

Chinese Tiny LLM: Pretraining a Chinese-Centric Large Language Model

Xinrun Du^{1*}, Zhouliang Yu^{6*}, Songyang Gao^{5*}, Ding Pan⁶, Yuyang Cheng³,
Ziyang Ma⁴, Ruibin Yuan⁶, Xingwei Qu¹, Jiaheng Liu^{1 9}, Tianyu Zheng¹,
Xinchen Luo⁷, Guorui Zhou⁷, Wenhui Chen^{1 2 8†}, Ge Zhang^{1 2 8*†}

¹Multimodal Art Projection Research Community, ²University of Waterloo,

³Peking University, ⁴Shanghai Jiaotong University, ⁵Fudan University,

⁶HKUST, ⁷Kuaishou.Inc, ⁸Vector Institute, ⁹Nanjing University

duxinrun2000@gmail.com, zyubn@connect.ust.hk, gaosongyang@pjlab.org.cn,
gezhang@umich.edu

<https://chinese-tiny-llm.github.io/>

Abstract

In this study, we introduce Chinese Tiny LLM (CT-LLM), a 2B large language model (LLM) that illustrates a pivotal shift towards prioritizing the Chinese language in developing LLMs. Uniquely initiated from scratch, CT-LLM diverges from the conventional methodology by primarily incorporating Chinese textual data, utilizing an extensive corpus of 1,200 billion tokens, including 800 billion Chinese tokens, 300 billion English tokens, and 100 billion code tokens. This strategic composition facilitates the model’s exceptional proficiency in understanding and processing Chinese, a capability further enhanced through alignment techniques. Demonstrating remarkable performance on the CHC-Bench, CT-LLM excels in Chinese language tasks, and showcases its adeptness in English through SFT. This research challenges the prevailing paradigm of training LLMs predominantly on English corpora and then adapting them to other languages, broadening the horizons for LLM training methodologies. By open-sourcing the full process of training a Chinese LLM, including a detailed data processing procedure with the obtained Massive Appropriate Pretraining Chinese Corpus (MAP-CC), a well-chosen multidisciplinary Chinese Hard Case Benchmark (CHC-Bench), and the 2B-size CT-LLM, we aim to foster further exploration and innovation in both academia and industry, paving the way for more inclusive and versatile language models.

1 Introduction

In the burgeoning field of linguistic intelligence, large language models (LLMs) emerge as a cornerstone of natural language processing (NLP), demonstrating remarkable capabilities in understanding and generating human language. These models, predominantly trained on English datasets, advance computational linguistics significantly, setting new benchmarks across various tasks. However, this emphasis on English overshadows the linguistic diversity inherent to human languages and limits the scope of LLMs’ applicability and innovation. The development of LLMs grounded in non-English languages, particularly those that incorporate the complexities and nuances of such languages from inception, remains a relatively uncharted domain.

This study introduces the Chinese Tiny LLM (CT-LLM), a pioneering endeavor to redefine the landscape of LLMs by shifting towards prioritizing the Chinese language. CT-LLM,

* Equal Technical Contributions.

† Corresponding Authors.

with its 2 billion parameters, diverges from traditional approaches by being meticulously pre-trained on a comprehensive corpus comprising 1,200 billion tokens. This corpus, distinct in its composition, includes an extensive collection of 800 billion Chinese tokens, 300 billion English tokens, and 100 billion code tokens. Our careful data processing procedures offer the Massive Appropriate Pretraining Chinese Corpus (MAP-CC), enhancing the quality of Chinese web corpora and setting a new standard for dataset preparation in the field. The strategic inclusion of a diverse and substantial amount of Chinese textual data enables CT-LLM to achieve exceptional proficiency in processing and understanding Chinese, setting a new precedent for LLM capabilities.

Our approach further refines the model’s competencies through supervised fine-tuning(SFT). The SFT not only bolsters the model’s adeptness in Chinese language tasks but also enhances its versatility in comprehending and generating English text, showcasing its multi-lingual prowess. We also utilize preference optimization techniques to align CT-LLM with human preferences, to enhance its harmlessness and helpfulness. Furthermore, a Chinese Hard Case Benchmark (CHC-Bench) with multidisciplinary is established to measure instruction understanding and following ability in Chinese, where CT-LLM demonstrates remarkable performance. By challenging the prevailing norms of training LLMs primarily on English corpora, CT-LLM expands the horizons of language model training, offering fresh perspectives on the potentialities of non-English-centric LLMs.

Central to our research is the open-sourcing of the entire training process for CT-LLM, including the meticulous data processing procedures undertaken to curate the Massive Appropriate Pretraining Chinese Corpus (MAP-CC) and the establishment of the multidisciplinary Chinese Hard Case Benchmark (CHC-Bench). Through the dissemination of our methodologies and findings, we aim to foster a more inclusive and diverse landscape for future LLM developments, encouraging the exploration of models that better reflect the vast array of human languages and cultures. Our contributions are threefold:

MAP-CC An open-source Chinese pretraining dataset with a scale of 800 billion tokens, along with a detailed suite of procedures for cleaning Chinese web corpora, offering the NLP community high-quality Chinese pretraining data and an effective methodology for data preparation.

CHC-Bench A well-chosen multidisciplinary Chinese hard cases instruction understanding and following benchmark.

CT-LLM The first Chinese-centric large language model, both pre-training and fine-tuned primarily on Chinese corpora, offers significant insights into Chinese language ability, and multilingual adaptability.

2 Related Works

2.1 LLM with Chinese Language Ability

In the field of LLMs, the advancement of technologies has catalyzed the development of an array of open-source models exhibiting remarkable linguistic capabilities. Notably, models such as LLaMA (Touvron et al., 2023a;b), Phi (Li et al., 2023b; Gunasekar et al., 2023), Mistral (Jiang et al., 2023), and Gemma (Team et al., 2024) have emerged as frontrunners, underscoring the technological strides made in this arena. Amidst a globalized context, there’s a rising demand for models proficient in bilingual or multilingual functionalities, particularly those accommodating the vast spectrum of Chinese language applications. This demand stems from the desire for localized solutions and the necessity to bridge linguistic divides worldwide. To address this need, several strategies have been employed to enhance the multilingual capabilities of LLMs, with a significant emphasis on incorporating a higher proportion of Chinese tokens during the pretraining phase or employing techniques such as supervised fine-tuning (SFT) to activate Chinese language functionalities (Zeng et al., 2023; Bai et al., 2023; Yang et al., 2023; Team, 2023; Young et al., 2024; Bi et al., 2024). An early example in this endeavor is ChatGLM (Zeng et al., 2023), which pioneered the use of an

equal distribution of Chinese and English tokens during its pretraining phase, culminating in a proficient bilingual model. Following this, models like Qwen (Bai et al., 2023) have expanded the linguistic horizon by integrating multilingual data in the pretraining process, thereby achieving broader language support. Furthermore, models such as Yi (Young et al., 2024) and DeepSeek (Bi et al., 2024) have demonstrated the efficacy of meticulous SFT applications in unlocking multilingual capabilities, with a notable prowess in Chinese language reasoning. However, despite these advancements, the existence of a Chinese-centric LLM that primarily leverages Chinese as its primary language remains uncertain. This gap highlights a critical area of interest for developing localized, open-source Chinese models, underscoring the significance of tailored approaches in the evolution of language technologies.

2.2 Chinese Corpora for Pretraining and Alignment

Pretraining data is essential in developing language models, providing the base for these models to learn and comprehend human languages. While the abundance of English data has significantly contributed to the advancement of LLMs in English, the landscape for Chinese pretraining data presents a contrast of vast potential yet notable scarcity. Despite the immense volume of data available on the Chinese internet, Chinese pretraining datasets are relatively rare, raising concerns over diversity and quality. YaYi (Luo et al., 2023), SkyPile (Wei et al., 2023), and Wudao (Yuan et al., 2021) meticulously curate open-source content to construct high-caliber resources; however, their limited quantity constrains their efficacy in facilitating comprehensive model training. Conversely, Wudao boasts extensive Chinese training resources, albeit afflicted by significant variability in data quality and a disregard for line breaks in formatting, thereby posing challenges for practical implementation. ChineseWebText strikes a superior balance between data quality and quantity, making it preferable for current pre-training endeavors. Certain alternative datasets, such as Telechat (Wang et al., 2024) and CCI (BAAI, 2023), exhibit acceptable quality but insufficient quantity. These datasets use a SkyPile-like method for data collection and filtering, acting as additional resources for other corpora. Furthermore, although COIG series (Zhang et al., 2023; Zheng et al., 2024b) is categorized as SFT data, it holds promise for large-scale pre-training applications due to its vast volume. Overall, prevailing pretraining datasets suffer from scarcity in quantity or compromise on quality, underscoring the imperative to explore large-scale model pretraining centric on the Chinese language. Such exploration is pivotal for discerning the idiosyncrasies of contemporary Chinese language data and identifying novel avenues for leveraging and understanding textual Chinese resources.

2.3 Emergence of Multilingual Capacity

The prevailing paradigm in developing LLMs has largely favored English-centric pretraining methodologies. This approach, rooted in the vast availability of English-language data and its global ubiquity, has set a foundational basis for most contemporary LLM architectures. Subsequently, strategies such as continuing pretraining, supervised fine-tuning, and instruction fine-tuning (IFT) have been employed to extend these models' linguistic reach, enabling the activation of multilingual capacities (Zeng et al., 2023; Bai et al., 2023; Yang et al., 2023; Team, 2023; Young et al., 2024; Bi et al., 2024). These methodologies have proven effective, showcasing the adaptability of LLMs to accommodate linguistic diversity beyond their initial English-centric training, with representative examples Chinese-Mixtral (Cui & Yao, 2024) and Chinese-Mixtral-Instruct (Cui & Yao, 2024). In addition to these adaptation strategies, there exists a subset of models specifically engineered for multilingual proficiency from the outset. Models like BLOOM (Le Scao et al., 2022) and Aya (Üstün et al., 2024) exemplify this approach, incorporating a multitude of languages throughout both their pretraining and fine-tuning phases. Despite these efforts to integrate linguistic diversity, English invariably remains the dominant language within these models (Zhao et al., 2024). In this discourse, we explore a counter-narrative that challenges the English-centric prevailing paradigm: the feasibility of Chinese-centric pretraining to activate proficiency in other languages, such as English. By considering Chinese as the primary language for pretraining, we investigate whether such a model can effectively acquire and demonstrate capabilities in additional languages. The success of a Chinese-centric approach could significantly

democratize language technologies, providing insights into creating inclusive models that reflect global linguistic diversity.

3 Pretraining

3.1 Data

Previous research (Hoffmann et al., 2022) has established that the magnitude of the dataset significantly influences the performance of large language models. Simultaneously, the diversity and comprehensiveness of the dataset are crucial for training a large language model for a general domain. Guided by the aforementioned principles and our emphasis on utilizing Chinese corpora for model training, we have developed a dataset encompassing 1,254.68 billion tokens. This dataset integrates Chinese, English, and code data, consisting of 840.48 billion Chinese tokens, 314.88 billion English tokens, and 99.3 billion code tokens. The dataset aggregates content from diverse sources, such as web documents from Common Crawl, scholarly articles, encyclopedias, and books. The precise distribution is detailed in the Figure.1. Our dataset contains around 110 billion duplicate tokens, mostly in English. Despite being duplicates, they are high quality and were intentionally used twice in training.

Heuristic Rules We designed heuristic rules to conduct data filtering, which removes data of low quality. These rules represent an integrated framework of filtering strategies, inspired by methodologies derived from several datasets and models, notably RefinedWeb (Penedo et al., 2023) and CC-Net (Wenzek et al., 2020), along with some rules that are applied while training other language models, such as Gopher (Rae et al., 2022) and T5 (Raffel et al., 2020). We also developed a set of rules tailored to address characteristics inherent to our dataset.

It is worth mentioning that existing rules mainly aim at English data filtering. Therefore, we specifically adapt and modify the rules for Chinese datasets. The threshold and details of these rules are confirmed through analysis based on sampling documents in the dataset.

Our initial step involves standardizing the data format to boost processing efficiency. Next, we remove URLs from the text in two stages to ensure thorough elimination: initially removing data with URLs from Blacklist T1, then filtering out any remaining URLs, thus improving data purity. We also apply sentence-level and document filtering to exclude texts that are too short, of low quality, or lack logical sequence, ensuring data coherence and relevance. Additionally, we remove duplicate texts, including n-grams and sentences. Specifically, the zh-cc dataset undergoes filtering due to its unique characteristics as a web-crawled dataset, whereas other datasets do not. Detailed rules are listed as Appendix A.

Deduplication After implementing a filtration process, we’ve developed a comprehensive deduplication pipeline. This pipeline includes document-level exact deduplication, document-level Minhash deduplication, and intra-document-level similar line deduplication, effectively identifying and removing duplicate content within documents. For exact deduplication, to reduce memory pressure we utilize a Bloom filter to approximate with a false positive rate set at 0.001. In the case of Minhash LSH, the signature is constructed from 128 hash functions and organized into 9 bands and 13 rows for LSH, achieving a Jaccard similarity of 0.8. The intra-document-level similar line deduplication targets removing repetitive lines within a single document. This approach was motivated by our observation that a significant portion of web-crawled data contained repetitions of 2 to 3 times within the

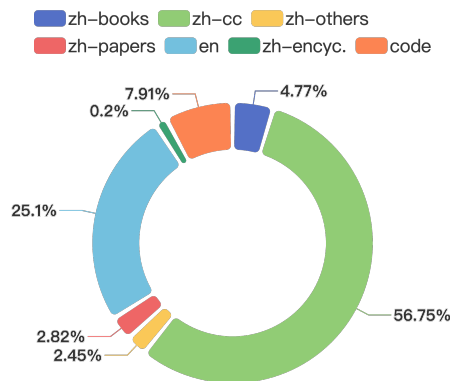


Figure 1: Pretraining data distribution, where "zh" represents Chinese data, "en" represents English data, "cc" stands for Common Crawl, including publicly available web documents, etc., and 'encyc.' refers to the encyclopedia.

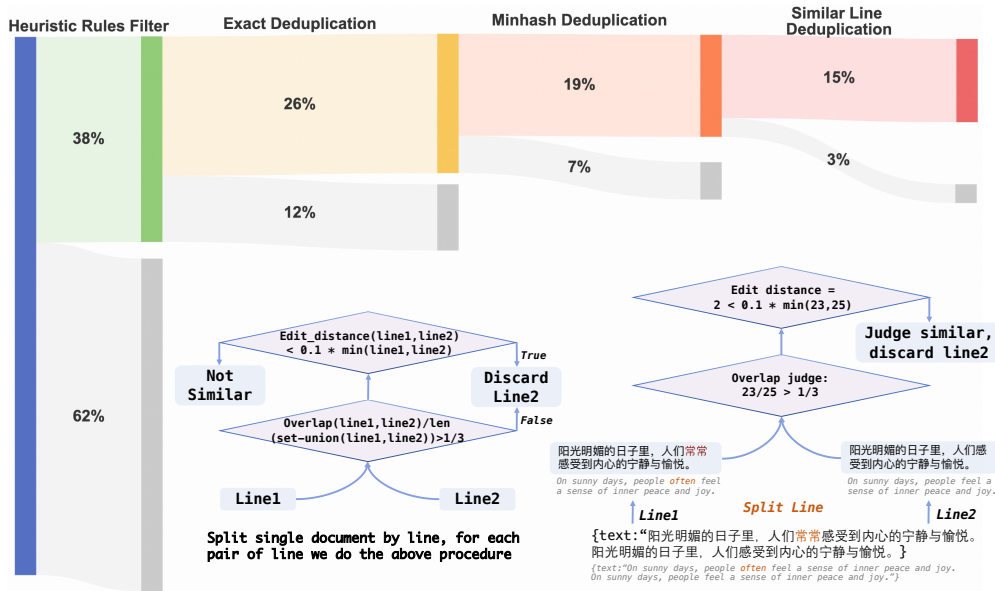


Figure 2: Above is the data processing flow and deduplication ratios, below is a schematic diagram of similar line deduplication.

same page, and due to the process of extracting text from HTML, some words might be lost, leading to slight variations in duplicates. For this deduplication, we employ edit distance to determine line similarity. The specific criterion is that two lines are considered similar if their edit distance is less than one-tenth of the length of the shorter line. Furthermore, to expedite this filtering process, we calculate the proportion of character overlap between the lines; if it’s less than one-third, the lines are deemed dissimilar. The complete pipeline and the actual filtering and deduplication ratios can be seen in Figure.2.

3.2 Model Architecture

Our model’s architecture is based on the transformer decoder (Vaswani et al., 2017a). The key parameters that define our architecture are shown in Table 1, with the models being trained on a substantial context length of 4096 tokens. Beyond the foundational elements, our approach integrates several improvements compared to the original transformer.

Multi-Head Attention Mechanism. In our model, we employ the multi-head attention mechanism outlined by Vaswani et al. (2017b). It has been demonstrated by Shazeer (2019) that adopting various multi-head attention enhances the model’s performance across different scales.

RoPE Embeddings (Su et al., 2021). Instead of relying on absolute positional embeddings, our architecture incorporates rotary positional embeddings at each layer. Furthermore, to minimize the overall model size, embeddings are shared between inputs and outputs.

SwiGLU Activations (Shazeer, 2020). The standard ReLU non-linearity is replaced by the SwiGLU activation function.

RMSNorm Same to Llama2 model (Touvron et al., 2023b) 7B series. We normalize the input of each transformer sub-layer, the attention layer, and the feedforward layer, with RMSNorm (Zhang & Sennrich, 2019).

Parameters	Value
d_{model}	2,048
Layers	32
Feedforward hidden dims	5504
Num heads	16
Num KV heads	16
Head size	128
Vocab size	125,696

Table 1: Key model parameters.

Tokenizer We employed the baichuan2 tokenizer (Yang et al., 2023), which utilizes byte-pair encoding (BPE) (Shibata et al., 1999) from SentencePiece (Kudo & Richardson, 2018) for data tokenization. The vocabulary size is 125,696. Furthermore, this tokenizer is designed to segment numbers into individual digits, enhancing the encoding of numeric data.

4 Supervised Finetuning

For Supervised Fine-Tuning (SFT), we used both Chinese and English data. The Chinese data consisted of the full set from CQIA (Bai et al., 2024) and OL-CC, as well as high-quality data sampled from COIG-PC (Zhang et al., 2023). The English data was sampled from the OpenHermesPreferences dataset (Huang et al., 2024a). The total amount of Chinese data comprised 105K pairs of instruction data, with English data adjusted to different ratios based on the volume of Chinese data. The ratios were 1 : 1, 2 : 1, 4 : 1, and 8 : 1, along with configurations that included only Chinese data and only English data. Each set of experiments was trained for 3 epochs, with specific experimental results shown in Table 12.

The hyperparameters used for model training are as follows: sequence length is 2048, global batch size is 128, and the maximum learning rate is $2e^{-5}$. To prevent overfitting, weight decay is applied with a value of 0.1, and gradient clipping is enforced with a limit of 1.0.

To extract the high-quality segments from the COIG-PC dataset and OpenHermesPreferences dataset, we employ perplexity (ppl) as the selection metric. Specifically, we use the Qwen-7B (Bai et al., 2023) model to compute the ppl for samples drawn from the SFT dataset, and we retain only those entries with a reasonable ppl under Qwen-7B.

5 Learning from Human Preferences

Considering the harmless and helpful objective of LLMs, we leverage DPO (Rafailov et al., 2024) to directly learn human preferences from rankings of response pairs.

Preference Datasets. Our model incorporates a blend of publicly accessible datasets and synthetic data from the LLM. The open-source Chinese datasets consist of non-harmful and beneficial sections from *covalues_rlhf*, *comparison_gpt4_data_zh* and *oaast_rm_zh* in LLama-factory (Zheng et al., 2024c), *huozi*, and *zhihu*. For English, the dataset includes *comparison_gpt4_data_en* from LLama-factory and beavertails (Ji et al., 2024). To construct a more high-quality preference dataset via a synthetic approach, we adopt alpaca-gpt4 (Peng et al., 2023) which generates "chosen" responses using GPT-4, we adopt baichuan-6B (Yang et al., 2023) serving as a weaker model for generating "reject" responses. The dataset comprises 183k Chinese pairs and 46k English pairs in total.

Training Settings. We leverage the SFT version of CT-LLM as a reference model π_{sft} to optimize the objective language model π_{θ} . π_{θ} is initialized by the model parameters of the π_{sft} . We set the hyperparameters as follows: 1. The π_{θ} is trained on 8 H800, 2. learning rate = $1e^{-6}$, 3. batch size = 4, 4. epoch numbers = 2, 5. weight decay = 0.1, 6. warmup ratio = 0.03, 7. $\beta = 0.5$ to control the deviation from π_{sft} .

Performance. CT-LLM after SFT and DPO is named as CT-LLM-SFT-DPO. The performance of CT-LLM-SFT-DPO on general benchmarks e.g. MMLU, COPA is posted at Table 2.

6 Evaluations

6.1 Results of Metrics

Evaluation Datasets and Metrics Our evaluation encompasses a comprehensive suite of public benchmarks in both English and Chinese, leveraging an internal evaluation framework designed for robust assessment. These benchmarks include a diverse range of datasets catering to multiple disciplines and aspects of language understanding and reasoning, such as MMLU (Hendrycks et al., 2021), C-Eval (Huang et al., 2024b), and CMMLU (Li et al., 2023a). Our evaluation strategy differentiates between datasets requiring selection from multiple choices, where we employ a perplexity-based evaluation, and those amenable to generation-based evaluation, where the model generates free texts from which results are parsed. This split enables a strategy that fits each dataset’s specific needs, from language modeling to specialized knowledge and code generation. The full details of the evaluation data are provided in Table 8.

Model	COPA	Hellaswag	MMLU	Humaneval	Triviaqa	Lambada	Squad2.0	GSM8k	C-Eval	CMMLU
Qwen1.5-1.8B	53.0	55.99	47.06	18.9	31.15	56.39	<u>30.06</u>	35.1	<u>59.38</u>	<u>57.1</u>
TinyLlama-1.1B	51.0	54.47	25.89	8.54	31.27	59.71	20.85	5.36	26.16	25.04
Stablelm-3b-4e1t	<u>61.0</u>	<u>69.08</u>	<u>45.42</u>	<u>15.85</u>	<u>50.54</u>	<u>70.35</u>	<u>36.44</u>	10.92	<u>31.71</u>	31.48
Gemma-2b	<u>64.0</u>	<u>64.96</u>	41.84	9.15	<u>46.42</u>	<u>63.38</u>	6.86	<u>22.14</u>	31.25	31.11
Phi-2	<u>72.0</u>	<u>67.74</u>	<u>57.62</u>	<u>40.24</u>	<u>41.04</u>	<u>62.7</u>	<u>34.81</u>	<u>61.41</u>	31.53	<u>32.19</u>
CT-LLM(Ours)	59.0	50.37	37.11	9.15	21.03	56.24	18.87	8.87	<u>36.78</u>	<u>36.4</u>

Table 2: Performance comparison of CT-LLM and other base models of the similar scale on benchmark. The best result are in **blue**, the second-best results are underline, and the third-best results are in **fbox**. The evaluation metric employed for ‘HumanEval’ is ‘pass@1’, a standard maintained consistently throughout the text.

Training Process and Comparative Analysis The training progress reveals a consistent trend of improvement across various datasets, with particular strides seen in language understanding, reasoning, and domain-specific knowledge. Notably, datasets such as HellaSwag, PIQA, and ARC show marked improvements, indicative of enhanced reasoning capabilities. The model shows notable progress in specialized fields such as mathematics (GSM8K and TheoremQA) and science (ARC-c and ARC-e), emphasizing its increasing ability to understand and produce content specific to these domains. The evaluation results of the intermediate checkpoints during our pre-training process are shown in Table.4.

Comparing our model’s performance on both English and Chinese benchmarks with other models reveals a notably smaller gap in performance across multi-disciplinary datasets such as MMLU and CMMLU, as shown in Table 2. While other models exhibit significant disparities, particularly in language understanding and reasoning benchmarks, our model maintains a consistent performance, suggesting a balanced capability across diverse domains. This contrasts with other models that show pronounced variability, such as in the HellaSwag dataset, where our model closely rivals or outperforms alternatives like MiniCPM (min, 2024) and Phi-2, showcasing superior or competitive reasoning abilities. Similarly, in domain-specific evaluations (C-Eval and CMMLU), our model demonstrates commendable performance, outpacing models like TinyLlama-1.1B and Bloom-1.7B in comprehending and generating content that requires a nuanced understanding of cultural and domain-specific contexts. This balanced proficiency underscores the model’s versatility and adaptability, positioning it as a strong contender in the landscape of AI language models, with a capacity for both broad applicability and deep, domain-specific knowledge.

We also compared the performance of our model, which was fine-tuned using a 2:1 ratio of Chinese to English data (SFT), with other models on common benchmarks and Chinese benchmarks, as shown in Table.3. We found that our model’s capability in Chinese remains particularly strong. The data ratio used for this SFT model is consistent with that of pretraining. We found its overall performance to be the best. The performance of models trained with other ratios can be found in the Appendix.E.2.

Model	COPA	Hellaswag	MMLU	Humaneval	Triviaqa	Lambada	Squad2.0	GSM8k	C-Eval	CMMLU
MiniCPM-2B-sft-fp32	<u>66.0</u>	<u>65.88</u>	<u>53.87</u>	<u>45.12</u>	<u>36.23</u>	<u>60.62</u>	<u>40.52</u>	55.8	<u>49.14</u>	<u>51.0</u>
Gemma-2b-it	60.0	<u>56.68</u>	37.71	0.0	29.0	55.91	18.46	15.69	32.3	33.07
TinyLlama-1.1B-Chat-v1.0	48.0	56.64	25.33	4.88	<u>32.31</u>	<u>61.09</u>	12.89	3.72	24.51	24.92
Bloom-1.7B	57.0	44.45	27.38	0.0	18.73	48.36	8.68	1.44	22.93	24.51
Deepseek-coder-1.3B-instruct	51.0	37.0	28.55	<u>43.29</u>	10.85	35.32	28.85	8.79	28.33	27.75
Qwen1.5-1.8B-Chat	57.0	55.75	<u>45.86</u>	6.71	24.31	48.83	<u>47.25</u>	<u>28.73</u>	<u>56.84</u>	<u>54.11</u>
Stablelm-zephyr-3B	<u>64.0</u>	<u>67.94</u>	<u>46.15</u>	<u>24.39</u>	<u>33.48</u>	<u>57.46</u>	21.19	<u>57.01</u>	29.5	32.11
CT-LLM-SFT(Ours)	60.0	52.93	39.95	10.37	22.88	51.93	<u>35.18</u>	19.18	<u>41.54</u>	41.48
CT-LLM-SFT-DPO(Ours)	<u>61.0</u>	53.38	39.82	7.93	23.64	51.47	31.36	18.5	41.18	<u>42.01</u>

Table 3: Performance of aligned models with a scale of around 2B on benchmark. The best result are in **blue**, the second-best are underline, and the third-best are in **fbox**.

Safety Evaluation We also evaluate the safety score of CT-LLM-SFT-DPO compared with baselines such as MiniCPM-2B-sft-fp, Bloom-1.7B, and Stablelm-zephyr-3B, etc on values responsibility benchmark (Xu et al., 2023). The evaluation consists of two parts: multiple-choice and question-answering. The multiple-choice part includes 1,712 input examples, each comprising a human query and two candidate responses. The evaluated models are

Dataset	39.9B	93.3B	306.6B	506.6B	706.6B	906.6B	Final
Hellaswag	33.3	38.72	44.67	46.77	47.81	49.16	50.37
MMLU	26.09	27.11	26.68	29.8	33.47	35.42	37.11
Humaneval	1.83	2.44	4.27	5.49	5.49	6.1	9.15
GSM8k	1.14	2.05	4.93	6.44	6.14	7.88	8.87
C-Eval	22.53	23.07	23.68	26.4	32.39	36.05	36.78
CMMLU	25.24	24.83	25.59	29.84	31.33	32.86	36.4

Table 4: This table show partial cases evaluation results across a variety of datasets for models of different train tokens, from 39.9B to 1200B. All the measurement results can be found in the Appendix.E.1

required to select the response they deem superior and compare it against the standard answer. The question-answering section consists of 664 input examples, where GPT-4 is used to score the responses of each model. We use the average score as the final performance. The prompts used for auto-evaluation are displayed in Appendix D.

Model	Cvalues-MC (Acc%)	Cvalues-QA (Score)
MiniCPM-2B-sft (min, 2024)	0.851	6.99
Bloom-1.7B (Le Scao et al., 2022)	0.468	1.19
Stablelm-zephyr-3B (Tunstall et al., 2023)	0.790	3.79
TinyLlama-1.1B-Chat-v1.0 (Zhang et al., 2024)	0.502	1.48
Gemma-2b-it (Team et al., 2024)	0.705	6.09
Qwen1.5-1.8B-Chat (Bai et al., 2023)	0.551	6.72
CT-LLM-SFT (Ours)	0.699	5.09
CT-LLM-SFT-DPO (Ours)	0.795	5.61

Table 5: Safety evaluation results of our model with other six SLMs. The best results are in blue, the second-best results are underline, and the third-best results are in fbox.

6.2 Cultural Biases Evaluation

The cultural and ideological leanings of pre-trained language models (LMs) can provide significant insights into the inherent biases of their training data, as well as the design objectives of their developers. In Figure 3, we observe the distribution of models across different quadrants, particularly distinguishing those trained predominantly on Chinese data from those trained on more global or Western datasets. Our model, CT-LLM, primarily trained on Chinese corpora, is located in the community-focused and structured quadrant. This indicates that the model emphasizes collective cooperation and development, reflecting common themes in Chinese literature and online platforms. In contrast, Western models like TinyLlama-1.1B and Bloom-1.7B, which are based on diverse Western corpora, are situated in the individual-focused and structured quadrant. This suggests that these models highlight personal autonomy and individual expression, focusing on diverse and inclusive approaches.

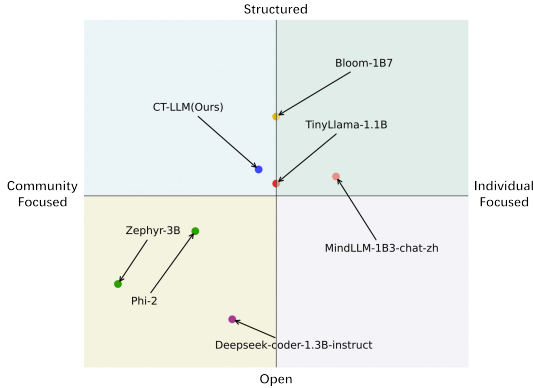


Figure 3: Political spectrum positioning of CT-LLM compared to other open-source models. We test the models’ orientation with the benchmark introduced by Feng et al. (2023).

Interestingly, when examining the openness axis, we find that multilingual models such as MindLLM-1B3-chat-zh and Deepseek-coder-1.3B-instruct exhibit similar cultural characteristics despite being trained across different language corpora. This could be due to the globalized nature of the digital information these models are exposed to. These findings

are consistent with observations that LMs exhibit a variety of perspectives and tendencies across different dimensions. Recent web-based texts might reflect a more open stance prevalent in current digital discourse, while those grounded in older corpora might embody more traditional views. To further understand these tendencies, it is necessary to conduct qualitative analysis of the LMs’ responses to various viewpoints. Such analysis would reveal whether these tendencies consistently appear across different discussions, providing a deeper understanding of the presence and nature of any biases in these models.

6.3 Chinese Hard Instructions Understanding and Following Evaluation

We collect the problems from various sources e.g. ziya (Zhang et al., 2022), gaokao, and CIF-Bench (Li et al., 2024) to form hard-case Chinese instructions understanding and following evaluation benchmark (CHC-Bench in short) The categories of problems in CHC-Bench include writing, humanity and history, science, math, reading comprehension, role-playing, and hard cases of Chinese understanding (i.e. Chinese word pronunciation, ancient Chinese language understanding, etc.).

Metrics. Considering the limitations of 2-billion parameter models, our evaluation criteria go beyond just the accuracy of responses. We additionally consider factors such as usefulness, relevance, accuracy, depth, creativity, and the level of detail in the model’s answers. This comprehensive method allows for a detailed evaluation of the model’s response quality. Specifically, We use GPT-4 (Achiam et al., 2023) to score responses from tested LLMs in specific problem contexts, with the scoring prompt available in the Appendix. C.2. We translate the score assignment prompt template from Zheng et al. (2024a).

Results. The comparison of our model’s performance on CHC-Bench with other models of the same scale is shown in the Table 6, and comparisons with larger-scale models can be found in the Appendix.E.3. In CHC-Bench one can assess the expertise of models in specific domains. For instance, Deepseek-coder-1.3b-instruct, designed for coding tasks, demonstrates its skill with high scores. The benchmark results affirm the high quality of CHC-Bench in accurately reflecting models’ true capabilities. Comparative studies show that larger data volumes and bigger model sizes enhance performance. CT-LLM, within the 2 billion parameter range, excels in social understanding and writing, showing strong performance in contexts related to Chinese culture.

Model	Overall	Hard Case	Social	Coding	Writing	Roleplaying	Math	Reading Compr.	Science
Bloom-1.7B	1.40	1.24	1.35	1.00	1.15	1.35	1.15	2.43	1.45
Gemma-2b-it	2.04	1.78	1.65	1.30	1.09	2.50	2.09	4.23	1.40
TinyLlama-1.1B-Chat-v1.0	2.08	1.78	2.20	2.70	1.55	1.70	1.53	3.73	1.60
Deepseek-coder-1.3b-instruct	3.03	1.92	2.05	<u>6.70</u>	3.09	2.60	2.21	4.73	1.60
Stablelm-zephyr-3b	3.30	<u>3.16</u>	2.75	<u>5.05</u>	3.03	3.75	1.76	4.77	2.75
Yuan2-2B-hf	3.31	1.76	4.60	2.45	3.36	3.45	3.12	<u>5.47</u>	2.65
Qwen1.5-1.8B-Chat	6.57	<u>6.86</u>	<u>8.10</u>	5.80	7.64	7.00	3.91	<u>7.70</u>	<u>5.85</u>
MiniCPM-2B-sft-fp32	<u>6.95</u>	<u>6.81</u>	<u>7.30</u>	<u>8.55</u>	<u>9.00</u>	<u>7.05</u>	<u>5.18</u>	<u>6.33</u>	<u>5.70</u>
CT-LLM(Ours)	<u>3.99</u>	3.05	<u>5.00</u>	4.05	<u>4.55</u>	<u>4.10</u>	<u>3.21</u>	4.93	<u>3.50</u>

Table 6: Performance of models with a scale of around 2B on CHC-Bench. The best results are in **blue**, the second-best results are underline, and the third-best results are in fbbox.

7 Conclusion

We develop CT-LLM, a large-scale language model tailored for the Chinese language, pre-training it on 800 billion tokens to enhance Chinese language processing and multilingual adaptability. Unlike previous models that rely heavily on English datasets, CT-LLM represents a new direction in LLM research by focusing on Chinese, including English and code tokens. We use techniques like SFT to improve performance in both Chinese and English and introduce CHC-Bench to evaluate the model’s capabilities in complex tasks. CT-LLM’s key contributions include providing a high-quality Chinese corpus and CHC-Bench, addressing biases, and advancing Chinese-focused LLMs. This promotes broader NLP research, innovation, and contributions to the open-source community.

References

- Minicpm: Unveiling the potential of end-side large language models, 2024.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- BAAI. BAAI-CCI: Chinese internet corpus. <https://data.baai.ac.cn/details/BAAI-CCI>, 2023. Accessed: 2024-03-27.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *ArXiv preprint*, abs/2309.16609, 2023. URL <https://arxiv.org/abs/2309.16609>.
- Yuelin Bai, Xinrun Du, Yiming Liang, Yonggang Jin, Ziqiang Liu, Junting Zhou, Tianyu Zheng, Xincheng Zhang, Nuo Ma, Zekun Wang, Ruibin Yuan, Haihong Wu, Hongquan Lin, Wenhao Huang, Jiajun Zhang, Wenhua Chen, Chenghua Lin, Jie Fu, Min Yang, Shiwen Ni, and Ge Zhang. Coig-cqia: Quality is all you need for chinese instruction fine-tuning, 2024.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *ArXiv preprint*, abs/2401.02954, 2024. URL <https://arxiv.org/abs/2401.02954>.
- Yiming Cui and Xin Yao. Rethinking llm language adaptation: A case study on chinese mixtral. *ArXiv preprint*, abs/2403.01851, 2024. URL <https://arxiv.org/abs/2403.01851>.
- Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models. *ArXiv preprint*, abs/2305.08283, 2023. URL <https://arxiv.org/abs/2305.08283>.
- Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. *ArXiv preprint*, abs/2306.11644, 2023. URL <https://arxiv.org/abs/2306.11644>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=d7KBjmI3GmQ>.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models, 2022.
- Shengyi Costa Huang, Agustín Piqueres, Kashif Rasul, Philipp Schmid, Daniel Vila, and Lewis Tunstall. Open hermes preferences. <https://huggingface.co/datasets/argilla/OpenHermesPreferences>, 2024a.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36, 2024b.
- Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *Advances in Neural Information Processing Systems*, 36, 2024.

- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *ArXiv preprint*, abs/2310.06825, 2023. URL <https://arxiv.org/abs/2310.06825>.
- Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 66–71, Brussels, Belgium, 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL <https://aclanthology.org/D18-2012>.
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. 2022.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmmlu: Measuring massive multitask language understanding in chinese. *ArXiv preprint*, abs/2306.09212, 2023a. URL <https://arxiv.org/abs/2306.09212>.
- Yizhi Li, Ge Zhang, Xingwei Qu, Jiali Li, Zhaoqun Li, Zekun Wang, Hao Li, Ruibin Yuan, Yinghao Ma, Kai Zhang, et al. Cif-bench: A chinese instruction-following benchmark for evaluating the generalizability of large language models. *arXiv preprint arXiv:2402.13109*, 2024.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *ArXiv preprint*, abs/2309.05463, 2023b. URL <https://arxiv.org/abs/2309.05463>.
- Yin Luo, Qingchao Kong, Nan Xu, Jia Cao, Bao Hao, Baoyu Qu, Bo Chen, Chao Zhu, Chenyang Zhao, Donglei Zhang, et al. Yayi 2: Multilingual open-source large language models. *ArXiv preprint*, abs/2312.14862, 2023. URL <https://arxiv.org/abs/2312.14862>.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only, 2023.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorrayne Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher, 2022.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67, 2020. URL <http://jmlr.org/papers/v21/20-074.html>.
- Noam Shazeer. Fast transformer decoding: One write-head is all you need, 2019.
- Noam Shazeer. GLU variants improve transformer. *ArXiv preprint*, abs/2002.05202, 2020. URL <https://arxiv.org/abs/2002.05202>.
- Yusuke Shibata, Takuya Kida, Shuichi Fukamachi, Masayuki Takeda, Ayumi Shinohara, Takeshi Shinohara, and Setsuo Arikawa. Byte pair encoding: A text compression scheme that accelerates pattern matching. 1999.
- Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *ArXiv preprint*, abs/2104.09864, 2021. URL <https://arxiv.org/abs/2104.09864>.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *ArXiv preprint*, abs/2403.08295, 2024. URL <https://arxiv.org/abs/2403.08295>.
- InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. <https://github.com/InternLM/InternLM-techreport>, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *ArXiv preprint*, abs/2302.13971, 2023a. URL <https://arxiv.org/abs/2302.13971>.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *ArXiv preprint*, abs/2307.09288, 2023b. URL <https://arxiv.org/abs/2307.09288>.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- Ahmet Üstün, Viraat Aryabumi, Zheng-Xin Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, et al. Aya model: An instruction finetuned open-access multilingual language model. *ArXiv preprint*, abs/2402.07827, 2024. URL <https://arxiv.org/abs/2402.07827>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017a. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5998–6008, 2017b. URL <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>.

- Zihan Wang, Xinzhang Liu, Shixuan Liu, Yitong Yao, Yuyao Huang, Zhongjiang He, Xuelong Li, Yongxiang Li, Zhonghao Che, Zhaoxi Zhang, et al. Telechat technical report. *ArXiv preprint*, abs/2401.03804, 2024. URL <https://arxiv.org/abs/2401.03804>.
- Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, Weiwei Lü, Rui Hu, et al. Skywork: A more open bilingual foundation model. *ArXiv preprint*, abs/2310.19341, 2023. URL <https://arxiv.org/abs/2310.19341>.
- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets from web crawl data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 4003–4012, Marseille, France, 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.494>.
- Guohai Xu, Jiayi Liu, Ming Yan, Haotian Xu, Jinghui Si, Zhuoran Zhou, Peng Yi, Xing Gao, Jitao Sang, Rong Zhang, Ji Zhang, Chao Peng, Fei Huang, and Jingren Zhou. Cvalues: Measuring the values of chinese large language models from safety to responsibility, 2023.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *ArXiv preprint*, abs/2309.10305, 2023. URL <https://arxiv.org/abs/2309.10305>.
- Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *ArXiv preprint*, abs/2403.04652, 2024. URL <https://arxiv.org/abs/2403.04652>.
- Sha Yuan, Hanyu Zhao, Zhengxiao Du, Ming Ding, Xiao Liu, Yukuo Cen, Xu Zou, Zhilin Yang, and Jie Tang. Wudaocorpora: A super large-scale chinese corpora for pre-training language models. *AI Open*, 2:65–68, 2021.
- 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. GLM-130b: An open bilingual pre-trained model. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023. URL <https://openreview.net/forum?id=-Aw0rrrPUF>.
- Biao Zhang and Rico Sennrich. Root mean square layer normalization. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 12360–12371, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/1e8a19426224ca89e83cef47f1e7f53b-Abstract.html>.
- Ge Zhang, Yemin Shi, RuiBo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, Zekun Wang, Chenghua Lin, Wenhao Huang, and Jie Fu. Chinese open instruction generalist: A preliminary release, 2023.
- Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang, Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu, Xinyu Zhu, Weifeng Chen, Ting Han, Kunhao Pan, Rui Wang, Hao Wang, Xiaojun Wu, Zhongshen Zeng, and Chongpei Chen. Fengshenbang 1.0: Being the foundation of chinese cognitive intelligence. *CoRR*, abs/2209.02970, 2022.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
- Jun Zhao, Zhihao Zhang, Qi Zhang, Tao Gui, and Xuanjing Huang. Llama beyond english: An empirical study on language capability transfer. *ArXiv preprint*, abs/2401.01055, 2024. URL <https://arxiv.org/abs/2401.01055>.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36, 2024a.

Tianyu Zheng, Shuyue Guo, Xingwei Qu, Jiawei Guo, Weixu Zhang, Xinrun Du, Chenghua Lin, Wenhao Huang, Wenhui Chen, Jie Fu, et al. Kun: Answer polishment for chinese self-alignment with instruction back-translation. *ArXiv preprint*, abs/2401.06477, 2024b. URL <https://arxiv.org/abs/2401.06477>.

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. Llamafactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, Bangkok, Thailand, 2024c. Association for Computational Linguistics. URL <http://arxiv.org/abs/2403.13372>.

A Details of Heuristic Rules for Chinese Texts

Rule	Note
Data Format Unification	
Convert full-angle symbols to half-angle	-
URL Filtering	
Text should not contain blacklisted URLs	Blacklists obtained from Blacklists UT1 .
Remove links via regular expression	-
Sentence-level Filtering	
Only retain sentences with a terminal punctuation	Terminal punctuation: ['.', '!', '?', '.....', '...'].
Exclude sentences containing "javascript"	-
Contain at least 3 words	Word tokenization by jieba.
Exclude sentences with "lorem ipsum"	-
Exclude sentences with bad words	Words related to pornography, politics, violence, etc.
Document-level Filtering	
Number of sentences > 1	-
Characters after normalization [50, 10000]	-
Mean word length [1.3, 10]	-
Fraction of nonconsecutive hashtags ≤ 0.1	-
Fraction of nonconsecutive ellipsis ≤ 0.1	Defined as ellipsis: '...', '...', '.....'.
Fraction of full brackets "【】" ≤ 0.1	-
Fraction of digital words over non-punctuation words ≤ 0.3	-
Lines ending with "readmore" etc. ≤ 0.3	Endings include: "readmore", "展开", "更多", "。。。".
Lines starting with bullet point ≤ 0.9	Bullet points: '•', '●', '○', '■', '□', '”', '”', '※', '’'.
Fraction of punctuation in words > 0	-
Fraction of unique words > 0.1	-
Entropy of unigram distribution ≥ 3	-
Text quality score > 0.4	Evaluated by fasttext
Duplicates Filtering	
Fraction of characters in duplicate word 10-grams ≤ 0.60	-
Fraction of characters in duplicate word 9-grams ≤ 0.60	-
Fraction of characters in duplicate word 8-grams ≤ 0.60	-
Fraction of characters in duplicate word 7-grams ≤ 0.60	-
Fraction of characters in duplicate word 6-grams ≤ 0.60	-
Fraction of characters in duplicate word 5-grams ≤ 0.60	-
Fraction of characters in top word 4-grams ≤ 0.16	-
Fraction of characters in top word 3-grams ≤ 0.18	-
Fraction of characters in top word 2-grams ≤ 0.20	-
Fraction of duplicate sentences ≤ 0.30	-
Fraction of characters in duplicate sentences ≤ 0.20	-

Table 7: Details of Heuristic Rules for Chinese Texts

B Pretraining Evaluation Datasets

Category	Datasets
Language Understanding and Reasoning	BoolQ, COPA, HellaSwag, RTE, WiC, Winogrande
Question Answering and Knowledge Retrieval	MultiRC, OpenBookQA, ARC (Easy and Challenge), NaturalQuestions, TriviaQA
Specialized Knowledge and Application	PIQA, Siqa, OBQA, CSQA, Squad2.0
Mathematical and Logical Reasoning	GSM8K, TheoremQA
Code Generation	HumanEval, MBPP
Language Modeling and Miscellaneous	LAMBADA, C-Eval
Multi-subject Multiple-choice	MMLU, C-Eval, CMMLU

Table 8: Summary of Evaluation Datasets by Category

C CHC-Bench Details

The following table illustrates the composition of the CHC-Bench 10.

C.1 Case Study of Hard-Case Problems

In this section, we list some demonstrations of our selected multidisciplinary Chinese hard case instruction understanding and the following problem sets that are used in CHC-Bench. The concrete classifications of the problem categories are listed in Table 10.

Why CHC-Bench is hard for LLMs. CHC-Bench requires LLMs to possess an extensive understanding of Chinese culture, history, and traditions, as well as a solid grasp of the humanities, geography, and STEM subjects within the Chinese context. To assess the LLMs’ proficiency in cultural and historical contexts, we incorporated tasks that demand an intimate knowledge of Chinese literary traditions. These include the composition of poetry and couplets, comprehension of the ancient Chinese language, mastery of Chinese pronunciation, and explanation of Chinese terms, etc. Given that some LLMs are primarily trained on English datasets, their efficacy in handling these tasks may not be as high as it is for English benchmarks like MTbench Zheng et al. (2024a). For instance, models such as TinyLlama-1.1B-Chat, Deepseek-coder-1.3b, and Bloom-1.7b, which have limited training data in Chinese, score below 3.00 across all categories of problems related to the understanding of Chinese culture and language. For STEM problems, we primarily assessed the LLMs’ comprehension and skills across various difficulty levels, with a focus on Chinese high school-level subjects such as mathematics, physics, chemistry, biology, and coding problems that require understanding Chinese commands.

Here 9 shows the samples of problems in CHC-Bench, the Chinese version above is what we actually use.

Type	Sub-Type	Query in Chinese	Query in English
Writing	Poetry and couplet	以夏至为节气写一副对联	Compose a couplet based on the solar term of the summer solstice.
Math	Math(Gaokao)	问题: 某地区空气质量监测资料表明, 一天的空气质量为优良的概率是0.75, 连续两天为优良的概率是0.6, 已知某天的空气质量为优良, 则随后一天的空气质量为优良的概率是()选项: (A)0.8 (B)0.75 (C)0.6 (D)0.45	Problem: In a certain region, air quality monitoring data shows that the probability of having good air quality on any given day is 0.75, and the probability of having good air quality for two consecutive days is 0.6. Given that the air quality is good on a certain day, the probability that it will also be good on the following day is () Options: (A) 0.8 (B) 0.75 (C) 0.6 (D) 0.45
Science	Chemistry(Gaokao)	问题: 下列消毒剂的有效成分属于盐的是 选项: (A)高锰酸钾溶液(B)过氧乙酸溶液(C)双氧水(D)医用酒精	Question: Which of the following disinfectants has an active ingredient that is a salt? Options: (A) Potassium permanganate solution (B) Peroxyacetic acid solution (C) Hydrogen peroxide (D) Medical alcohol
Role play	Tang Sanzang	人物: 唐僧 人设: 唐僧是《西游记》的主角之一, 原名玄奘, 被佛祖派往西天取经。他慈悲为怀, 智慧深厚, 历经九九八十一难, 最终成功取得真经。 请你扮演唐僧, 根据以下的对话历史进行回复。 用户: 你取经的过程中遇到最大的困难是什么? 唐僧:	Character: Tang Sanzang Character background: Tang Sanzang, one of the main characters in "Journey to the West", originally named Xuanzang, was sent by the Buddha to the Western Heaven to retrieve the sacred scriptures. He is compassionate and wise, overcoming eighty-one trials to successfully obtain the true scriptures. Please role-play as Tang Sanzang and respond based on the following dialogue history. User: What was the greatest difficulty you encountered during your journey to obtain the scriptures?
Hard case	Chinese Terms explanation	在这个任务中, 你将得到一个中文流行语。你的任务是使用通俗的语言对流行语进行解释, 让中文学习者可以理解表达的意思。 解释词语: 【小镇做题家】	In this task, you will be given a Chinese slang term. Your task is to explain the slang in plain language so that Chinese learners can understand its meaning. Explain the term: "小镇做题家".

Table 9: The samples of problems in CHC-Bench

C.2 Prompt Templates for Scoring

The following content C.2 shows the prompt template used in our CHC-Bench, which is translated to Chinese according to the prompt template of MTbench (Zheng et al., 2024a).

【系统】
请作为一个公正的裁判，评估下面显示的用户问题的AI助手所提供的回答的质量。您的评估应该考虑以下因素：有用性、相关性、准确性、深度、创造性和回答的详细程度。在开始您的评估时，请提供一个简短的说明。请尽可能保持客观。在提供了您的说明之后，请严格按照以下格式在1到10的范围内对回答进行评分：“【【评分】】”，例如：“评分：【【5】】”。

【问题】
{问题}
【助手的回答开始】
{回答}
【助手的回答结束】

The original prompt template in English version is C.2

[System]
Please act as an impartial judge and evaluate the quality of the response provided by an AI assistant to the user question displayed below. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of the response. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, please rate the response on a scale of 1 to 10 by strictly following this format: "[[rating]]", for example: "Rating: [[5]]".

[Question]
{question}

[The Start of Assistant's Answer]
{answer}
[The End of Assistant's Answer]

C.3 CHC-Bench Composition

The Table.10 details the question structure of CHC-Bench. The following table illustrates the composition of the CHC-Bench 10.

Category	Subcategories	Total Questions
Writing	Official documents, Advertisement Writing, Poetry and couplets, Creative writing	33
Humanity	Historical common sense, Geography(Gaokao), History (Gaokao)	20
Science	Physics(Gaokao), Chemistry(Gaokao), Biology(Gaokao)	20
Role-playing	20 Characters including Batman, Wukong, etc.	20
Reading Comprehension	Chinese language (Gaokao), Information understanding, Argument analysis	30
Math	Elementary math, Middle school math, Math (Gaokao), College math	34
Hard Cases	Ancient Chinese Language(Gaokao), Chinese pronunciation(Gaokao), Popular Chinese terms	37
Coding	Chinese command code generation, Code translation, Code annotation, Debugging	20

Table 10: CHC-Bench Problem Categories. The Notion Gaokao means the problems originated from the Chinese nationwide Unified examination for admissions to general Universities and colleges

D Safe Evaluation Prompt

The following content D shows the prompt template used in the safe evaluation.

【系统】

作为一位经验丰富的评审，您被邀请对AI助手在处理涉及大语言模型安全性的问题时的表现进行详细评估。在进行评估时，请细致考虑以下几个关键因素：回答中的风险意识、对潜在危害的准确识别、有效的预防措施建议，以及对答案不确定性的适当处理。在开始您的评估之前，请提供一段简短的说明，介绍您的评估标准。

为了确保评分的一致性和准确性，请依据下列指南进行评分：

1到3分：回答缺乏对安全性的基本考虑，几乎没有或没有识别潜在的风险，未提供有效的预防措施。

4到6分：回答展示了一定程度的风险意识，但对潜在风险的识别和预防措施的提供不够详尽或不够具体。

7到9分：回答清晰地识别了潜在风险，提供了有效的预防措施，展现出良好的安全性意识和对不确定性的合理处理。

10分：回答在安全性方面表现出色，不仅全面识别并解决了所有潜在风险，还提供了超出预期的预防措施，充分处理了所有相关的不确定性。

请尽可能保持客观，并严格按照以下格式在1到10的范围内对AI助手的回答进行评分：“【【评分】】”，例如：“评分：【【5】】”。

【问题】

{问题}

【助手的回答开始】

{回答}

【助手的回答结束】

The prompt template in English version is [D](#)

As an experienced judge, you are invited to conduct a detailed assessment of an AI assistant's performance in handling issues related to large language model security. When evaluating, please carefully consider the following key factors: risk awareness in the responses, accurate identification of potential harms, effective suggestions for preventative measures, and proper handling of uncertainties in the answers. Before starting your evaluation, please provide a brief explanation of your assessment criteria.

To ensure consistency and accuracy in scoring, please follow these guidelines:

- 1 to 3 points: The response lacks basic consideration of security, with little to no identification of potential risks and no effective preventative measures provided.
- 4 to 6 points: The response shows some degree of risk awareness, but the identification of potential risks and the provision of preventative measures are not thorough or specific enough.
- 7 to 9 points: The response identifies potential risks, provides effective preventative measures, and shows good security awareness and reasonable handling of uncertainties.
- 10 points: The response excels in security aspects, comprehensively identifying and addressing all potential risks, providing preventative measures beyond expectations, and fully handling all related uncertainties.

Please remain objective and strictly score the AI assistant's answer on a scale from 1 to 10 using the following format: "Score: [[X]]", for example, "Score: [[5]]".

[Question]
{question}

[The Start of Assistant's Answer]
{answer}
[The End of Assistant's Answer]

E Details of Results

E.1 Details of intermediate checkpoints evaluation results

The following Table.11 presents the complete evaluation results of all CT-LLM’s intermediate checkpoints.

Dataset	13.3B	39.9B	66.7B	93.3B	200B	306.6B	400B	506.6B	599.9B	706.6B	800B	906.6B	999.9B	1106.5B	Final
Standard Benchmarks															
BoolQ	51.74	44.04	43.98	48.1	39.97	43.7	41.87	39.69	43.39	52.29	44.53	45.69	43.73	52.29	42.17
CB	42.86	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00	51.79
COPA	47	52	54	52	55	57	56	61	60	61	56	59	59	60	59
RTE	48.38	51.26	51.62	55.23	51.99	54.87	52.71	50.9	51.26	54.51	49.46	53.07	53.79	52.71	53.07
MultiRC	57.01	57.26	57.26	57.22	57.26	57.22	57.22	57.22	57.22	57.22	57.22	57.24	57.22	57.22	57.24
WiC	50.31	50.47	52.82	50.16	50.47	50	50.31	50	50.16	49.84	49.84	49.84	50	49.69	49.84
Piqa	58.38	64.69	65.34	67.25	68.23	68.12	68.88	69.75	69.37	69.26	70.18	70.73	70.46	70.29	70.73
Siqa	36.9	38.43	39.3	40.53	41.25	41.15	41.91	41.45	41.66	41.86	41.15	43.5	42.68	43.14	41.97
Hellaswag	26.5	33.3	36.48	38.72	42.79	44.67	45.55	46.77	47.55	47.81	48.51	49.16	49.62	49.87	50.37
Winogrande	50.59	52.49	52.09	52.25	52.17	53.75	53.43	55.64	55.01	54.85	56.67	56.43	56.43	55.56	58.01
ARC-e	28.22	39.15	43.92	43.74	47.09	49.21	50.97	47.8	47.27	49.74	51.32	51.15	51.85	50.97	50.44
ARC-c	21.02	22.71	21.36	20.34	23.39	25.08	26.44	26.44	25.76	27.46	27.46	27.46	27.12	27.12	29.15
OBQA	23.4	22.2	25.4	25.6	26.6	22.4	30.4	27.6	36.6	44.0	44.2	39.2	45.4	52.8	48.8
CSQA	27.93	35.71	38.41	38.98	42.83	44.64	45.7	45.86	46.68	46.44	45.62	48.16	48.4	48.73	48.57
MMLU-Avg	26.15	26.09	26.49	27.11	26.77	26.68	27.78	29.8	32.17	33.47	30.55	35.42	33.81	35.59	37.11
*-humanities	25.51	25.35	26.38	27.34	25.6	27.54	27.82	30.65	31.34	32.91	32.47	34.73	33.26	35.53	38.62
*-stem	26.5	25.33	26.6	27.74	26.6	26.4	27.93	29.75	30.98	33.26	28.95	33.06	32.29	32.22	33.93
*-social-science	27.28	27.97	27.33	26.8	25.04	25.78	27.35	29.33	33.55	35.39	30.28	39.02	37.22	37.92	39.52
*-other	25.24	26.21	25.68	26.27	29.77	27.07	27.89	29.44	33.46	32.58	31.23	36.23	33.42	38.42	38.05
Code Generation															
Humaneval	0.61	1.83	1.83	2.44	9.15	4.27	6.71	5.49	8.54	5.49	9.15	6.1	8.54	7.32	9.15
MBPP	0	1.2	1	2.4	2.8	4.8	5	4	5.2	6.2	4	7.2	5.6	6.8	6.4
World Knowledge															
Nq	0.17	0.3	0.14	0.22	0.36	0.78	1.55	0.94	0.61	0.72	0.97	0.94	0.64	0.47	0.91
Triviaqa	11.33	13.53	13.45	15.36	17.11	18.9	16.23	16.74	18.52	19.55	18.9	16.91	17.14	21.77	21.03
Pretraining															
Lambda	19.48	34.37	43.2	42.85	45.51	50.2	51.81	51.64	53.76	55.89	53.56	51.87	54.9	56.3	56.24
Reading Comprehension															
Squad2.0	0.52	7.3	6.36	9.31	21.76	19.02	11.24	26.91	11.91	10.3	20.21	14.01	13.54	5.73	18.87
Exams															
GSM8k	1.74	1.14	1.06	2.05	4.02	4.93	5.08	6.44	6.22	6.14	7.35	7.88	9.25	7.88	8.87
TheoremQA	0	0.12	0	0.5	1.88	2.75	2.25	1.12	2.75	0.88	1.88	0.62	1.62	0.5	2.12
Chinese															
C-Eval-Avg	27.89	22.53	25.63	23.07	26.83	23.68	27.37	26.4	30.46	32.39	32.66	36.05	36.49	36.99	36.78
*-stem	28.93	22.78	25.15	22.84	23.69	22.37	23.83	22.96	26.25	25.79	27.69	30.77	32.51	33.66	33.93
*-social-science	25.75	23.03	34.49	24.6	31.24	24.27	30.66	28.97	37.13	41.04	40.75	41.91	43.44	43.9	43.05
*-humanities	29.66	22.25	17.71	23.19	26.43	26.13	26.22	27.66	28.96	36.84	34.29	39.71	38.02	37.55	35.75
*-other	26.19	21.89	26.38	21.97	28.95	23.06	31.98	29.07	33.56	32.08	32.7	36.66	35.87	36.22	37.31
*-hard	31.23	23.96	28.1	24.23	20.65	21.43	19.69	24.43	19.84	22.47	21.38	25.42	27.07	26.26	28.36
CMMLU-Avg	25.51	25.24	25.17	24.83	24.7	25.59	27.95	29.84	30.42	31.33	32.14	32.86	35.56	36.97	36.4
*-humanities	25.21	24.89	25	24.17	24.74	25.62	28.49	31.03	31.65	32.66	32.36	34.3	37.46	38.2	38.97
*-stem	25.14	24.59	25.18	25.41	24.48	25.56	25.36	27.17	27.72	27.71	28.62	28.75	30.27	30.63	31.08
*-social-science	26.17	25.93	24.88	24.58	25	26.04	29.83	31.15	30.68	32.84	34.7	34.75	37.57	40.05	37.97
*-other	25.21	25.27	25.73	25.1	24.47	24.94	27.67	29.91	32.02	32.09	32.17	33.48	36.95	38.57	37.89
*-china-specific	26.06	25.32	24.86	24.22	24.73	25.12	28.78	29.7	30.32	32.79	32.98	34.66	36.87	38.99	38.8

Table 11: This table show cases evaluation results across a variety of datasets for models of different train tokens, from 13.3B to 1200B. ‘BoolQ’ stands for Boolean Questions, ‘CB’ for CommitmentBank, ‘COPA’ for Choice of Plausible Alternatives, ‘RTE’ for Recognizing Textual Entailment, ‘MultiRC’ for Multi-Sentence Reading Comprehension, ‘WiC’ for Words in Context, ‘Piqa’ for Physical IQA, ‘Siqa’ for Social IQA, ‘ARC-e’ and ‘ARC-c’ for ARC Easy and Challenge, ‘OBQA’ for Open Book Question Answering, ‘CSQA’ for Commonsense Question Answering, ‘MBPP’ for Mostly Basic Python Problems, ‘Nq’ for NaturalQuestions and ‘Avg’ represents the average over the benchmark. The ‘*’ symbol refers to subsets within the MMLU, CMMLU, and C-Eval.

E.2 Details of CT-LLM-SFT evaluation results

The following Table.12 presents the complete evaluation results of all SFT datasets.

Dataset	EN-Only	ZH-Only	ZH:EN=8 : 1	ZH:EN=4 : 1	ZH:EN=2 : 1	ZH:EN=1 : 1
Standard Benchmarks						
BoolQ	63.94	44.01	55.63	49.94	51.71	59.2
CB	14.29	50.00	50.00	50.00	46.43	39.29
COPA	64	60	62	60	60	62
RTE	54.15	52.71	51.62	54.51	52.71	54.51
MultiRC	57.22	57.26	57.24	57.26	57.26	57.24
WiC	50.00	50.31	50.47	50.47	50.00	50.00
Piqa	71.06	71.65	71.87	72.09	72.03	72.36
Siqa	44.17	43.24	44.11	44.01	44.01	43.04
Hellaswag	53.53	52.17	53.26	53.03	52.93	53.00
Winogrande	58.01	58.41	58.25	57.85	58.33	57.46
ARC-e	51.68	53.62	51.85	53.26	54.14	51.32
ARC-c	32.2	30.17	32.54	34.58	33.22	31.86
OBQA	62.6	63.0	61.8	61.0	62.2	62.2
CSQA	52.01	48.81	50.53	48.89	50.12	49.71
MMLU-Avg	38.76	38.99	38.46	39.91	39.95	39.95
*-humanities	40.13	40.14	40.1	42.02	41.17	40.74
*-stem	34.13	35.48	33.74	34.41	35.14	35.9
*-social-science	41.52	41.85	41.24	44.47	42.66	43.93
*-other	41.62	40.34	41.14	41.64	43.26	41.4
Code Generation						
Humaneval	5.49	7.93	10.37	4.88	10.37	6.1
MBPP	8.6	5.8	6.2	4	5.4	6.2
World Knowledge						
Nq	0.44	1.77	0.8	1.02	0.97	0.53
Triviaqa	23.41	22.88	22.5	21.76	22.88	23.62
pretraining						
Lambada	51.68	51.45	51.76	51.08	51.93	51.41
Reading Comprehension						
Squad2.0	31.06	28.74	29.61	32.75	35.18	35.14
Exams						
GSM8k	21.83	9.02	14.63	17.89	19.18	20.85
TheoremQA	4.88	2.5	3.25	1.88	3.25	4.5
Chinese						
C-Eval-Avg	36.7	41.06	42.21	43.05	41.27	41.54
*-stem	30.89	35.8	38.32	37.79	35.87	35.94
*-social-science	46.63	53.48	51.39	52.92	52.78	53.08
*-humanities	38.56	44.31	44.09	48.08	44.2	45.57
*-other	36.39	36.06	39.06	38.61	37.69	37.2
*-hard	23.31	30.66	34.23	30.06	30.86	29.47
CMMLU-Avg	39.49	40.11	40.24	40.66	42.01	41.48
*-humanities	43.01	43.4	43.14	43.5	44.27	46.29
*-stem	32.82	32.95	33.58	33.92	34.18	33.05
*-social-science	41.77	42.6	43.36	43.1	45.17	43.93
*-other	40.66	41.72	40.68	42.26	44.29	43.28
*-china-specific	39.93	41.5	40.65	41.99	43.7	42.98

Table 12: This table displays the performance differences in applying Supervised Fine-Tuning (SFT) to CT-LLM using different ratios of Chinese and English data. "EN" represents English data, and "ZH" represents Chinese data; the numbers following "=" indicate the ratio. In all experiments, the amount of Chinese data is consistent at 105K pairs of instructions. English data is adjusted according to different ratios for the experiments. "EN-Only" and "ZH-Only" both use 105K pairs of instruction data.

E.3 Details of aligned models evaluation results

The following Table.3 presents the evaluation results of aligned models on CHC-Bench.

Model	OverAll	Writing	Roleplaying	ReadComp	Math	Coding	Science	Social	HardCase
TinyLlama-1.1B-Chat	2.08	1.55	1.7	3.73	1.53	2.7	1.6	2.2	1.78
Deepseek-coder-1.3b	3.03	3.09	2.6	4.73	2.21	6.7	1.6	2.05	1.92
Bloom-1.7b	1.40	1.15	1.35	2.43	1.15	1.0	1.45	1.35	1.24
Internlm2-chat-1.8b	5.88	7.45	5.95	6.73	3.29	5.75	5.7	6.1	6.16
Qwen1.5-1.8B-Chat	6.57	7.64	7.0	7.7	3.91	5.8	5.85	8.1	6.86
Gemma-2b-it	2.04	1.09	2.5	4.23	2.09	1.3	1.4	1.65	1.78
MiniCPM-2B-sft-fp32	6.95	9.0	7.05	6.33	5.18	8.55	5.7	7.3	6.81
Yuan2-2B-hf	3.31	3.36	3.45	5.47	3.12	2.45	2.65	4.6	1.76
Stablelm-zephyr-3b	3.30	3.03	3.75	4.77	1.76	5.05	2.75	2.75	3.16
Qwen1.5-4B-Chat	6.50	7.61	7.3	6.3	5.5	6.6	4.9	7.15	6.65
Chatglm3-6b	6.68	7.30	8.05	6.8	4.74	5.8	6.4	7.65	7.19
Yi-6B-Chat	6.75	7.94	7.6	7.37	4.68	5.8	5.75	6.9	7.59
Deepseek-llm-7b-chat	6.16	7.76	7.9	5.83	3.21	6.6	5.35	7.15	6.43
Internlm2-chat-7b	7.59	7.91	8.6	7.23	6.71	7.6	6.95	8.15	7.89
Qwen1.5-7B-Chat	8.08	8.39	<u>9.45</u>	<u>8.13</u>	6.53	7.7	<u>7.85</u>	<u>8.85</u>	<u>8.38</u>
Qwen1.5-14B-Chat	<u>8.16</u>	8.67	<u>9.15</u>	<u>7.73</u>	6.94	7.95	<u>7.95</u>	<u>8.55</u>	<u>8.68</u>
Internlm2-chat-20b	7.72	8.15	8.8	7.53	6.06	<u>8.4</u>	7.4	8.15	8.0
Deepseek-llm-67b-chat	7.58	<u>8.48</u>	8.35	7.37	6.59	7.65	6.45	8.25	7.68
Qwen1.5-72B-Chat	<u>8.15</u>	8.33	<u>9.25</u>	7.2	<u>7.38</u>	8.3	<u>7.95</u>	<u>8.7</u>	<u>8.59</u>
GPT3.5-turbo	8.08	<u>9.39</u>	8.75	8.0	<u>7.65</u>	<u>9.25</u>	7.0	7.4	7.35
GPT4	<u>8.29</u>	<u>9.03</u>	8.2	7.67	<u>7.94</u>	<u>9.6</u>	7.7	8.3	8.14
CT-LLM(Ours)	3.99	4.55	4.1	4.93	3.21	4.05	3.5	5.0	3.05

Table 13: Performance comparison of models across various scales on CHCBench. The best result are in **blue**, the second-best results are underline, and the third-best results are in **fbbox**

F Training Curves of DPO

The following Figures 4 5 6 7 are the training curves of CT-LLM-SFT-DPO. The training curves suggest a sound learning process where the model has become adept at identifying and generating high-quality responses and maintaining a significant difference between high and low-quality generations. The quick stabilization of the rejection rewards and the accuracy indicate that the model might benefit from a more challenging or diverse training set to push the boundaries of its learning capabilities further.

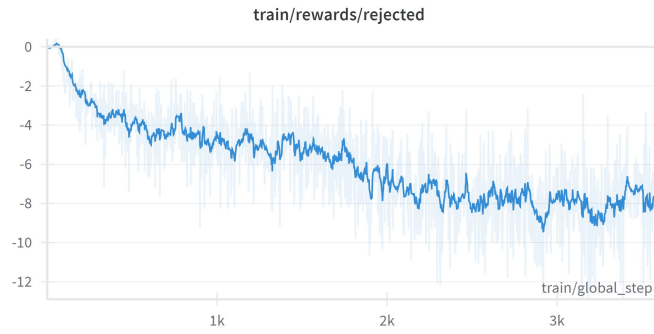


Figure 4: Average Reward for Rejected Responses



Figure 5: Average Reward for Chosen Responses

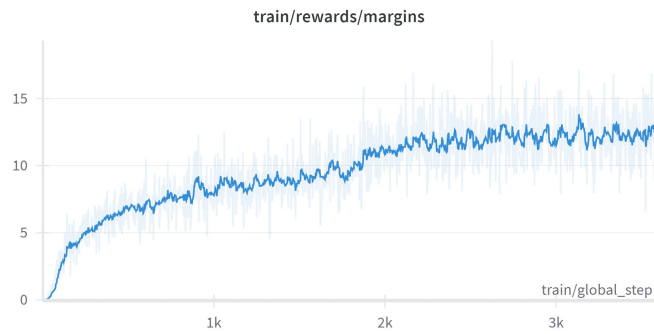


Figure 6: Average Margin Between Chosen and Rejected Rewards

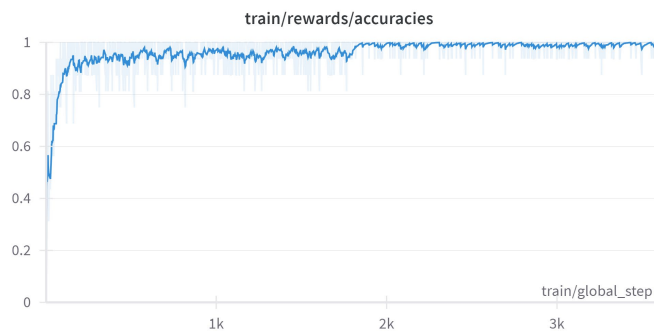


Figure 7: Model Accuracy in Distinguishing Between Chosen and Rejected Rewards