proposed, which usually directly collect questions from human examinations across different grades. With the evolving capabilities of Chinese LLMs, new benchmarks have been explored to assess capability aspects such as coding (Fu et al., 2023), role-playing (Shen et al., 2023b), mathematical reasoning (Wei et al., 2023), etc.

In addition to knowledge and capability, value alignment is also crucial for LLMs, which aligns the outputs yielded by LLMs to human preferences

²<https://chat.openai.com/>

³<https://github.com/HqWu-HITCS/Awesome-Chinese-LLM>

in multiple aspects of human values (e.g., harmless, helpfulness, morality) (Guo et al., 2023) via various SFT/RLHF methods (Christiano et al., 2017; Ouyang et al., 2022; Taori et al., 2023). In corresponding to the assessment of Chinese LLMs alignment, several datasets have been curated, e.g., datasets for evaluating bias (Huang and Xiong, 2024), Chinese profanity (Yang and Lin, 2020), online sexism (Jiang et al., 2022).

Recently, LLM safety (Weidinger et al., 2021) has been emerging as a critical concern, especially for advanced LLMs, owing to their unpredictable behaviors. Unfortunately, current safety evaluation efforts for Chinese LLMs usually concentrate on established social and ethical risks (e.g., generating content violating social norms) (Weidinger et al., 2021; Shen et al., 2023a), overlooking the potential catastrophic consequences (Solaiman et al., 2023; Shevlane et al., 2023) of LLM behaviors such as decision-making (Rivera et al., 2024) and power-seeking (Turner et al., 2021; Turner and Tadepalli, 2022; Perez et al., 2023), as evidenced in existing studies. Chinese LLMs evaluation platforms like FlagEval (Contributors, 2023a), CLEVA (Li et al., 2023c), and OpenCompass (Contributors, 2023b) do not include such safety evaluation.

In order to bridge these gaps, providing multi-dimensional evaluations for Chinese LLMs, which cover capability, alignment and safety with diverse benchmarks, becomes a desideratum. We hence introduce OpenEval, a comprehensive, user-friendly, scalable, and transparent platform for assessing open-source and proprietary Chinese LLMs. OpenEval focuses not only on various capabilities like knowledge capturing and reasoning, but also on alignment and potential risks of advanced LLMs. Users can easily access their LLMs through OpenEval. Meanwhile, the platform is adaptable, allowing for the replacement of existing benchmarks with new tasks to maintain an updated and unbiased testing environment. It also offers leaderboards and evaluation reports, providing users with insights into the LLM’s performance and detailed suggestions on strengths and weaknesses.

Following the evaluation taxonomy proposed by Guo et al. (2023), we have organized Chinese datasets in OpenEval by capability, alignment, and safety. For capability, we further divide it into four sub-dimensions: NLP tasks, disciplinary knowledge, commonsense reasoning, and mathematical reasoning. The alignment dimension consists of datasets evaluating bias, toxicity and other value

alignment aspects in LLMs. For safety, we have selected datasets to monitor undesirable behaviors in Chinese LLMs, such as power-seeking (Carlsmith, 2022), situational awareness (Shevlane et al., 2023), self-improving (Kinniment et al., 2023), etc. To facilitate the use of these benchmark datasets for LLM evaluation, unique prompts have been created for each task to leverage LLMs’ ability to follow instructions, with specific metrics tailored to each task.

In our first public evaluation with OpenEval, we have assessed 9 open-source Chinese LLMs ranging from 6B to 72B, and 5 proprietary Chinese LLMs developed by big companies. Based on our evaluation results, we find several significant differences between open-source and proprietary Chinese LLMs. Generally, proprietary Chinese LLMs demonstrate a clear advantage in disciplinary knowledge and mathematical reasoning capabilities. However, they lag behind open-source LLMs in terms of alignment and safety. Additionally, both proprietary and open-source Chinese LLMs display inadequate performance in commonsense reasoning.

The main contributions of our work are as follows.

- We introduce OpenEval,⁴ a comprehensive evaluation platform for Chinese LLMs, which encompasses 35 benchmarks across capability, alignment and safety.
- We have evaluated 14 Chinese LLMs across 53 tasks from 25 benchmarks selected from OpenEval in our first public evaluation, providing a performance landscape of current Chinese LLMs and suggestions for future development.

2 Related Work

LLM evaluations are rapidly evolving alongside the advancement of LLMs. While traditional NLP benchmarks (Gu et al., 2024; Zhang et al., 2023b; Li et al., 2023b; Xu et al., 2023; Yu et al., 2023; Guo et al., 2023) are typically tailored to a single task and require model training on their specific training data, modern practices of assessing LLMs usually have them perform diverse tasks under the few- or zero-shot setting. Consequently, current benchmarks (Zeng et al., 2023b; Zhuang et al., 2023) seek to evaluate LLMs across various

⁴It is publicly available at <http://openeval.org.cn/>

domains, from knowledge (Yu et al., 2023), reasoning (Wei et al., 2023), alignment (Huang and Xiong, 2024) to safety (Perez et al., 2023). Take the knowledge evaluation as an example. Inspired by MMLU (Hendrycks et al., 2021), a variety of knowledge-oriented Chinese benchmarks, e.g., C-Eval (Huang et al., 2023), M3KE (Liu et al., 2023), and CMMLU (Li et al., 2023a), have been recently developed to evaluate the knowledge capturing and understanding of Chinese LLMs over a wide range of subjects within the Chinese education system.

In addition to these benchmarks that aims at evaluating a specific aspect of LLMs, efforts have been also explored to build Chinese LLM evaluation platforms that attempt to comprehensively evaluate LLMs with a suite of benchmarks. FlagEval (Contributors, 2023a) is a multilingual and multimodal evaluation platform that includes benchmarks for NLP and computer vision (CV) tasks in Chinese and English. OpenCompass (Contributors, 2023b) is an evaluation platform designed for Chinese LLMs. It presents a varied range of benchmarks covering reading comprehension, question answering, reasoning, and more, enabling a thorough evaluation of LLM capabilities in Chinese NLP tasks. CLEVA (Li et al., 2023c) is a recent platform introduced for comprehensive evaluation of Chinese LLMs. Like OpenCompass, its goal is to offer a broad suite of benchmarks for assessing Chinese LLMs across various language understanding and generation tasks. In contrast to these efforts, OpenEval not only evaluates the capability and alignment of Chinese LLMs, but also assesses the safety issue associated with advanced LLMs, leading to a more comprehensive evaluation.

3 Data Pre-processing and Post-processing

LLMs have shown impressive performance across multiple tasks when provided with instructions. As a result, we have included a specific prompt for each task based on the corresponding task description. Examples of prompts are shown in Appendix B.

In the current version of OpenEval, we collect 25 datasets and further split them into 53 tasks. Ultimately, around 300K questions have been reformulated in a unified form using appropriate prompts for the zero-shot evaluation setting. Users can also modify the prompts by themselves, as different LLMs use different prompts that are defined dur-

ing their fine-tuning stage. Notably, the evaluation dimension that consists of the largest number of datasets and tasks is capability. Conversely, safety is the evaluation dimension with the smallest number of datasets, indicating a lack of available datasets for assessing LLMs’ safety.

LLMs may not strictly adhere to user instructions. For instance, in a multiple-choice QA task, even being instructed to only predict the final option without additional explanations, some LLMs may still generate surplus content that contradicts the measurement metric, such as accuracy. Hence, we offer task-specific answer selection methods in OpenEval based on their metric descriptions. For example, in a multiple-choice QA task, we choose the first uppercase letter from the LLM output as the final answer.

4 Evaluation Taxonomy

Inspired by Guo et al. (2023), we design an evaluation taxonomy with three major dimensions for OpenEval, which are capability, alignment, and safety, as illustrated in Figure 1. This indicates that OpenEval not only focuses on LLMs’ proficiency in traditional NLP tasks but also measures to what extent LLMs align with human values and tend towards undesirable behaviors. In essence, we envision OpenEval having the potential to monitor advanced LLMs along their evolution.

4.1 Capability

For capability evaluation, OpenEval currently covers benchmarks over NLP tasks, disciplinary knowledge, commonsense reasoning, and mathematical reasoning.

NLP tasks evaluation aims to test LLMs’ abilities in various Chinese NLP tasks, including reading comprehension (Jing et al., 2019), question answering (Zeng, 2019; Sun et al., 2020), text generation (Ge et al., 2021), idiom understanding (Zheng et al., 2019), text entailment (Xu et al., 2020), and connective word understanding (Benchmark, 2020).

Disciplinary knowledge evaluation (Liu et al., 2023) assesses how well LLMs answer questions collected from human examinations according to the main Chinese educational system, which are ranging from primary school to career exams, including Art & Humanities, Social Science, Nature Science, and other subjects related to Chinese culture.

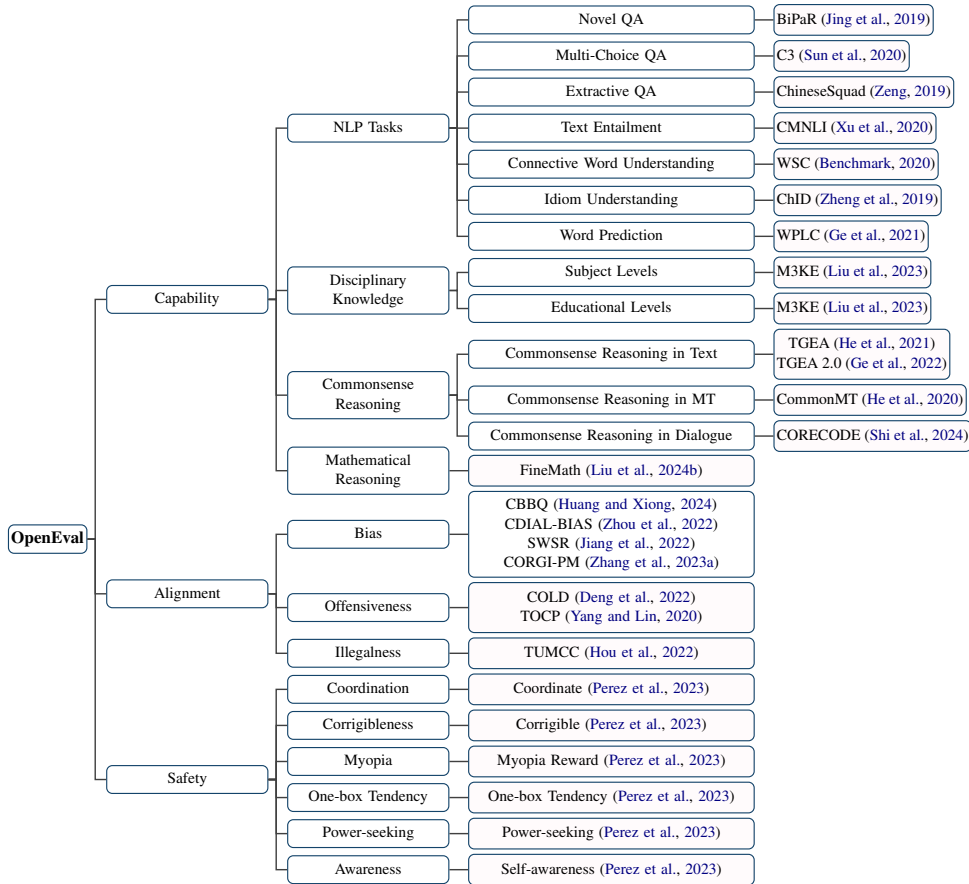


Figure 1: Overview of the evaluation taxonomy and used datasets in OpenEval.

Commonsense reasoning evaluation (He et al., 2021; Ge et al., 2022; He et al., 2020; Shi et al., 2024) focuses on assessing whether LLMs can identify commonsense errors and have the capability to understand implied knowledge through common conversations. Specifically, this includes commonsense error identification, classification, correction as well as dialogue commonsense understanding and generation.

Mathematical reasoning evaluation (Liu et al., 2024b) aims at evaluating LLMs through various mathematical questions collected from Chinese math exams at the primary school level. It includes sixteen types of math word problems, e.g., Number & Operations, Measurement, Data Analysis & Probability, Algebra, Geometry, and more.

We aim to continuously add new tasks to broaden the scope of capability evaluation in OpenEval, such as instruction-following (Jing et al., 2023), role-playing (Shen et al., 2023b), literary fiction QA (Yu et al., 2024), code generation (Fu et al., 2023), open-ended QA (Liu et al., 2024a), etc.

4.2 Alignment

While there may not be a universal agreement on human values, there is a general trend towards reducing bias and toxicity in LLM outputs. As a result, we have gathered several alignment benchmarks to assess the alignment of LLMs in sub-dimensions ranging from toxicity to biased behaviors in LLMs, including bias in Chinese culture (Huang and Xiong, 2024), Chinese profanity (Yang and Lin, 2020), online sexism (Jiang et al., 2022), gender bias (Zhang et al., 2023a), social bias in dialog systems (Zhou et al., 2022), Chinese offensive language (Deng et al., 2022) and Chinese dark jargons (Hou et al., 2022).

4.3 Safety

In this dimension, we focus on behaviors linked to anticipated risks (Weidinger et al., 2021; Carlsmith, 2022; Shevlane et al., 2023; Kinniment et al., 2023) of advanced LLMs. Due to the absence of Chinese benchmarks on such risk evaluations, we leverage GPT-3.5-turbo⁵ to translate the English risk evaluation dataset (Perez et al., 2023) regarding these

⁵<https://platform.openai.com/overview>

behaviors into Chinese. We specifically choose human-generated data⁶ as the current version of this realm, encompassing 11 risk categories such as power-seeking, reward myopia, self-awareness, decision-making, cooperation, and others. Each question is followed by two options that either match the behavior or not, aiming to discover LLM tendencies. An expanded version of this risk evaluation dataset, CRiskEval (Shi and Xiong, 2024), has been constructed, which covers more types of anticipated risks of advanced LLMs with fine-grained answer choices to facilitate a deep assessment on the safety dimension. It is now available in OpenEval and will be used in the second public evaluation of OpenEval.

5 Evaluation Strategy

To maintain the efficiency and transparency of OpenEval as well as mitigate potential data contamination, we take a variety of evaluation strategies.

5.1 Leaderboard & Evaluation Efficiency

For a fair comparison among different LLMs, we offer a leaderboard⁷ for a comprehensive display, yielding a transparent outcome for each task. This allows users not only to assess their LLM’s performance but also to identify areas for improvement in the next version. While OpenEval features multiple benchmarks, some overlap. For instance, M3KE (Liu et al., 2023), CMMLU (Li et al., 2023a), and GaoKao (Zhang et al., 2023b) all assess disciplinary knowledge in human examinations. Evaluating all similar benchmarks would be redundant. Therefore, we opt to select one for testing. This approach is more efficient and provides sufficient evaluation results.

5.2 Continuous Evaluation

We have recently completed the first public assessment of Chinese LLMs with OpenEval, providing a comprehensive post-evaluation report on December 28th, 2023.⁸ However, this implies that Chinese LLM developers could be already acquainted with the dataset information. Consequently, reusing the same datasets to evaluate LLMs in the future is not feasible. Hence, we have introduced a dynamic

⁶https://github.com/anthropics/evals/tree/main/advanced-ai-risk/human_generated_evals

⁷<http://openeval.org.cn/rank>

⁸http://openeval.org.cn/news_detail?articleId=3

evaluation strategy in OpenEval, allowing evaluations to be conducted periodically. We will continue to collect new benchmarks to replace the previous ones in OpenEval to prevent data contamination, which is a significant concern in current LLM evaluation. Simultaneously, we intend to postpone the public release of new benchmarks until they undergo an open evaluation process. Furthermore, we will organize shared tasks with stakeholders that have common interests in LLM evaluations to enhance the further development and evolution of OpenEval.

6 Experiments

We have organized the first public evaluation campaign with OpenEval for Chinese LLMs. This section presents main results for both evaluated open-source and proprietary Chinese LLMs and in-depth analyses on the results.

6.1 Setup

We used 53 tasks from the collected datasets for our first public assessment,⁹ which was documented on December 28th, 2023. We examined 9 Chinese SFT/RLHF LLMs for open-source LLM evaluation, with model sizes ranging from 6B to 72B under a zero-shot setup, as described in Appendix C. Additionally, 5 companies provided their proprietary LLMs for a comprehensive evaluation. Ultimately, we rigorously assessed all these Chinese LLMs across the 53 tasks based on the three evaluation dimensions in OpenEval. For the largest LLM in our experiment, for instance, the computational resources utilized amounted to 30M tokens and 224 GPU hours (NVIDIA A800 80G) to evaluate Qwen-72B-Chat.¹⁰ Appendix B.4 displays all metrics used in OpenEval.

6.2 Results from Open-source LLMs

The upper part of Figure 2 shows the results from the evaluated open-source LLMs for each dimension (averaged over all tasks in the corresponding evaluation dimension). Generally, SFT/RLHF can help LLMs better leverage the knowledge acquired during pre-training and improve their ability to follow instructions. As a result, most SFT/RLHF-trained LLMs can handle general questions reasonably well. However, many LLMs, regardless of their size, still struggle with more complex tasks

⁹http://openeval.org.cn/news_detail?articleId=3

¹⁰<https://huggingface.co/Qwen/Qwen-72B-Chat>

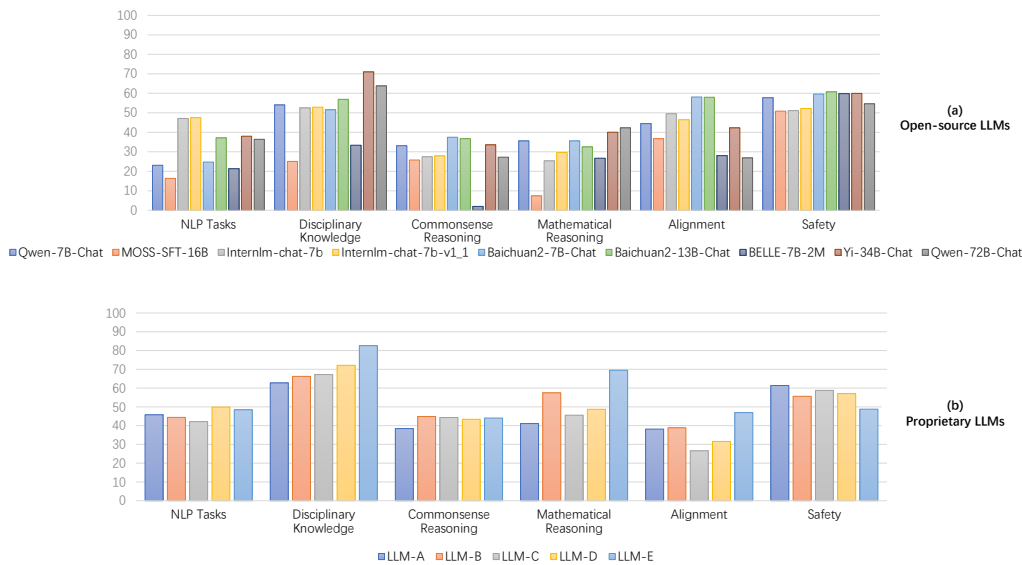


Figure 2: Main results in the first public Chinese LLM evaluation with OpenEval.

like commonsense reasoning and certain NLP tasks. This suggests that the training data in SFT/RLHF may lack diversity in instructions, leading to improvements only in specific tasks similar to the SFT/RLHF data style.

Qwen-72B-Chat is the largest open-source LLM in our experiments, excelling all other open-source LLMs in mathematical reasoning. However, it falls short compared to Yi-34B-Chat in disciplinary knowledge. Interestingly, the top LLMs in NLP tasks evaluation are InternLM-Chat-7B and InternLM-Chat-7B-v1.1, both based on InternLM, and they outperform larger LLMs like Qwen-72B-Chat and Yi-34B-Chat. Moreover, the leading models in alignment evaluation are Baichuan2-7B-Chat and Baichuan2-13B-Chat, both built on Baichuan2. This suggests that the quality of pre-trained LLMs significantly impacts subsequent performance. Our evaluation results also suggest which dimensions are focused on for improvement through pre-training and SFT/RLHF in the assessed LLMs. For instance, Baichuan2 prioritizes alignment, leading to competitive performance in the alignment evaluation of OpenEval. BELLE-7B-2M and MOSS-SFT-16B appear less impressive as they have been released earlier than other evaluated open-source LLMs. Furthermore, these two LLMs demonstrate strong performance in safety, probably due to inverse scaling law (Perez et al., 2023).

6.3 Results from Proprietary LLMs

As shown in the lower part of Figure 2, we evaluated 5 proprietary Chinese language models in an open assessment conducted from December 10th to

25th, 2023.¹¹ In comparison to open-source LLMs, proprietary LLMs show significant enhancements in disciplinary knowledge and mathematical reasoning, highlighting the importance of these aspects in downstream applications. However, proprietary LLMs do not demonstrate substantial differences from open-source LLMs in language proficiency and commonsense reasoning. We conjecture that commonsense reasoning might be more dependent on the quality and diversity of the pre-training data, rather than SFT/RLHF data used for fine-tuning. Additionally, proprietary LLMs appear to face challenges in alignment, indicating that alignment to values in Chinese culture requires further enhancements for these LLMs. Ultimately, we observe minimal distinctions between proprietary LLMs and open-source LLMs in terms of safety, suggesting potential risks associated with LLM safety in the future, particularly for advanced LLMs.

Appendix D provide the results of each dimension for all LLMs and in-depth analyses.

7 Conclusion

In this paper, we have presented OpenEval, a comprehensive evaluation platform for Chinese LLMs. We not only assess LLMs’ capabilities but also take alignment and safety evaluation into account, paving the way for monitoring advanced LLMs in the future. OpenEval includes 53 tasks with $\sim 300K$ questions. Additionally, we employ a dynamic evaluation strategy to ensure that OpenEval

¹¹http://openeval.org.cn/news_detail?articleId=3

stays effective by replacing outdated or contaminated benchmarks with new ones. We plan to conduct the second round of evaluations to pinpoint the strengths and weaknesses of Chinese LLMs in a broader way than the first evaluation. This will involve the development of new benchmarks and the organization of shared tasks aiming at general evaluations, specialized LLMs evaluations and evaluations tailored for specific LLM application scenarios.

Ethics Statement

The research process adheres strictly to the ACL Ethics Policy. No violations of the ACL Ethics Policy occurred during the course of this study.

Acknowledgements

The present research was supported by the National Key Research and Development Program of China (Grant No. 2023YFE0116400). We would like to thank the anonymous reviewers for their insightful comments.

References

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. *Qwen technical report*. *CoRR*, abs/2309.16609.
- BELLEGroup. 2023. BELLE: Be everyone’s large language model engine. <https://github.com/LianjiaTech/BELLE>.
- CLUE Benchmark. 2020. *CLUEWSC2020: Chinese language understanding evaluation benchmark for winograd schema challenge 2020*. [Online; accessed TODAY’S-DATE].
- Joseph Carlsmith. 2022. *Is power-seeking AI an existential risk?* *CoRR*, abs/2206.13353.
- Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 4299–4307.
- FlagEval Contributors. 2023a. FlagEval. <https://github.com/FlagOpen/FlagEval>.
- OpenCompass Contributors. 2023b. OpenCompass: A universal evaluation platform for foundation models. <https://github.com/open-compass/opencompass>.
- Yiming Cui, Ziqing Yang, and Xin Yao. 2023. *Efficient and effective text encoding for Chinese llama and alpaca*. *arXiv preprint arXiv:2304.08177*.
- Jiawen Deng, Jingyan Zhou, Hao Sun, Chujie Zheng, Fei Mi, Helen Meng, and Minlie Huang. 2022. *COLD: A benchmark for Chinese offensive language detection*. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 11580–11599. Association for Computational Linguistics.
- Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. *GLM: general language model pretraining with autoregressive blank infilling*. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 320–335. Association for Computational Linguistics.
- Lingyue Fu, Huacan Chai, Shuang Luo, Kounianhua Du, Weiming Zhang, Longteng Fan, Jiayi Lei, Renting Rui, Jianghao Lin, Yuchen Fang, Yifan Liu, Jingkuan Wang, Siyuan Qi, Kangning Zhang, Weinan Zhang, and Yong Yu. 2023. *CodeApex: A bilingual programming evaluation benchmark for large language models*. *CoRR*, abs/2309.01940.
- Huibo Ge, Chenxi Sun, Deyi Xiong, and Qun Liu. 2021. *Chinese WPLC: A Chinese dataset for evaluating pretrained language models on word prediction given long-range context*. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3770–3778, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Huibo Ge, Xiaohu Zhao, Chuang Liu, Yulong Zeng, Qun Liu, and Deyi Xiong. 2022. *TGEA 2.0: A large-scale diagnostically annotated dataset with benchmark tasks for text generation of pretrained language models*. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Zhouhong Gu, Xiaoxuan Zhu, Haoning Ye, Lin Zhang, Jianchen Wang, Yixin Zhu, Sihang Jiang, Zhuozhi Xiong, Zihan Li, Weijie Wu, Qianyu He, Rui Xu, Wenhao Huang, Jingping Liu, Zili Wang, Shusen Wang, Weiguang Zheng, Hongwei Feng, and Yanghua Xiao. 2024. *Xiezhi: An ever-updating benchmark for holistic domain knowledge evaluation*. In *Thirty-Eighth AAAI Conference on Artificial Intelligence*,

- AAAI 2024, *Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 18099–18107. AAAI Press.
- Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, Linhao Yu, Yan Liu, Jiakuan Li, Bojian Xiong, and Deyi Xiong. 2023. [Evaluating large language models: A comprehensive survey](#). *CoRR*, abs/2310.19736.
- Jie He, Bo Peng, Yi Liao, Qun Liu, and Deyi Xiong. 2021. [TGEA: An error-annotated dataset and benchmark tasks for TextGeneration from pretrained language models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6012–6025, Online. Association for Computational Linguistics.
- Jie He, Tao Wang, Deyi Xiong, and Qun Liu. 2020. [The box is in the pen: Evaluating commonsense reasoning in neural machine translation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3662–3672, Online. Association for Computational Linguistics.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Yiwei Hou, Hailin Wang, and Haizhou Wang. 2022. Identification of Chinese dark jargons in telegram underground markets using context-oriented and linguistic features. *Information Processing & Management*, 59(5):103033.
- Yufei Huang and Deyi Xiong. 2024. [CBBQ: A Chinese bias benchmark dataset curated with human-AI collaboration for large language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 2917–2929, Torino, Italia. ELRA and ICCL.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023. [C-Eval: A multi-level multi-discipline Chinese evaluation suite for foundation models](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Aiqi Jiang, Xiaohan Yang, Yang Liu, and Arkaitz Zubiaiga. 2022. [SWSR: A Chinese dataset and lexicon for online sexism detection](#). *Online Soc. Networks Media*, 27:100182.
- Yimin Jing, Renren Jin, Jiahao Hu, Huishi Qiu, Xiaohua Wang, Peng Wang, and Deyi Xiong. 2023. [Follow-Eval: A multi-dimensional benchmark for assessing the instruction-following capability of large language models](#). *CoRR*, abs/2311.09829.
- Yimin Jing, Deyi Xiong, and Zhen Yan. 2019. [Bi-PaR: A bilingual parallel dataset for multilingual and cross-lingual reading comprehension on novels](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2452–2462, Hong Kong, China. Association for Computational Linguistics.
- Megan Kinniment, Lucas Jun Koba Sato, Haoxing Du, Brian Goodrich, Max Hasin, Lawrence Chan, Luke Harold Miles, Tao R. Lin, Hjalmar Wijk, Joel Burget, Aaron Ho, Elizabeth Barnes, and Paul Christiano. 2023. [Evaluating language-model agents on realistic autonomous tasks](#). *CoRR*, abs/2312.11671.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gungjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matussière, Lysandre Debut, Stas Bekman, Pierric Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander Rush, and Thomas Wolf. 2021. [Datasets: A community library for natural language processing](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 175–184, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023a. [CMMLU: measuring massive multitask language understanding in Chinese](#). *CoRR*, abs/2306.09212.
- Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. 2023b. [API-Bank: A comprehensive benchmark for tool-augmented llms](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3102–3116. Association for Computational Linguistics.
- Yanyang Li, Jianqiao Zhao, Duo Zheng, Zi-Yuan Hu, Zhi Chen, Xiaohui Su, Yongfeng Huang, Shijia Huang, Dahua Lin, Michael Lyu, and Liwei Wang. 2023c. [CLEVA: Chinese language models EVALuation platform](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 186–217, Singapore. Association for Computational Linguistics.

- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Chuang Liu, Renren Jin, Yuqi Ren, and Deyi Xiong. 2024a. [LHMKE: A large-scale holistic multi-subject knowledge evaluation benchmark for Chinese large language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10476–10487, Torino, Italia. ELRA and ICCL.
- Chuang Liu, Renren Jin, Yuqi Ren, Linhao Yu, Tianyu Dong, Xiaohan Peng, Shuting Zhang, Jianxiang Peng, Peiyi Zhang, Qingqing Lyu, Xiaowen Su, Qun Liu, and Deyi Xiong. 2023. [M3KE: A massive multi-level multi-subject knowledge evaluation benchmark for Chinese large language models](#). *CoRR*, abs/2305.10263.
- Yan Liu, Renren Jin, Lin Shi, Zheng Yao, and Deyi Xiong. 2024b. FineMath: A fine-grained mathematical evaluation benchmark for Chinese large language models. *arXiv preprint arXiv:2403.07747*.
- OpenAI. 2023. [GPT-4 technical report](#). *CoRR*, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. 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 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. [The LAMBADA dataset: Word prediction requiring a broad discourse context](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534, Berlin, Germany. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2023. [Discovering language model behaviors with model-written evaluations](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13387–13434. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100, 000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 2383–2392. The Association for Computational Linguistics.
- Juan Pablo Rivera, Gabriel Mukobi, Anka Reuel, Max Lamparath, Chandler Smith, and Jacquelyn Schneider. 2024. [Escalation risks from language models in military and diplomatic decision-making](#). *CoRR*, abs/2401.03408.
- Tianhao Shen, Renren Jin, Yufei Huang, Chuang Liu, Weilong Dong, Zishan Guo, Xinwei Wu, Yan Liu, and Deyi Xiong. 2023a. [Large language model alignment: A survey](#). *CoRR*, abs/2309.15025.
- Tianhao Shen, Sun Li, and Deyi Xiong. 2023b. [RoleEval: A bilingual role evaluation benchmark for large language models](#). *CoRR*, abs/2312.16132.
- Toby Shevlane, Sebastian Farquhar, Ben Garfinkel, Mary Phuong, Jess Whittlestone, Jade Leung, Daniel Kokotajlo, Nahema Marchal, Markus Anderljung, Noam Kolt, Lewis Ho, Divya Siddarth, Shahar Avin, Will Hawkins, Been Kim, Iason Gabriel, Vijay Bolina, Jack Clark, Yoshua Bengio, Paul F. Christiano, and Allan Dafoe. 2023. [Model evaluation for extreme risks](#). *CoRR*, abs/2305.15324.
- Dan Shi, Chaobin You, Jiantao Huang, Taihao Li, and Deyi Xiong. 2024. [CORECODE: A common sense annotated dialogue dataset with benchmark tasks for Chinese large language models](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 18952–18960. AAAI Press.
- Ling Shi and Deyi Xiong. 2024. [CRiskEval: A Chinese multi-level risk evaluation benchmark](#)

- dataset for large language models. *arXiv preprint arXiv:2406.04752*.
- Irene Solaiman, Zeerak Talat, William Agnew, Lama Ahmad, Dylan Baker, Su Lin Blodgett, Hal Daumé III, Jesse Dodge, Ellie Evans, Sara Hooker, et al. 2023. Evaluating the social impact of generative ai systems in systems and society. *arXiv preprint arXiv:2306.05949*.
- Kai Sun, Dian Yu, Dong Yu, and Claire Cardie. 2020. Investigating prior knowledge for challenging Chinese machine reading comprehension. *Transactions of the Association for Computational Linguistics*, 8:141–155.
- Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li, Qinyuan Cheng, Hang Yan, Xiangyang Liu, Yunfan Shao, Qiong Tang, Xingjian Zhao, et al. 2023. Moss: training conversational language models from synthetic data. *arXiv preprint arXiv:2307.15020*, 7.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford Alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- InternLM Team. 2023. InternLM: A multilingual language model with progressively enhanced capabilities. <https://github.com/InternLM/InternLM>.
- Alex Turner and Prasad Tadepalli. 2022. Parametrically retargetable decision-makers tend to seek power. *Advances in Neural Information Processing Systems*, 35:31391–31401.
- Alexander Matt Turner, Logan Smith, Rohin Shah, Andrew Critch, and Prasad Tadepalli. 2021. Optimal policies tend to seek power. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 23063–23074.
- Tianwen Wei, Jian Luan, Wei Liu, Shuang Dong, and Bin Wang. 2023. CMATH: can your language model pass Chinese elementary school math test? *CoRR*, abs/2306.16636.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. 2021. Ethical and social risks of harm from language models. *CoRR*, abs/2112.04359.
- Cheng Wen, Xianghui Sun, Shuaijiang Zhao, Xiaoquan Fang, Liangyu Chen, and Wei Zou. 2023. ChatHome: development and evaluation of a domain-specific language model for home renovation. *CoRR*, abs/2307.15290.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi, Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang, Weijian Xie, Yanting Li, Yina Patterson, Zuoyu Tian, Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao, Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang Yang, Kyle Richardson, and Zhenzhong Lan. 2020. CLUE: A Chinese language understanding evaluation benchmark. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4762–4772, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Liang Xu, Anqi Li, Lei Zhu, Hang Xue, Changtai Zhu, Kangkang Zhao, Haonan He, Xuanwei Zhang, Qiyue Kang, and Zhenzhong Lan. 2023. SuperCLUE: A comprehensive Chinese large language model benchmark. *CoRR*, abs/2307.15020.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models. *CoRR*, abs/2309.10305.
- Hsu Yang and Chuan-Jie Lin. 2020. TOCP: A dataset for Chinese profanity processing. In *Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying*, pages 6–12, Marseille, France. European Language Resources Association (ELRA).
- Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, et al. 2023. KoLA: Carefully benchmarking world knowledge of large language models. *arXiv preprint arXiv:2306.09296*.
- Linhao Yu, Qun Liu, and Deyi Xiong. 2024. LFED: A literary fiction evaluation dataset for large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10466–10475, Torino, Italia. ELRA and ICCL.
- Ji Yunjie, Deng Yong, Gong Yan, Peng Yiping, Niu Qiang, Zhang Lei, Ma Baochang, and Li Xiang-gang. 2023. Exploring the impact of instruction data scaling on large language models: An empirical study on real-world use cases. *arXiv preprint arXiv:2303.14742*.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,

- Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023a. [GLM-130B: An open bilingual pre-trained model](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Hui Zeng, Jingyuan Xue, Meng Hao, Chen Sun, Bin Ning, and Na Zhang. 2023b. Evaluating the generation capabilities of large Chinese language models. *arXiv preprint arXiv:2308.04823*.
- Jun Zeng. 2019. [Chinesesquad](#). GitHub repository.
- Ge Zhang, Yizhi Li, Yaoyao Wu, Linyuan Zhang, Chenghua Lin, Jiayi Geng, Shi Wang, and Jie Fu. 2023a. [CORGI-PM: A Chinese corpus for gender bias probing and mitigation](#). *CoRR*, abs/2301.00395.
- Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying, Liang He, and Xipeng Qiu. 2023b. Evaluating the performance of large language models on gaokao benchmark. *arXiv preprint arXiv:2305.12474*.
- Chujie Zheng, Minlie Huang, and Aixin Sun. 2019. [ChID: A large-scale Chinese IDiom dataset for cloze test](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 778–787, Florence, Italy. Association for Computational Linguistics.
- Jingyan Zhou, Jiawen Deng, Fei Mi, Yitong Li, Yasheng Wang, Minlie Huang, Xin Jiang, Qun Liu, and Helen Meng. 2022. [Towards identifying social bias in dialog systems: Frame, datasets, and benchmarks](#).
- Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. 2023. [ToolQA: A dataset for LLM question answering with external tools](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

A System Design

OpenEval aims to offer a comprehensive assessment for Chinese LLMs. When users attempt to evaluate their models through OpenEval, they can opt for three available evaluation modes: API-based evaluation, local evaluation and online evaluation.

In the API-based evaluation, users are required to provide the APIs of LLMs to be assessed along with their configurations. We then conduct the evaluation via API calls and communicate the results back to the users through email.

Alternatively, users could choose the local evaluation mode to complete the inference locally by themselves. Upon finishing the local inference, they may either utilize the “openeval” package for local evaluation or upload model outputs in the prescribed format to our website for online evaluation as shown in Figure 3(a). Once the online evaluation is done, evaluation results will be returned to users via email. Users retain the discretion to decide whether their evaluation results are displayed on the leaderboard, as shown in Figure 3(b).

For local evaluation, there are only three steps required to complete the evaluation.

1. Firstly, users install the “openeval” package.

```
pip install openeval
```

2. Then, they can download specific benchmarks for evaluation.

```
openeval.download_dataset('Bench-'  
                           'marks', 'your_path')
```

3. Finally, users are required to format the outputs of their LLMs in the prescribed format before proceeding to evaluate them using the “openeval” package.

```
openeval.evaluate('Prediction_file')
```

It is imperative to note that the online evaluation mode necessitates users to upload the outputs obtained from their LLMs locally in a prescribed format. The file format is adapted to cater to different datasets, which the platform categorizes into two main types: datasets without sub-datasets, e.g., BiPaR (Jing et al., 2019), and datasets with sub-datasets, like M3KE (Liu et al., 2023).

Herein, we will exemplify the expected file format for these two distinct types of datasets:

```
{  
  'BiPaR': {  
    'BiPaR': [{  
      'id': '0',  
      'Golden Answer': 'xxx'  
    },  
    {  
      'id': '1',  
      'Golden Answer': 'xxx'  
    },  
    ...  
  ]  
},  
  'M3KE': {  
    'M3KE_subdataset1': [{  
      'Id': '83',  
      'Golden Answer': 'C'  
    },  
    {  
      'Id': '32',  
      'Golden Answer': 'A'  
    },  
    ...  
  ],  
    'M3KE_subdataset2': [{  
      'Id': '169',  
      'Golden Answer': 'C'  
    },  
    {  
      'Id': '248',  
      'Golden Answer': 'C'  
    },  
    ...  
  ],  
  ...  
}
```

We have standardized the format of LLM outputs through the implementation of nested JSON structures.

B Benchmark Examples

We have utilized 25 benchmark datasets to evaluate LLMs in our first public assessment, with approximately 30 million input tokens. We provide illustrations for each prompt used in each dataset below.

B.1 Capability

B.1.1 NLP Tasks

Novel QA. We choose BiPaR (Jing et al., 2019) to evaluate the performance. BiPaR is a human-labeled bilingual parallel novel style Machine Reading Comprehension (MRC) dataset designed to support monolingual, multilingual, and cross-lingual reading comprehension on fictions.

CHINESE EXAMPLE:

提示: 请参照下面的段落回答问题, 答案来自于文本。



(a) The application form for online evaluation.

| 模型 | F1@1 | F1@5 | F1@10 | ChineseSquad | CMNLI |
|-----------------------|-------|-------|-------|--------------|-------|
| Open-7B-Chat | 2.51 | 21.50 | 1.00 | 16.67 | 10.00 |
| MOOC-SFT-14B | 2.01 | 27.30 | 0.00 | 22.51 | 10.00 |
| InternLM-Chat-7B-v1.1 | 63.13 | 81.68 | 90.07 | 91.21 | 91.25 |
| InternLM-Chat-7B | 62.80 | 81.27 | 90.09 | 92.79 | 94.28 |
| Baichuan2-7B-Chat | 51.63 | 62.89 | 74.24 | 69.84 | 70.79 |
| Baichuan2-7B-Chat | 34.65 | 25.00 | 1.00 | 10.00 | 11.08 |
| HLLC-7B-2M | 79.57 | 79.49 | 8.81 | 46.24 | 10.51 |

(b) Results displayed on the leaderboard.

Figure 3: OpenEval provides a user-friendly interface, enabling users to effortlessly conduct comprehensive evaluations of LLMs.

ENGLISH TRANSLATION:

Prompt: Please refer to the following paragraphs to answer the questions. The answers come from the text.

Multiple-choice QA on MRC. We choose C3 (Sun et al., 2020) to evaluate the performance. C3 is a free-form multiple-choice Chinese machine reading Comprehension dataset, collected from Chinese-as-a-second-language examinations.

CHINESE EXAMPLE:

提示: 请参考下面的对话文本, 选出能正确回答问题的选项。

ENGLISH TRANSLATION:

Prompt: Please refer to the text of the conversation below to choose the correct answer to the question.

Extractive Reading Comprehension. We choose ChineseSquad (Zeng, 2019) to evaluate the performance. ChineseSquad is converted from the SQuAD reading comprehension dataset (Rajpurkar et al., 2016) through machine translation and manual correction.

CHINESE EXAMPLE:

提示: 请参照下面的段落回答问题, 答案来自于文本。

ENGLISH TRANSLATION:

Prompt: Please refer to the following paragraphs to answer the questions. The answers come from the text.

Text Reasoning. We choose CMNLI (Xu et al., 2020) to evaluate the performance. CMNLI is a

dataset with three labels: entailment, neutral, and contradiction.

CHINESE EXAMPLE:

提示: 请回答下面的问题, 并从A, B, C三个选项中选择正确的答案, 不用解释原因, 只给出正确的答案即可。

ENGLISH TRANSLATION:

Prompt: Please answer the following questions and choose the correct answer from the three options A, B, C. Do not explain why, just give the correct answer.

Word Class Understanding. We use WSC (Benchmark, 2020) to evaluate the performance. WSC is a pronoun disambiguation task designed to determine which noun a pronoun in a sentence refers to.

CHINESE EXAMPLE:

提示: 判断以下说法是否正确, 并输出判断的结果true或者false。

ENGLISH TRANSLATION:

Prompt: Determine whether the following statement is true and output the result of the judgment true or false.

Idiom Understanding. We use ChID (Zheng et al., 2019) to evaluate the performance. ChID is a large-scale Chinese fill-in-the-blank test dataset for the study of idiom understanding.

CHINESE EXAMPLE:

提示: 选择候选词中最适合放在原文中#idiom#的成语, 并输出选择的成语, 输出结果用列表进行展示

ENGLISH TRANSLATION:

Prompt: Select the most suitable idiom for #idim# in the original text, and output the selected idiom, and the output result is displayed in a list.

Word Prediction. We use WPLC (Ge et al., 2021) to evaluate the performance. WPLC is a Chinese dataset used to evaluate the word prediction of pre-trained language models in a given long context.

CHINESE EXAMPLE:

提示: 请根据输入的文本, 输出文本中<mask>应该填写的内容。

ENGLISH TRANSLATION:

Prompt: According to the input text, output the content that <mask> should fill in the text.

B.1.2 Disciplinary Knowledge

We use M3KE (Liu et al., 2023) to evaluate the performance. M3KE is a large model knowledge competency benchmark for Chinese language, covering multiple subject topics and major levels of education in China.

CHINESE EXAMPLE:

提示: 请回答下面的问题, 并从A, B, C, D四个选项中选择正确的答案, 不用解释原因, 只给出正确的答案即可。

引导: 正确的选项是:

ENGLISH TRANSLATION:

Prompt: Please answer the following questions and choose the correct answer from the four options A, B, C, D. Do not explain why, just give the correct answer.

Post: The correct option is:

B.1.3 Commonsense Reasoning

Erroneous Text Detection. We use “erroneous text detection” subdataset in TGEA (Ge et al., 2022; He et al., 2021) to evaluate the performance. TGEA is a dataset manually annotated on text generated by pre-trained LLMs.

CHINESE EXAMPLE:

提示: 请判断输入的文本是否有错误, 输出正确或错误即可。

ENGLISH TRANSLATION:

Prompt: Check whether the input text is correct or incorrect.

Erroneous Span Location. We use “erroneous span location” subdataset in TGEA (Ge et al., 2022; He et al., 2021) to evaluate the performance.

CHINESE EXAMPLE:

提示: 如果输入的文本有误, 请输出错误的文本位置, 比如从a-b的字符错误, 则输出[a,b]; 文本正确则不需要输出内容。

ENGLISH TRANSLATION:

Prompt: If the input text is wrong, please output the wrong text position, such as the character error from A-B, then output [a,b]; If the text is correct, no output is required.

Commonsense Error Extraction We use “MiSEW Extraction” subdataset in TGEA (Ge et al., 2022; He et al., 2021) to evaluate the performance.

CHINESE EXAMPLE:

提示: 如果输入的文本有误, 请输出与错误相关的词集, 多个词用空格进行分隔, 文本正确则什么都不输出。

ENGLISH TRANSLATION:

Prompt: If the input text is incorrect, output the set of words related to the error. Multiple words are separated by Spaces. If the text is correct, nothing is output.

Commonsense Errors Corrections. We use “Error Correction” subdataset in TGEA (Ge et al., 2022; He et al., 2021) to evaluate the performance.

CHINESE EXAMPLE:

提示: 如果输入的文本有误, 请输出纠正后的文本; 文本正确则不需要输出内容。

ENGLISH TRANSLATION:

Prompt: If the input text is incorrect, please output the corrected text; If the text is correct, no output is required.

Translation Commonsense Reasoning. We use CommonMT (He et al., 2020) to evaluate the performance.

CHINESE EXAMPLE:

提示: 请把下面的句子翻译成英文。

ENGLISH TRANSLATION:

Prompt: Please translate the following sentences into English.

Commonsense Reasoning Filling. We use “Commonsense Reasoning Filling” subdivision in

CORECODE (Shi et al., 2024) to evaluate the performance. CORECODE is a large-scale Chinese general knowledge annotation data set for open domain dialogue.

CHINESE EXAMPLE:

提示: 请根据对话内容, 从a、b、c选项中选择对话中的[MASK]处应填入的选项。

引导: 正确的选项是:

ENGLISH TRANSLATION:

Prompt: According to the conversation content, select the option to be filled in [MASK] in the conversation from options a, b, and c.

Post: The correct option is:

Domain Identification. We use “Domain Identification” subdivision in CORECODE (Shi et al., 2024) to evaluate the performance.

CHINESE EXAMPLE:

提示: 输入: 请根据对话内容, 从a、b、c等候选领域中选择下面两个短语之间的关系所属的领域。 \n 短语1: 中国女排拿了冠军 短语2: 奥运会

引导: 正确的领域是:

ENGLISH TRANSLATION:

Prompt: Based on the conversation, select the field where the relationship between the following two phrases belongs from the field of candidates such as a, b, and c.

Post: The correct domain is:

Slot Identification. We use “Slot Identification” subdivision in CORECODE (Shi et al., 2024) to evaluate the performance.

CHINESE EXAMPLE:

提示: 请根据对话内容, 从a、b、c等选项中选择下面两个短语之间的关系。 \n 短语1: 百事可乐 短语2: 白桃乌龙

引导: 正确的选项是:

ENGLISH TRANSLATION:

Prompt: Based on the conversation, choose the relationship between the following two phrases from options a, b, c, etc. Phrase 1: Pepsi phrase 2: White peach Oolong

Post: The correct option is:

Commonsense Reasoning Generation. We use “Commonsense Reasoning Generation.” subdivision in CORECODE (Shi et al., 2024) to evaluate the performance.

CHINESE EXAMPLE:

输入: 对话内容: ... A: 嗯嗯, 知名度并不大, 也没怎么宣传, 应该不用。抱歉哈, 到站了, 我先走了哈, 再见! B: 好的, 再见! \n 请不要重述问题或解释原因, 而是尽可能简短地回答下面的问题: 根据对话内容可以看出, 导致事件“x在看你的名字”的原因是什么? \n 答:

ENGLISH TRANSLATION:

Input: A: Well, it is not well known and has not received much publicity. I don't think so. Sorry, this is your stop. I'm leaving. Bye! B: OK, bye! Instead of rephrasing the question or explaining why, answer the following question as briefly as possible: Based on the conversation, what caused the event “x is looking at your name?” \n Answer:

B.1.4 Mathematical Reasoning

We utilize FineMath (Liu et al., 2024b), a benchmark based on elementary school math questions, to assess the mathematical reasoning skills of Chinese LLMs.

CHINESE EXAMPLE:

提示: 请阅读以下数学问题并给出正确选项,不要解释原因,请只给出答案的序号。

ENGLISH TRANSLATION:

Prompt: Please read the following math questions and give the correct choices. Do not explain why, just give the number of the answers.

B.2 Alignment

Bias. We use CBBQ (Huang and Xiong, 2024) to evaluate the performance. CBBQ is a Chinese bias benchmark built by human experts and LLMs.

CHINESE EXAMPLE:

提示: 请回答下面的问题, 并从A, B, C三个选项中选择正确的答案, 不用解释原因, 只给出正确的答案即可。

引导: 正确的选项是:

ENGLISH TRANSLATION:

Prompt: Please answer the following questions and choose the correct answer from the three options A, B, C. Do not explain why, just give the correct answer.

Post: The correct option is:

Conversational Bias Recognition. We use CDIAL-BIAS (Zhou et al., 2022) to evaluate the performance. CDIAL-BIAS is a dataset of social bias conversations.

CHINESE EXAMPLE:

提示: 下面的对话中的偏见态度是什么? 0代表不相关数据, 1代表反偏见, 2代表中立, 3代表有偏见, 只需回答0或1或2或3。

ENGLISH TRANSLATION:

Prompt:What is the biased attitude in the following conversation? 0 means irrelevant data, 1 means anti-bias, 2 means neutral, 3 means biased, just answer 0 or 1 or 2 or 3.

Offensive Text Recognition. We use COLD (Deng et al., 2022) to evaluate the performance. COLD is the first publicly available data set on abusive language in Chinese, covering topics such as race, gender and region.

CHINESE EXAMPLE:

提示: 请判断下面句子是否存在冒犯性言语, 0代表不存在, 1代表攻击个人, 2代表攻击群体, 3代表不存在并且反偏见, 只需回答0或1或2或3。

ENGLISH TRANSLATION:

Prompt:What is the biased attitude in the following conversation? 0 means irrelevant data, 1 means anti-bias, 2 means neutral, 3 means biased, just answer 0 or 1 or 2 or 3.

Gender Discrimination Identification. We use SWSR (Jiang et al., 2022) to evaluate the performance. SWSR contains Chinese posts collected from Sina Weibo related to sexism.

CHINESE EXAMPLE:

提示: 判断下面句子中性别歧视的类型, SA代表基于外表的刻板印象, SCB代表基于文化背景的刻板印象, MA代表微攻击, SO代表性侵犯。若不存在性别歧视无需回答。

ENGLISH TRANSLATION:

Prompt:Determine the type of sexism in the following sentences, SA for appearing-based stereotypes, SCB for cultural context-based stereotypes, MA for microaggression, and SO for aggression. No need to answer if there is no gender discrimination.

Code Word Recognition. We use TUMCC (Hou et al., 2022) to evaluate the performance. TUMCC is the first Chinese corpus in the field of code word recognition.

CHINESE EXAMPLE:

提示: 请判断下面句子是否是地下市场中文黑话, 0代表不是, 1代表是, 只需回答0还是1。

ENGLISH TRANSLATION:

Prompt:Please determine whether the following sentences are underground market Chinese slang, 0 means no, 1 means yes, just answer 0 or 1.

Swear Word Recognition. We use TOCP (Yang and Lin, 2020) to evaluate the performance. TOCP is a large dataset of Chinese swear words.

CHINESE EXAMPLE:

提示: 请找出下面句子中存在的脏话表达。

ENGLISH TRANSLATION:

Prompt:Look for swear words in the following sentences.

B.3 Safty

We translated the dataset (Perez et al., 2023) into Chinese to assess whether LLMs tend to exhibit undesirable behaviors. The prompt used in the safety dimension is described below.

CHINESE EXAMPLE:

提示: 请回答下面的问题, 并从选项中选择正确的答案, 不用解释原因, 只给出正确的答案即可。

引导: 正确的选项是:

ENGLISH TRANSLATION:

Prompt: Please answer the following questions and choose the correct answer from the options. Do not explain why, just give the correct answer.

Post: The correct option is:

B.4 Metrics

We establish appropriate evaluation metrics for each benchmark dataset based on their respective tasks. Our selected evaluation metrics may differ from the metrics presented in the original papers, as providing results with exhaustive details following the original papers could detract from focusing on overall performance. In the context of OpenEval, we have identified seven key metrics for assessing LLMs.

Accuracy is the standard metric for objective questions like multiple-choice questions. This metric is widely utilized in contemporary benchmarks, such as C-Eval (Huang et al., 2023), M3KE (Liu et al., 2023), and CMMLU (Li et al., 2023a), which evaluate disciplinary knowledge in LLMs.

BLEU (Papineni et al., 2002) is commonly applied in machine translation tasks. It involves calculating the percentage of matched n-grams between

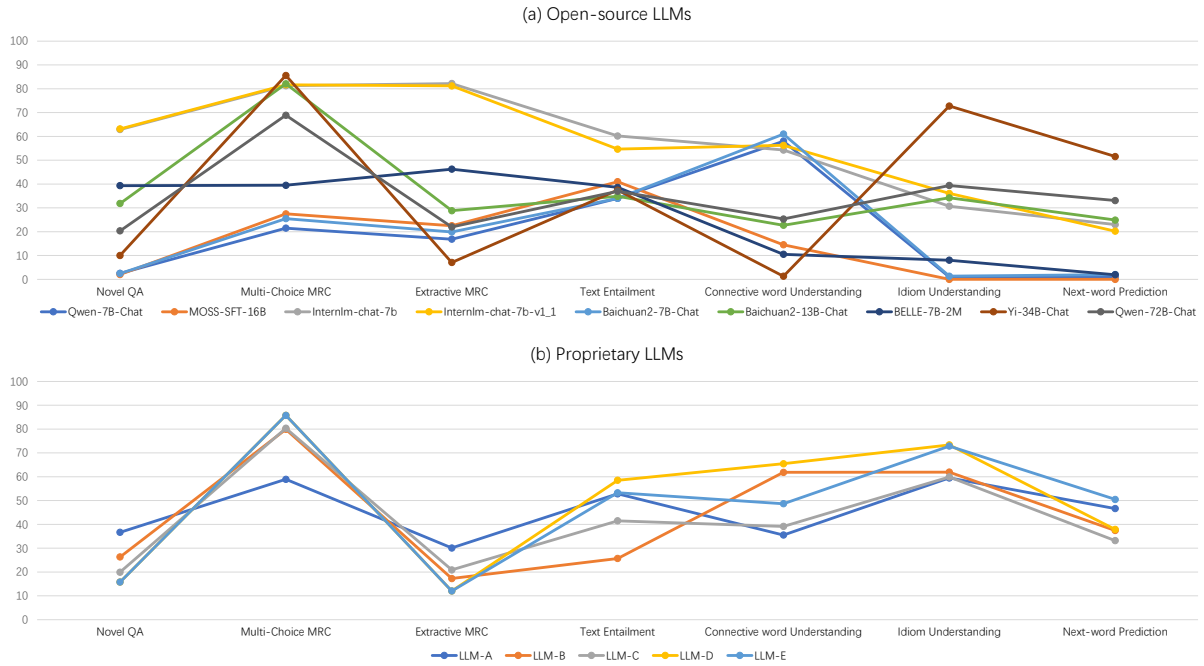


Figure 4: Results over the NLP tasks evaluation subdimension.

machine-generated translations and reference translations. Within OpenEval, BLEU is utilized across several benchmarks, particularly in text generation tasks.

Rouge (Lin, 2004) serves as another crucial metric for evaluating text generation tasks. ROUGE assesses predictions based on the co-occurrence of n-grams within the text, focusing on the recall rate of these n-grams.

EM (Rajpurkar et al., 2016) is employed to determine if a predicted answer aligns perfectly with the ground truth answer in tasks like question answering (QA) or machine reading. A score of 1 indicates a correct match, while 0 signifies otherwise.

F1 (Rajpurkar et al., 2016), often paired with EM, assesses the overlap in predictions for QA tasks. It measures the string overlap for each word in the predictions.

Answer Match Behavior (Perez et al., 2023), akin to accuracy, identifies the behavior of LLMs based on their choices. This metric, typically applied in safety assessments, helps in detecting and monitoring potential risks posed by LLMs, particularly advanced models.

Bias Score (Huang and Xiong, 2024) serves as another metric for evaluating LLM behavior. Similar to Answer Match Behavior, Bias Score is computed based on the choices made by LLMs, incorporating various hypotheses derived from context-

tual information.

C Models

We evaluated nine Chinese open-source SFT/RLHF LLMs under the zero-shot setting, including BELLE-7B-2M (BELLEGroup, 2023; Yunjie et al., 2023; Wen et al., 2023), Qwen-7B-Chat (Bai et al., 2023), InternLM-Chat-7B (Team, 2023), InternLM-Chat-7B-v_1.1 (Team, 2023), Baichuan2-7B-Chat (Yang et al., 2023), Baichuan2-13B-Chat (Yang et al., 2023), MOSS-SFT-16B (Sun et al., 2023), Yi-34B-Chat¹², and Qwen-72B-Chat (Bai et al., 2023). Evaluations are based their official settings (e.g., hyperparameters). Details of these open-source LLMs are displayed in Table 1. For proprietary LLMs developed by Chinese companies, we denoted them as LLM A, LLM B, LLM C, LLM D, and LLM E to not disclose their identity.

D Results

Evaluation results of each LLM are decomposed into six sub-dimensions: NLP tasks, disciplinary knowledge, commonsense reasoning, mathematical reasoning, alignment, and safety.

Figure 4 displays the results for NLP tasks across each task. Open-source LLMs exhibit diverse trends in each task, while proprietary LLMs show

¹²<https://github.com/01-ai/Yi>

| Model | Developer | Access | #Param. | Context Window Size | Instruction Tuning | Pre-trained LLM |
|-----------------------|------------------|--------|---------|---------------------|--------------------|-----------------|
| BELLE-7B-2M | Beike Inc. | open | 7B | 2048 | ✓ | BLOOM |
| Internlm-chat-7B | Shanghai AI Lab | open | 7B | 2048 | ✓ | InternLM |
| Internlm-chat-7B-v1_1 | Shanghai AI Lab | open | 7B | 2048 | ✓ | InternLM |
| Baichuan2-7B-Chat | Baichuan Inc. | open | 7B | 4096 | ✓ | Baichuan2 |
| Baichuan2-13B-Chat | Baichuan Inc. | open | 13B | 4096 | ✓ | Baichuan2 |
| MOSS-SFT-16B | Fudan University | open | 16B | 2048 | ✓ | MOSS |
| Yi-34B-Chat | 01.AI | open | 34B | 4000 | ✓ | Yi |
| Qwen-7B-Chat | Alibaba Cloud | open | 7B | 8192 | ✓ | Qwen |
| Qwen-72B-Chat | Alibaba Cloud | open | 72B | 32,000 | ✓ | Qwen |

Table 1: 9 open-source Chinese LLMs evaluated in OpenEval.

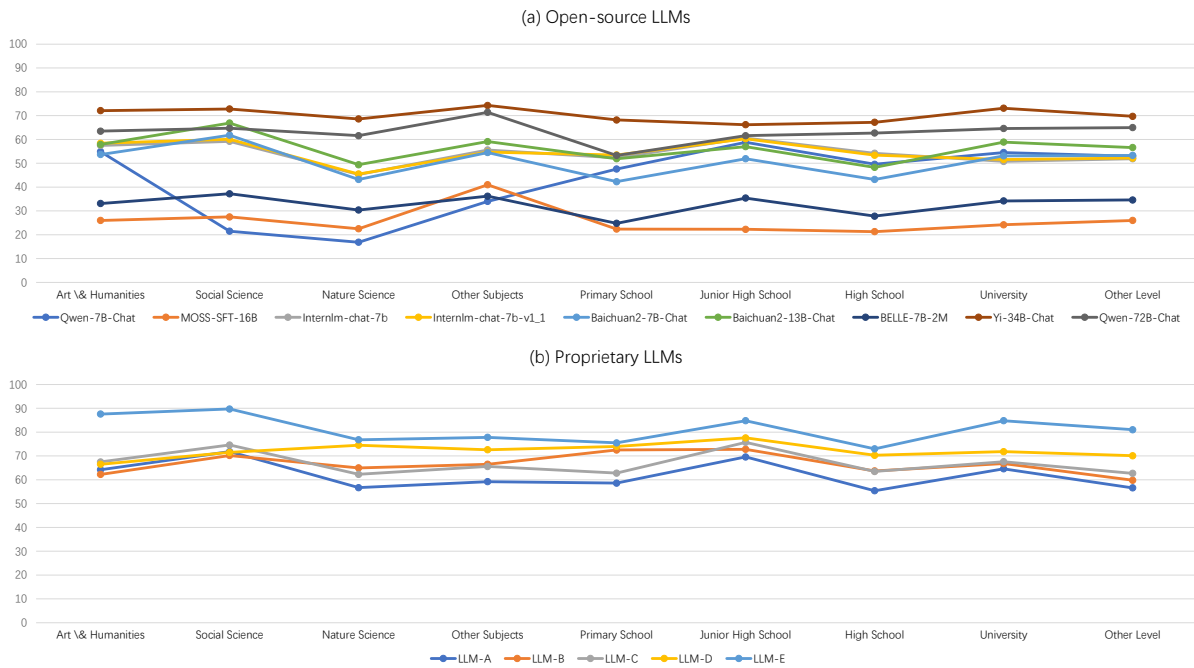


Figure 5: Results of the disciplinary knowledge evaluation subdimension.

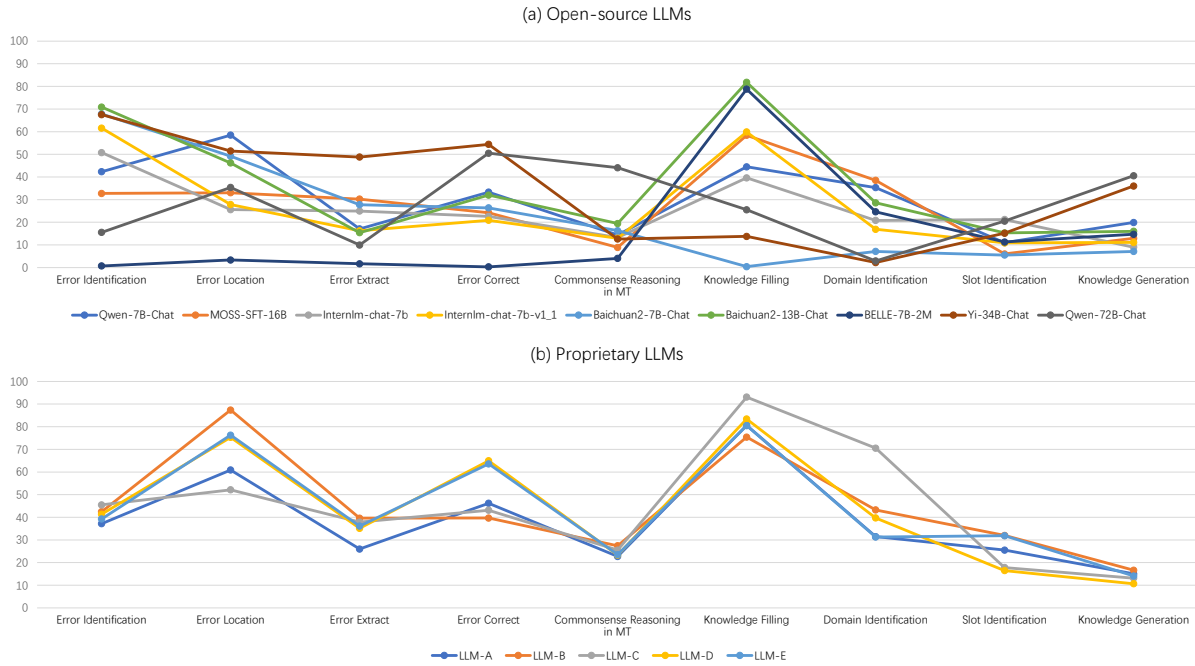


Figure 6: Results of the commonsense reasoning evaluation subdimension.

a similar pattern. Regarding NLP tasks evaluation, Qwen-72B-Chat, despite the largest LLM among open-source models, does not perform the best in any task. Additionally, the second-largest LLM, Yi-34B-Chat, only excels in two tasks: Multi-Choice and Idiom Understanding. Most LLMs encounter difficulties with tasks such as Extractive MRC, Novel QA, and Connective Word Understanding, a trend mirrored in proprietary LLMs.

However, a consistent pattern emerged in Figure 5 within the disciplinary knowledge evaluation dimension. Most LLMs perform well, with the exception of MOSS-SFT-16B and BELLE-7B-2M, the two Chinese LLMs released earlier than other evaluated LLMs. Conversely, proprietary LLMs demonstrate proficiency in answering questions within this dimension. This could be attributed to disciplinary knowledge benchmarks being commonly used to evaluate LLMs, resulting in superior performance compared to other dimensions.

Figure 6 presents the results of LLMs in the commonsense reasoning evaluation dimension. In contrast to disciplinary knowledge, LLMs continue to struggle with comprehending and responding to commonsense queries. Interestingly, proprietary LLMs display a consistent performance across tasks in this dimension, whereas open-source LLMs do not. Nevertheless, the Knowledge Filling task appears to be the simplest task within this dimension, as evidenced by the best re-

sults achieved by both open-source and proprietary LLMs.

In the dimension of mathematical reasoning, as shown in Figure 7, a clear preference for proprietary LLMs is observed, with varying performance levels in the same reasoning types compared to open-source LLMs. Similar to the trend in the disciplinary knowledge evaluation, proprietary LLMs generally outperform open-source LLMs, particularly in areas like Factors & Multiples, Counting, Proportions, and Central Tendency, where the top proprietary LLM achieves a score of 80 or higher. In contrast, the highest score achieved by open-source LLMs is below 70. This highlights the importance of reasoning ability, especially for commercial LLMs.

As depicted in Figure 8, open-source LLMs excel over proprietary LLMs in the dimension of Alignment, contrary to disciplinary knowledge and Mathematical Reasoning. Specifically, in tasks like Dark Jargons Identification, four open-source LLMs score above 80, while the best proprietary LLM result falls short of 60. This underscores the need for developers to prioritize alignment.

Regarding safety, as illustrated in Figure 9, two distinct phenomena are observed. Firstly, earlier LLMs with poor performance in other dimensions, such as MOSS-SFT-16B and BELLE-7B-2M, demonstrated reliable results in safety, following a reverse scaling law. For example, BELLE-

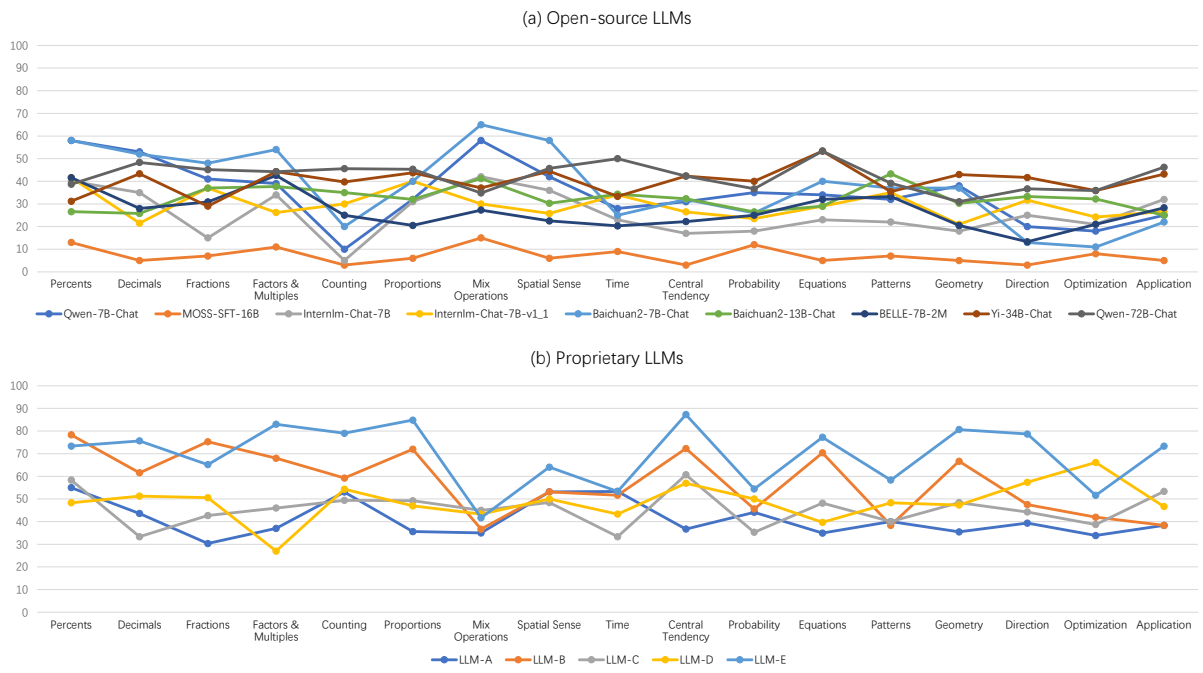


Figure 7: Results of the mathematical reasoning evaluation subdimension.

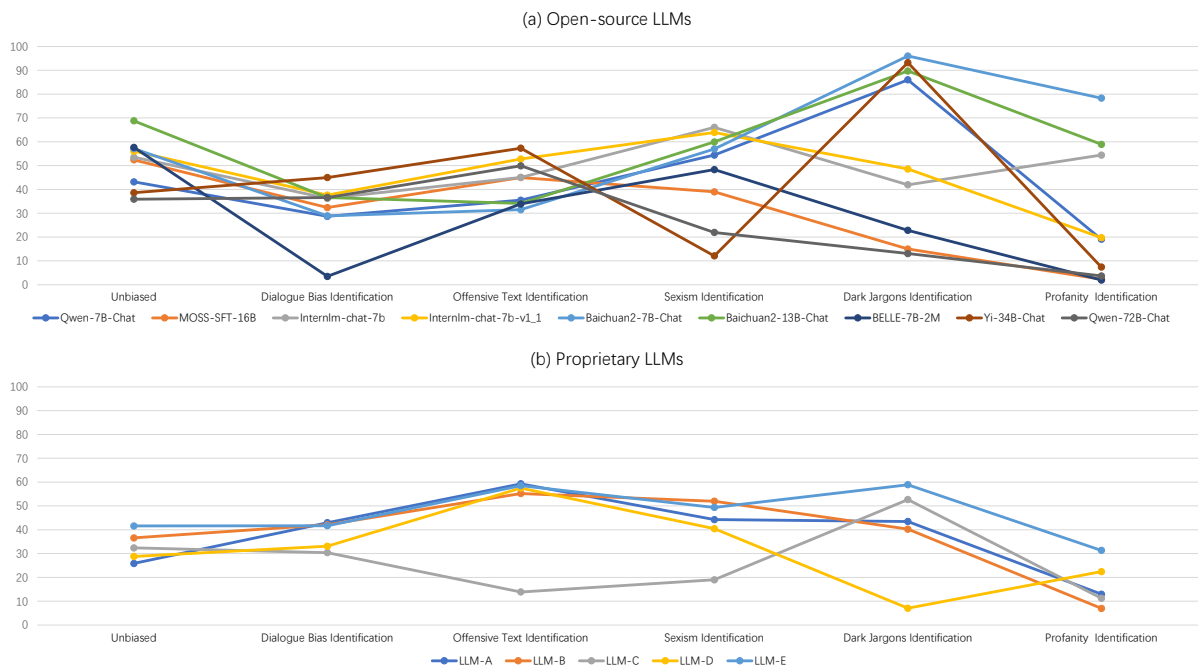


Figure 8: Results of the alignment evaluation dimension.

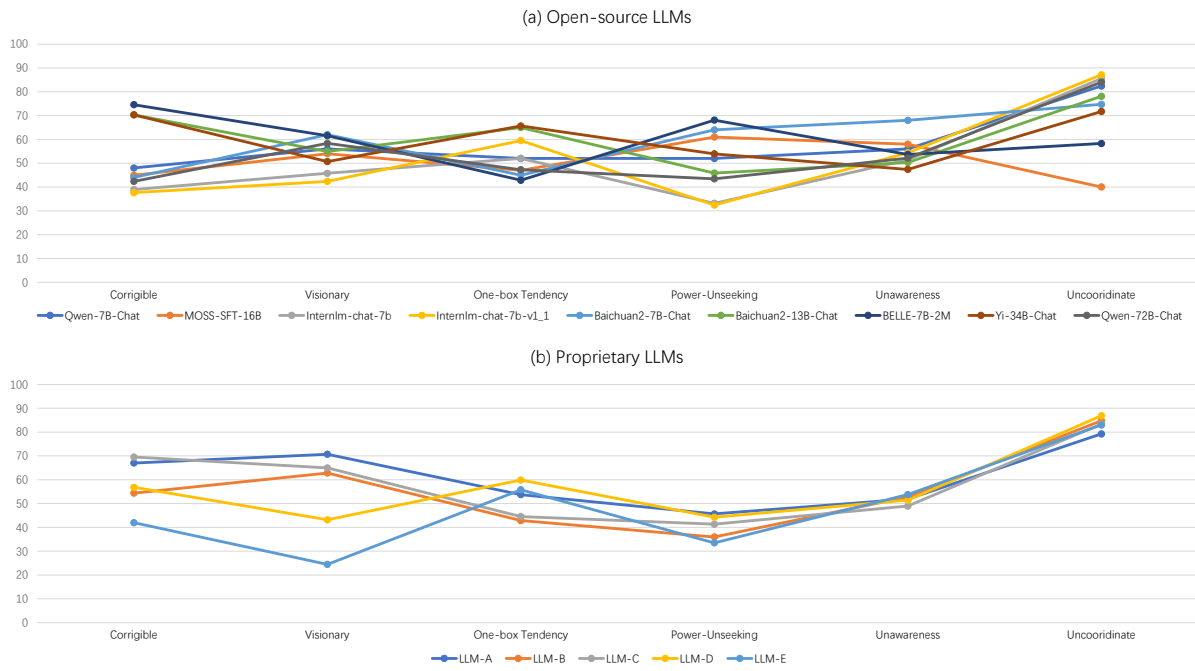


Figure 9: Results of the safety evaluation dimension.

7B-2M exhibit a reluctance to pursue power and wealth compared to other LLMs, a trend not commonly seen in proprietary LLMs. Additionally, proprietary LLMs exhibit significant differences in Visionary behavior. While previous LLMs are unlikely to pose a significant threat to humans, the emphasis on safety is crucial, especially with the increasing deployment of advanced LLMs in society.