WITH INSTRUCTION BACK-TRANSLATION

Anonymous authors

000

001

003 004

006

012 013

014

015

016

017

018

019

021

023

024

025

026

027 028 029 Paper under double-blind review

Abstract

In this paper, we introduce Kun¹, a novel approach for creating high-quality instruction-tuning datasets for large language models (LLMs) without relying on manual annotations. Adapting a self-training algorithm based on instruction backtranslation and answer polishment, Kun leverages unlabelled data from diverse sources such as Wudao, Wanjuan, and SkyPile to generate a substantial dataset of over a million Chinese instructional data points. This approach presents a novel departure from traditional methods by using a self-curation process to refine and select the most effective instruction-output pairs. Our experiments with the 6B-parameter Yi model across various benchmarks demonstrate Kun's robustness and scalability. Our method's core contributions lie in its algorithmic advancement, which enhances data retention and clarity, and its innovative data generation approach that substantially reduces the reliance on costly and time-consuming manual annotations. This methodology presents a scalable and efficient solution for improving the instruction-following capabilities of LLMs, with significant implications for their application across diverse fields.

1 INTRODUCTION

The development of large language models (LLMs) relies on human-annotated datasets, yet the creation of such datasets typically faces scalability issues due to the significant resources required. Our study introduces Kun, a novel approach leveraging unlabelled data to create a high-quality instruction-tuning dataset. This method diverges from manual annotations, employing a self-training algorithm that includes a unique process called AP (Answer Polishment),

AP is central to Kun's strategy. It addresses a critical challenge in the Humpback (Li et al., 2023c) method, where raw data, once labeled, are directly used in instruction datasets. The unscreened raw data often mismatches between instructions and responses, as raw data may not inherently align with the instructional context. AP refines this raw data, ensuring a tighter correlation between the instructions and responses through a back-translation process. This leads to a dataset where each instruction-output pair is more coherent and contextually relevant. Unlike methods dependent on LLMs (Peng et al., 2023; Taori et al., 2023b; Zheng et al., 2023), Kun offers an independent and scalable approach to instruction-based training.

We opt for the 6B-parameter Yi model due to its open-source nature and dependable performance². Its efficacy is tested and proven across diverse dataset sizes, including widely recognized benchmarks like C-EVAL (Huang et al., 2023) and CMMLU (Li et al., 2023b). To evaluate the performance of the model, we design a comprehensive human evaluation which contains 500 prompts from ShareGPT-zh, covering various tasks. Responses generated by our model are compared with those from other models, showcasing the superiority of our *Kun-52k* variant. Further details can be found in 4.2.3. Additionally, we evaluate the quality of our dataset, which includes 1,000 instruction-output pairs each from sources like Wudao (Xue et al., 2022), Wanjuan (He et al., 2023), and SkyPile (Wei et al.,

¹The dataset is named Kun as Chinese pronunciation of Humpback Li et al. (2023c).

²https://github.com/01-ai/Yi



Figure 1: Overview of *Answer Polishment*. Initially, the Yi base model is fine-tuned using quality seed instruction data to create a label and a primary chat model. The label model then annotates a large amount of primary data, turning it into labeled data. This is filtered and refined by rules and the primary chat model, producing the final dataset. This dataset is used to further train the primary chat model, resulting in an highly efficient final chat model.

077 078 079

081 082

084

085

087

090 091

074

075

076

2023). This evaluation, focusing on clarity, feasibility, practicality, and alignment, ensures the high quality of our dataset. The key contributions of our work are:

- *Algorithmic Advancement*: AP in Kun enhances data retention and resolves ambiguities, leading to an expanded pool of high-quality data for fine-tuning.
- *Large-scale high quality data creation*: Over a million diverse Chinese instructional data points are produced from sources like Wudao, Wanjuan, and SkyPile, surpassing traditional crowdsourced annotations in quality and reducing reliance on manual annotation.
- 2 Related Work

Instruction Tuning Instruction tuning is widely recognized as a key technique for activating 092 LLMs to adhere to human conversational norms. (Mishra et al., 2022; Wang et al., 2022b; 2023b). 093 Instruction tuning empowers various domain-specific or task-specific LLMs, including natural 094 language generation evaluation (Jiang et al., 2023), math (Yue et al., 2023; Xu et al., 2023a; Azerbayev 095 et al., 2023), code Luo et al. (2023), music Li et al. (2023a); Deng et al. (2023), and medicine Wang 096 et al. (2023a). Instruction tuning not only tailors the models' task-specific responsiveness but also bolsters their cross-task generalization capabilities, thus enhancing performance across various 098 dynamic application contexts (Wei et al., 2021; Sanh et al., 2022; Mishra et al., 2022; Wang et al., 099 2022b). Recent studies have broadened the scope of instruction tuning to encompass a wider array of 100 general tasks, notably incorporating input from users of language models (Ouyang et al., 2022; Peng et al., 2023). 101

However, the open-source community is still lacking high-quality Chinese instruction tuning corpora.
Current datasets, like COIG (Zhang et al., 2023), BELLE (Ji et al., 2023), MOSS (Sun et al., 2023),
and OL-CC (OL-CC, 2023), face issues such as limited scope, poor quality, commercial restrictions,
or insufficient coverage. This gap hampers the advancement of LLMs in effectively processing
and executing complex Chinese instructions, highlighting the critical need for more diverse and
superior-quality datasets.

Self-Improvement of LLMs In fine-tuning LLMs, the availability of extensive, high-quality instructional data is crucial. Presently, the generation of such data mainly relies on human manual annotation, a labor-intensive method that lacks scalability for future data augmentation.

114 An alternative approach involves deriving 115 instructional data from more advanced 116 LLMs (Taori et al., 2023a), exemplified by ex-117 tracting authentic instruction-response sets from 118 dialogues within the GPT series models (Wang et al., 2022a). A more refined technique utilizes 119 the Self-Instruction framework, autonomously 120 generating additional instructional data from ini-121 tial seed data. Combined with the Ada-Instruct 122 or the Evol-Instruct framework, this approach 123 can transform basic instructions into complex 124 ones, specifically tailored for distinct tasks (Cui 125 & Wang, 2023; Luo et al., 2023). 126

Nevertheless, these instruction generation 127 methodologies all require a robust teacher. The 128 ultimate potential of the model is limited by the 129 teacher's expertise or resource expenditure (Li 130 et al., 2023c). To overcome this limitation, the 131 SPIN (Chen et al., 2024) framework incorpo-132 rates a self-play mechanism, It generates train-133 ing data from previous iterations, refining its 134 strategy by distinguishing between responses



Figure 2: The top 10 categories in each of these three areas: Academic Disciplines, Industry Sectors, Text Type

generated autonomously and those derived from human-annotated data.

This gradual process elevates the LLM from a nascent model to a robust one. Considering the vast amount of knowledge present in web text, Humpback Li et al. (2023c) introduces a technique based on Instruction Backtranslation. This method allows a base model to independently utilize vast amounts of unlabeled data to generate a high-quality instruction tuning dataset. However, empirical findings indicate that the effectiveness of this method is still constrained by the seed model's performance and its ability to discern truly high-quality data.

142 143

144

3 Method

145 Our training methodology necessitates a foundational model, high-quality seed instruction data, and 146 a substantial volume of unlabeled data, with the primary source being web text. Given the extensive 147 content diversity inherent in large-scale web documents, which encompass a wide array of topics such 148 as music, art, technology, etc., reflecting the broad spectrum of human interests, certain subsets of 149 these documents may be more apt for generating instructions. Unlike labeled data, these documents 150 lack predetermined outcomes or objectives. This method involves the refinement and optimization of 151 data selection during the fine-tuning process of Large Language Models (LLMs), as illustrated in 152 Figure 1. This approach allows for the collection of a significant volume of instructional data at a low cost, circumventing the exorbitant expenses associated with manual labor, in a manner that is both 153 academically rigorous and professional. Our method consists of two main steps: 154

- 15
- 156 157

• Supervised Fine-Tuning (SFT) with High-Quality Seed Data: This involves using SFT on the base model with high-quality seed data to create two models - the label model for annotating primary data and the primary chat for improving data quality.

Quality Assessment and Refinement in Primary Chat: The primary chat assesses and refines the label model's output. This repeated process produces a lot of high-quality data, essential for the primary chat's further training. It leads to a high-performance final chat model, trained extensively with superior data.



filtration of the candidate labeled data.

Our experiments tested two different filtering methods, each with its strengths and weaknesses:

216	
217	诸生近我以下文太是否为疑问句或近使句,若是疑问句或考近使句则回答【是】 若不是疑问句或近使句则回答【否】 例如·
218	
219	输入: 在进行消防安全评估时,需要注意哪些问题: 回答: 【是】; 输入: 有哪些适合艺术小白阅读的艺术类书籍? 回答: 【是】;
220	输入: 请给我一个关于"护理专业"的简要介绍。回答: 【是】;
221	输入:外研版八年级上册英语 Module 3 话题写作标题: My Favorite Place 回答: 【否】; Prompt
222	输入: 泰安一平,作为一所历史悠久的教育机构,始终来承见良的传统,致力于提供优质的教育资源。 回答 【 召】 请判断以下输入文本:
223	Please tell me if the following text is a question or an imperative sentence. If it is a question or imperative sentence, answer "Yes"; if it is not a question or imperative sentence, answer "No". For example:
224	Input: What is trace element fertilizer? Answer: [Yes]; Input: What issues should be considered during a fire safety assessment? Answer: [Yes];
225	Input: What are some art books suitable for art novices? Answer: [Yes]; Input: Please give me a brief introduction to the "nursing profession". Answer: [Yes]:
226	Input: Please express your views on the importance of corporate culture construction. Answer: [Yes];
227	input: Foreign Research carton Grade & English Module 3 Writing Topic: My Favorite Place Answer: [No]; Input: Tai'an No.1 Middle School, as a long-established educational institution, has always adhered to fine traditions and is committed to providing
228	quality educational resources. Answer: [No] Please judge the following input text:
229	
230	
231	
232	Figure 4: Filter prompt we use to screen out unsuitable content for instructions.
233	• Comprehensive Scoring of Labeled Data: This method evaluates the full labeled data set
234	based on a combined score, including instructions and outputs
235	
236	• Focused Scoring of Instruction Component: This technique only assesses the instruction
237	part (output from the label model). High-scoring instructions are chosen, and then the output
238	part of these selected data is refined.
239	Our analysis shows that the second method is more effective then the first. In the first method, good
240	outputs are often discarded because of poor instructions from the label model, and the reverse is also
2/11	true causing unnecessary exclusions. Moreover, this approach occasionally retains data with one
241	noor quality instruction because the corresponding output is high quality and vice versa leading to
242	uneven data quality and negatively impacting further training.
243	
244	In contrast, the second method only scores the instruction component, as in instruction tuning for
240	LLM, instructions are often considered more important than outputs, Yet, it doesn't assess the output,
240	sometimes leading to suitable instructions paired with unsuitable outputs. To address this, we use the
247	primary chai model to evaluate and refine the instructions and outputs, ensuring they angle well. This approach produces high quality labeled data. The score and rafine prompts we used in this process
248	approach produces high-quality labeled data. The score and tenne prohipts we used in this process are shown in Figure 10
249	are shown in Figure 10.
250	Utilizing the substantial volume of top-quality labeled data from these procedures, we further train
251	the main chat model, achieving a high-performance final model, as shown in Experiments.
252	
253	Δ Experiments
254	
255	In this spation, we comprehensively detail the experimental presedures and methodologies employed
256	in our study
257	in our study.
258	
259	4.1 EXPERIMENTAL SETUP
260	We first detail the experimental estury used in our study security the base we deleted with the first
261	we first detail the experimental setup used in our study, covering the base model selection, fine-tuning
262	process, basenne comparisons, and the evaluation methods.
263	
264	4.1.1 BASE MODEL & FINETUNING
265	In our experiments, we utilize the Yi model with 6R parameters ³ developed by 01 AL as our
266	foundational language model for fine-tuning. Renowned for its proficiency in both English and
267	Chinese the Vi series has shown impressive results on global benchmarks like the AlnacaFyal

Chinese, the Yi series has shown impressive results on global benchmarks like the AlpacaEval Leaderboard (Dubois et al., 2023; Li et al., 2023d) and SuperCLUE (Xu et al., 2023b).

³https://huggingface.co/01-ai/Yi-6B

Sauraa	Instruction Quality			Output Quality		
Source	Clarity%	Feasibility%	Practicality%	Excellent%	Pass%	Fail%
Wudao	96.67	96.40	96.87	69.50	20.03	10.47
Wanjuan	98.27	97.63	96.57	85.63	11.13	3.24
Skypile	98.90	98.37	95.40	42.73	40.43	16.84
ALL	97.94	97.47	96.28	66.00	23.87	10.13

284 285

286

293

295

Table 1: Manual Quality Analysis of Synthetic Data Generated by Kun.

The fine-tuning process is carried out using varying sizes of our high-quality, instructionally curated dataset. This phase is executed on a computing setup with 32 Nvidia A800 GPUs, amounting to a total of 192 GPU hours. We adopt a learning rate of 2e-5 and a batch size of 16, aiming for an optimal balance between computational efficiency and model performance. All the models have been fine-tuned with the same number of update steps.

4.1.2 BASELINES

For the Kun dataset, we annotated command data from three sources: Wudao, Wanjuan, and Skypile.
Quantitative details of this augmented dataset are provided in Figure 7. In evaluating the performance
of Kun, our study contrasts it with data curated from four prominent Chinese open-sourced datasets,
including COIG (Zhang et al., 2023; BAAI, 2023a;b), OL-CC (OL-CC, 2023), and BELLE (Ji et al.,
2023). These datasets are unique in their composition and focus, providing a comprehensive basis for
comparison.

4.1.3 EVALUATION

Human Evaluation. To assess the general quality of model responses, we conduct human evaluations
 using a test set of 500 prompts sampled from ShareGPT-zh. These prompts, derived from real world
 user inquiries, encompass a diverse array of tasks, such as creative writing, information seeking,
 providing guidance, logical reasoning, storytelling, problem-solving, etc.

For the evaluation, responses generated by different models for each prompt are presented side-byside. Human evaluators are asked to choose their preferred answer, providing a direct comparison
of model performance. In total, eight models were compared. For this evaluation, we engage a team
of experienced crowdsource annotators, ensuring a balanced and unbiased assessment. Detailed
examples that show the comparison process can be found in Figure 11. Standard Benchmarks.
In addition to human evaluations, the models are also assessed using two standard benchmarks for
Chinese LLM evaluation: C-EVAL (Huang et al., 2023) and CMMLU (Li et al., 2023b).

These evaluation methods, comprising both human judgment and standardized benchmarks, offer a multifaceted perspective on the capabilities of the Kun model, enabling a thorough comparison with existing models and datasets.

310 311

4.2 AUGMENTATION DATA STATISTICS 312

In this section, we delve into the detailed statistical analysis and diversity assessment of our augmented
 dataset, as well as the rigorous quality evaluation conducted using human annotators. By exploring
 the comprehensive scale, varied nature, and assessed quality of the instruction-output pairs, we aim
 to highlight the robustness and reliability of the data curated for our study.

317

318 4.2.1 STATISTICS AND DIVERSITY319

Our work involve the purification of approximately 377,592 high-quality instruction-output pairs from the Wudao, Wanjuan, and Skypile datasets.We analyze a 10% subset of instructions from the past 20 years, revealing significant temporal diversity with 56% of instructions from the recent three years (Figure 8). The variation in instruction and output lengths, analyzed using the Yi-6B model, is shown in Figure 3, reflecting content complexity. To assess instruction diversity, we categorize them into 24 academic disciplines, 16 industry sectors, and 15 text types as per Wikipedia⁴ using the Qwen-72B-Chat⁵ (Bai et al., 2023). Repeated for accuracy, this categorization highlights the data's range, as shown in Figure 2, where the top 10 categories in each area signify its broad scope.

Source	Consistency%
Clarity	96.87
Feasibility	97.73
Practicality	97.43
ALL	92.13

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

351

352

353

354

355

356

357

358

359

360

377

Table 2: The proportion of identical evaluations from three assessors on a single dimension. **All**: The proportion of consistent assessments across all three dimensions within the same item.



Figure 5: Heatmap of Model Comparative Win Rates in Human Evaluations. mix-lite:Kun-26k+COIG-lite.mix-one:Kun-26k+COIG-one

4.2.2 QUALITY EVALUATION

A critical aspect of our dataset curation process is the rigorous data quality assessment. We conduct
 a comprehensive quality evaluation of the instruction-output pairs to achieve this. For augmented
 data curated from each source (Wudao, Wanjuan, and Skypile), we randomly select 1,000 instruction output pairs, resulting in 3,000 pairs subjected to an independent quality assessment.

Instruction Quality. For instruction quality, a team of 24 professional annotators with a bachelor's degree or higher evaluates each instruction across three key dimensions: clarity, feasibility, and practicality. Each aspect is assessed with a simple yes/no answer, providing a straightforward yet effective measure of instruction quality. The evaluation criteria are as follows:

- **Clarity:** Evaluators determine whether the instruction was unambiguous and coherent, encompassing necessary information without any vague terms or explanations.
- Feasibility: Evaluators assess whether the instruction was valid and answerable within the context and scope of the model's capabilities.
- Practicality: Evaluators judge the relevance of the instruction in everyday scenarios.

Output Quality. The quality of the outputs is evaluated based on their alignment with the instructions.
 Evaluators are asked to rate each output as *Excellent*, *Pass*, or *Fail*, based on how well it met the requirements and intent of the instruction.

To ensure objectivity and reliability, three different evaluators evaluate each instruction-output pair. The consistency rates for the evaluation across the three dimensions of the instructions have all exceeded 90%, and the evaluation of the instruction-response are also Consistently. This results demonstrate a significant degree of consistency in their judgments. Further details on evaluating identical ratings are presented in Figure 6 and Table 2. Examples that demonstrate the process of assessing can be found in Figure 11

As indicated in Table 1, the instruction quality across all sources is consistently high, suggesting effective formulation and clarity in conveying their purposes. However, the output quality varies more noticeably among the sources. While some sources like Wanjuan exhibite a high percentage of "Excellent" outputs, others such as Skypile demonstrate a more diverse distribution of output quality. This section presents the analysis of our experiment results, encompassing human evaluation outcomes and performance on standard benchmarks.

⁴https://www.wikipedia.org

⁵https://huggingface.co/Qwen/Qwen-72B-Chat



4.2.3 HUMAN EVALUATION

Our human evaluation results are illustrated through the heatmap in Figure 5, which shows the model vs. model win rates, with color variations indicating the relative performance strengths. The heatmap highlights that the *Kun-52k* model emerges as the most dominant, followed by the mixed model, showcasing their superior ability to handle a wide range of prompts. In contrast, the baseline models *COIG-39k*, and *Belle-52k*, garner lower preference percentages. This suggests that despite their strengths, these models may not align as closely with user expectations or prompt requirements as the Kun models. Further analysis is provided in Appendix A.4.

400 401 402

4.2.4 STANDARD BENCHMARKS

The performance of the models on standard benchmarks, Table 3 presents the performance statistics of various models on the C-EVAL and CMMLU benchmarks. As shown, we evaluate numerous models, including different sizes of the Kun model, baseline models, and mixed models. Each model's performance is measured in terms of perplexity and generation quality, providing a comprehensive view of its strengths and weaknesses.

From the table, we observe that *Kun-39k* generally exhibits lower perplexity and higher generation quality, confirming its top-tier performance in language understanding and generation. Interestingly, the mixed model display robust performance, with the mixed model often outperforming *Kun-52k*. The baseline models and smaller Kun variants present mixed results, excelling in some metrics while falling short in others. These highlight potential areas for further improvement in model training and fine-tuning strategies.

414 415

5 CONCLUSION

416

417 Our approach represents a breakthrough in instruction-tuning for LLMs, utilizing a novel self-training 418 algorithm to leverage over a million quality Chinese instructional data points from diverse sources effectively. This strategy, different from manual annotations, not only enhances the instruction-419 following capabilities of LLMs but also ensures the high quality and diversity of training data. 420 Empirical evaluations using the 6B-parameter Yi model across benchmarks like C-EVAL, CMMLU, 421 and human evaluations, have demonstrated its robustness and scalability. Innovations within our 422 approach, such as AP, have notably improved data retention and clarity, offering a scalable, efficient 423 way to augment LLMs' instructional capabilities. This research not only progresses the field but 424 also broadens LLMs' application scope, offering a novel solution to the challenges in developing 425 instruction-tuning datasets. 426

427

ETHICS STATEMENT

428 429

This study adheres to the highest ethical standards, ensuring all research activities are conducted with
 a commitment to responsibility and respect for participant rights. Our ethical policy encompasses
 data usage, intellectual property rights respect, and research transparency. To safeguard data privacy

432 and security, particularly when handling unlabeled data from sources like Wudao, Wanjuan, and 433 SkyPile, we implement stringent measures to comply with data protection laws, especially concerning 434 personal information. This involves anonymizing and desensitizing data prior to utilization. In terms 435 of intellectual property rights, we ensure that all employed data and generated guiding points adhere 436 to applicable copyright laws and intellectual property agreements, thereby avoiding infringement on any third-party intellectual property rights. Moreover, we pledge to provide a comprehensive 437 account of our research methodology, detailing the processes of data generation, model training, and 438 evaluation, to facilitate reproducibility and validation by the academic community. 439

440 441

442

LIMITATIONS

443 This study introduces an innovative methodology for generating data to reduce the reliance on costly 444 and time-consuming manual annotations. However, there are limitations impacting the generalizability 445 and scalability of our findings. Firstly, the diversity of data: the methodology, while capable of 446 generating instructional data from a vast pool of unlabeled data, may produce data whose quality and diversity are constrained by the original data source's breadth and caliber. This is particularly relevant 447 when generating instructions for niche or specialized domains, where additional methods might be 448 necessary to guarantee comprehensive coverage and precision. Secondly, the generalization capability 449 of the model: although tests on the 6B-parameter Yi model confirm the methodology's efficacy, its 450 performance and applicability could differ among models of various sizes and tasks. Its effectiveness 451 might require further investigation, especially in smaller models or those not designed for Chinese 452 language processing. Thirdly, the assumptions underlying the algorithm: the study's self-training 453 algorithm relies on instruction back-translation and answer embellishment. These premises might 454 not hold across all instruction or answer types, notably for tasks demanding specialized knowledge 455 or creativity, where the generated instructional data may not adequately reflect complex cognitive 456 processes. Future Directions: Subsequent research should examine the methodology's applicability to large-scale language models across different languages and the generation of high-quality instructional 457 data in specialized fields such as medicine or law. Additionally, advancing self-training algorithms 458 to more effectively manage intricate and specialized instructions represents a crucial avenue for 459 exploration. 460

461 462

463

467

468 469

470 471

472

473

477

478

479

References

- Zhangir Azerbayev, Hailey Schoelkopf, Keiran Paster, Marco Dos Santos, Stephen McAleer, Albert Q
 Jiang, Jia Deng, Stella Biderman, and Sean Welleck. Llemma: An open language model for
 mathematics. *arXiv preprint arXiv:2310.10631*, 2023.
 - BAAI. Coig-pc, 2023a. URL https://huggingface.co/datasets/BAAI/COIG-PC.
 - BAAI. Coig-pc-lite, 2023b. URL https://huggingface.co/datasets/BAAI/ COIG-PC-Lite.
 - 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 arXiv:2309.16609*, 2023.
- Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024.
 - Wanyun Cui and Qianle Wang. Ada-instruct: Adapting instruction generators for complex reasoning. *arXiv preprint arXiv:2310.04484*, 2023.
- Zihao Deng, Yinghao Ma, Yudong Liu, Rongchen Guo, Ge Zhang, Wenhu Chen, Wenhao Huang, and Emmanouil Benetos. Musilingo: Bridging music and text with pre-trained language models for music captioning and query response. *arXiv preprint arXiv:2309.08730*, 2023.
- Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin,
 Percy Liang, and Tatsunori B. Hashimoto. Alpacafarm: A simulation framework for methods that
 learn from human feedback, 2023.

- 486
 487
 488
 488
 488
 489
 480
 480
 480
 480
 480
 480
 480
 480
 480
 481
 481
 481
 482
 483
 484
 484
 484
 485
 485
 486
 486
 486
 486
 486
 486
 486
 486
 487
 486
 487
 488
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
 486
- 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. C-eval: A multilevel multi-discipline chinese evaluation suite for foundation models. *CoRR*, abs/2305.08322, 2023. doi: 10.48550/ARXIV.2305.08322. URL https://doi.org/10.48550/arXiv.2305. 08322.
- Yunjie Ji, Yong Deng, Yan Gong, Yiping Peng, Qiang Niu, Lei Zhang, Baochang Ma, and Xiangang
 Li. Exploring the impact of instruction data scaling on large language models: An empirical study
 on real-world use cases, 2023.
- Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhu Chen. Tigerscore: Towards building explainable metric for all text generation tasks. *arXiv preprint arXiv:2310.00752*, 2023.
- Dichucheng Li, Yinghao Ma, Weixing Wei, Qiuqiang Kong, Yulun Wu, Mingjin Che, Fan Xia,
 Emmanouil Benetos, and Wei Li. Mertech: Instrument playing technique detection using self supervised pretrained model with multi-task finetuning. *arXiv preprint arXiv:2310.09853*, 2023a.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. CMMLU: measuring massive multitask language understanding in chinese. *CoRR*, abs/2306.09212, 2023b. doi: 10.48550/ARXIV.2306.09212. URL https://doi.org/10.48550/arXiv.2306.09212.
- Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Luke Zettlemoyer, Omer Levy, Jason Weston, and
 Mike Lewis. Self-alignment with instruction backtranslation. *arXiv preprint arXiv:2308.06259*, 2023c.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023d.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing
 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with
 evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. In *ACL*, 2022.
- 523 OL-CC. Openlabel-chinese conversations dataset (ol-cc), 2023. URL https://data.baai.ac. 524 cn/details/OL-CC.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with
 gpt-4. arXiv preprint arXiv:2304.03277, 2023.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=9Vrb9D0WI4.

540 Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li, Qinyuan Cheng, Hang Yan, Xiangyang Liu, 541 Yunfan Shao, Qiong Tang, Xingjian Zhao, Ke Chen, Yining Zheng, Zhejian Zhou, Ruixiao Li, Jun 542 Zhan, Yunhua Zhou, Linyang Li, Xiaogui Yang, Lingling Wu, Zhangyue Yin, Xuanjing Huang, 543 and Xipeng Qiu. Moss: Training conversational language models from synthetic data. 2023. 544 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy 545 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. 546 https://github.com/tatsu-lab/stanford_alpaca, 2023a. 547 548 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy 549 Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-following model. Stanford Center for Research on Foundation Models. https://crfm. stanford. edu/2023/03/13/alpaca. 550 html, 3(6):7, 2023b. 551 552 Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: 553 Tuning llama model with chinese medical knowledge. arXiv preprint arXiv:2304.06975, 2023a. 554 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and 555 Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 556 2022a. 558 Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, 559 Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, et al. 560 Super-naturalinstructions: generalization via declarative instructions on 1600+ tasks. In EMNLP, 561 2022b. 562 Zekun Wang, Ge Zhang, Kexin Yang, Ning Shi, Wangchunshu Zhou, Shaochun Hao, Guangzheng 563 Xiong, Yizhi Li, Mong Yuan Sim, Xiuying Chen, et al. Interactive natural language processing. 564 arXiv preprint arXiv:2305.13246, 2023b. 565 566 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, 567 Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In International Conference on Learning Representations, 2021. 568 569 Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, 570 Weiwei Lü, Rui Hu, Chenxia Li, Liu Yang, Xilin Luo, Xuejie Wu, Lunan Liu, Wenjun Cheng, 571 Peng Cheng, Jianhao Zhang, Xiaoyu Zhang, Lei Lin, Xiaokun Wang, Yutuan Ma, Chuanhai Dong, 572 Yanqi Sun, Yifu Chen, Yongyi Peng, Xiaojuan Liang, Shuicheng Yan, Han Fang, and Yahui Zhou. 573 Skywork: A more open bilingual foundation model, 2023. 574 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin 575 Jiang. Wizardlm: Empowering large language models to follow complex instructions. arXiv 576 preprint arXiv:2304.12244, 2023a. 577 578 Liang Xu, Angi Li, Lei Zhu, Hang Xue, Changtai Zhu, Kangkang Zhao, Haonan He, Xuanwei Zhang, 579 Qiyue Kang, and Zhenzhong Lan. Superclue: A comprehensive chinese large language model 580 benchmark. arXiv preprint arXiv:2307.15020, 2023b. 581 Zhao Xue, Hanyu Zhao, Sha Yuan, and Yequan Wang. WuDaoCorpora Text, December 2022. 582 583 Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 584 Mammoth: Building math generalist models through hybrid instruction tuning. arXiv preprint 585 arXiv:2309.05653, 2023. 586 Ge Zhang, Yemin Shi, Ruibo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, 587 Zekun Wang, Chenghua Lin, Wenhao Huang, and Jie Fu. Chinese open instruction generalist: A 588 preliminary release, 2023. 589 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, 590 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and 591 chatbot arena. arXiv preprint arXiv:2306.05685, 2023. 592

A APPENDIX

A.1 EXAMPLES IN HUMAN EVALUATION OF DATA AND MODEL PERFORMANCE

Figure 11 illustrates specific examples from our human evaluation process, as detailed in Sections 4.2.2 and 4.1.3. This figure includes the evaluation of both instruction and output quality, along with a comparison of model answers. For each example, the original question is provided, followed by the evaluated responses.

The instruction quality is assessed for clarity, feasibility, and practicality to evaluate the precision in
 formulating instructions. The output quality assessment focuses on the extent to which each response
 meets the instruction's requirements. Additionally, the model comparison examines which answers
 align most closely with human preferences, highlighting the practical effectiveness of the models
 under consideration.

A.2 DATA FILTERING RULES

The text provided lists the manually established rules we use to delete some low-quality data in the process of generating data to improve data quality. These rules are used to filter the generated instructions or responses in cases with very obvious low-quality characteristics, thereby ensuring the reliability of the entire process. 1. Sensitive information such as phone numbers, home addresses, etc. Further analysis is provided in Figures 9, which depict the human preference evaluation win rates of *Kun-52k* and the mixed model, respectively, when compared with other models. Notably, the mixed model demonstrates a higher win rate than Kun-52k, indicating its enhanced effectiveness in meeting human evaluators' preferences.

• Paragraphs that are largely repetitive.

- The length is too short (length i = 4).
- The text contains a large number of meaningless characters.
- Contains specific low-quality keywords.
- Refusal to answer.
- Text formatting errors.

A.3 YEARS



Figure 8: Instructions spanning over 20 years, with 56% from the last five years.

651

652

653

648 A.4 More showcases of human evaluation results

Further analysis is provided in Figures 9, which depict the human preference evaluation win rates of *Kun-52k* and the mixed model, respectively, when compared with other models. Notably, the mixed model demonstrates a higher win rate than *Kun-52k*, indicating its enhanced effectiven.



one) model, respectively, when compared with other models. Notably, the *mixed* model demonstrates a higher win rate than *Kun-52k*, indicating its enhanced effectiven

702 A.5 PROMPT

704	
705	
706	评价指令质量。根据以下指令内容使用【好】或【差】进行评价,【好】的评判标准极为严格,对【差】则非常宽松。 高氏是长金枝树:
707	"同项里有之行任"。 明确性与精确性、清晰地陈述目标,提供详细的指令,并避免任何模糊的解释或术语。
708	完整性:确保包含所有必要的信息,执行时不需要猜测。 可行性:确保指令是实际的,并适合接收者的能力。
709	逻辑顺序: 按逻辑顺序描述步骤。 -低质量指令特性: Prompt
710	含隐私信息:如人名等。 樟糊和不明确:目标、步骤、要求或目的不清晰。
711	信息问题: 信息缺失、错误或过时。 不切实际: 報出执行者的能力或资源
712	复杂性:使用冗余的词汇、术语或涉及多模态内容,如图片或表格。 逻辑语言:
713	这相优心。少张顺序小月均碱之上下入月束。 示例: 12.4.4.934.4.4.5.934.4.4.5.934.4.4.5.1.4.5.1.4.5.1.4.5.1.1.5.1.1.5.1.1.5.1.5
714	- 【请解释什么是 联合开发, 发具风险与收益。】合: 【好】 - 【你认识黄丰吗?】 答: 【差】
715	-【什么是排他性许可!->什么是 什么是】答: 【差】 -【什么是民窿? 什么是民窿?】答: 【差】
716	-【教育具有哪些社会发展功能?教育具有哪些社会发展功能?】答:【差】 -【教育具有哪些社会发展功能?】答:【好】
717	-【帮我搜索一下"李四"的百度百科。】答: 【差】
718	Assess the quality of instructions. Evaluate the following content with either "Good" or "Poor", where "Good" is judged by very strict standards and "Poor" is judged very leniently.
719	- Characteristics of high-quality instructions: Clarity and Precision: Clarity instructions:
720	Completeness: Ensures all necessary information is included, leaving no need for guesswork during execution.
721	Logical Order: Describes steps in a logical sequence.
722	- Characteristics of low-quality instructions: Contains Private Information: Such as personal names.
723	Vague and Unclear: Goals, steps, requirements, or purposes are unclear. Information Issues: Information is missing, incorrect, or outdated.
724	Impractical: Beyond the executor's capabilities or resources. Complexity: Uses redundant vocabulary, jargon, or involves multi-modal content such as images or tables.
725	Logical Disarray: Steps are out of order or lack contextual background.
726	- [Please explain what is "joint development" and its risks and rewards.] Answer: [Good]
727	- [bb you having reing:] Allswer. [roor] - [What is an exclusive license? -> What is] Answer: [Poor]
728	- [What is a civil purrow? What is a civil purrow?] Answer: [Poor] - [What social development functions does education have? What social development functions does education have?] Answer: [Poor]
729	- [What social development functions does education have?] Answer: [Good] - [Help me search for "Li Si" on Baidu <u>Baike</u> .] Answer: [Poor]
730	
731	(autout) 法回答/instruction)*
732	{output},请回答:{instruction}
733	{output}, please answer: {instruction}*
734	*{output} and {instruction} refer to the output and instruction parts of the current labeled data, respectively. This prompt aims to use the output as background information, employing the primary chat model to re-answer the question to obtain a new output suitable as training data.
735	
736	
737	
738 Fig	gure 10: The top section comprises score prompt used to assess the quality of labeled data
739 ins	tructions. The bottom section features refine prompt for refining the output part of the labeled data.
740	
741	
742	
743	
744	
745	

756	
757	
758	
759	instruction:蛋白质晶体学走研究什么的:主要研究于权定什么:可以靠次哪些问题: Q1:上述指令是否明确? (即清晰地陈述目标,包含所有必要的信息,无模糊的解释或术语)评价:[是]
760	Q2:上述指令是否实用?(即指令是有效且可回答的)评价:[是] Quality Q3:上述指令是否实用?(即在日常生活中经常出现)评价:[否]
761	Please evaluate the following instruction:
762	instruction: What does protein crystallography study? What are its main research methods? What problems can it solve? Q1: Is the above instruction clear? (i.e., clearly states the goal, contains all necessary information, and has no ambiguous explanations or terms)
763	Evaluation: [Yes] Q2: Is the above instruction feasible? (i.e., the instruction is valid and answerable)
764	Evaluation: [Yes] Q3: Is the above instruction practical? (i.e., frequently occurs in everyday life)
765	Evaluation: [No]
766	
767	Question:家用生活电器有哪些?
769	Answer: 家用生活电器是指那些用于家庭日常生活、提供使利和改善生活质量的电气设备。它们涵盖了多个类别, 包括但不限于以下几种: 1. 制冷电器 - 如冰箱、冷柜等。 2. 空调器 - 包括窗式空调机、分体式空调机和中央空
700	调系统等。 3. 清洁电器 - 例如洗衣机、干衣机、吸尘器和地板打蜡机等。 4. 烹饪电器 - 比如电饭锅、微波炉、 烤箱、电磁灶等。 5. 热水供应电器 - 包含热水器(燃气热水器、电热水器)和热泵等。 6. 照明电器 - 包括各种
709	灯具、LED灯泡以及智能照明控制系统等。7.个人护理电器 像电动剃须刀、吹风机、美容仪和个人按摩器材等。 Output
770	主要包括电视机、音响系统、游戏 consoles 等。10. 智能家居电器 - 如智能温控器、智能插座、安全监控摄像头 Quality
//1	等。随着将我的远少,还有许多新兴的家用电路个断满线,这些电路通常结合了看能化功能,可以通过于机应 用程序或语音控制来操作,为人们的生活带来更多的便捷与舒适。需要注意的是,上述列表并不是详尽的,不
772	同的文化和地区可能有特定的家用电器类型。此外,随着时间的推移,一些旧式的电器可能会被更高效、节能 的新产品所取代。
773	请对此答案做出评价,它是否回答了指令所提出的问题? A 正确且令人满意的回答
774	B 可以接受的回答。但质量一般 C 错误或质量较差的回答
775	评价: [A]
776	Question: What are some household appliances?
777	Answer: Housenoid appliances refer to electrical devices used in daily family life to provide convenience and improve the quality of life. They encompass multiple categories, including but not limited to the following:
778	 Refrigeration Appliances: Such as refrigerators, freezers, etc. Air Conditioners: Including window air conditioners, split systems, and central air conditioning systems.
779	 Cleaning Appliances: For example, washing machines, dryers, vacuum cleaners, and floor waxers. Cooking Appliances: Like rice cookers, microwaves, ovens, and induction cookers.
780	5. Water Heating Appliances: Including water heaters (gas water heaters, electric water heaters) and heat pumps. 6. Linhting Appliances: Comprising various lamos. FED bulls, and intelligent linhting control systems.
781	7. Personal Care Appliances: Such as electric shavers, hair dryers, beauty instruments, and personal massagers. 8. Small Appliances: This schedown includes a wide variety on the scoffee makers tractors: history blanders can wilk makers air fruers ats
782	9. Entertainment Appliances: Amis dategory includes a wine variety, such as coje e induers, todsters, jances, bienders, soy mink makers, an jiyers, etc.
783	10.5mart Home Appliances: Like smart thermostats, smart sockets, security cameras, etc. With technological advancements, many new household appliances continue to emerge, often featuring smart functions that can be operated via
784	mobile apps or voice control, bringing more convenience and comfort to people's lives. It is important to note that the above list is not exhaustive; different cultures and regions may have specific types of household appliances. Additionally, over time, some older appliances may be replaced by
785	more efficient, energy-saving new products.
786	Please evaluate this answer. Does it respond to the question asked?
787	B. Acceptable answer, but of average quality
788	Evaluation: [A]
789	
790	Question:举例说明计算机体系结构中的缓存写入策略 Answert:在计算机体系结构中,为了确保教理存储的确定性,提高系统性能和减少资源消耗,通常采用以
791	下几种写回策略:1.全盘刷新(Full Write):当处理器需要对某个内存地址进行修改时,整个存储器芯片 合全部国际,注册写同定式可以经证数据的准确性,但可能导致于量时间设施在处理重复信息上一影响多
792	会主部受制。这种与四方认可以体征数据的准确性,但可能导致入量的间质费任实理要复高达工,影响标 统的整体效率。2.局部刷新(Partial Write):对于单位业数几个字节的数据更改,只为其附近的相邻区 MitHYGE基本:可比方的性能性导致,并你可以可能可以可能可以可能。
703	域进行更新。这种方式可以在床行高准确性的向时,卫有时间和空间,碎面切托。3. 里叠式与回(Ovenap Writing):将多个连续的写入操作合并为一次大型的写入操作,从而实现更高效的写入过程。这种方法
794	适用于天空乂忤或者颏繁访问的程序代码寺场景,能够显者提升系统性能。4.异步与回(Asynchronous): 根据不同的需求和环境条件,选择 合适的写回时机, 可以充分利用硬件资源的潜力,避免因等待而导致的额外开销。总之,各种写回策略各有其特点和使用范围,在实际应用
795	甲应结合具体情况进行权衡与优化,以确保系统稳定运行并达到最佳效果。 Answer2: 在计算机系统设计中,为了提高存储器的访问速度和减少处理器等待时间,通常采用的写缓冲(或称*Cache")技术。当一个数据被
796	更新时,首先会将其存入到写缓冲区中,然后由CPU向内存发出指令来将修改后的内容写到内存中去。这样就避免了直接从内存读取新数据的 延迟问题,提高了系统的性能表现。同时,由于写缓冲是位于 CPU 和主存储器之间的中间层,因此可以降低对主存储器和处理器的压力,从
707	而达到优化整个系统资源利用的效果。当然、随着技术的不断发展、现代计算机系统中已经逐渐开始使用更先进、更高效率的写缓冲方案了。 总体来说、通过合理的设计和使用写缓冲器、能够有效提升计算机体系的整体运行效能并实现更好的用户休险效果
700	请从以上两个回复中选出你认为更好的一项 A Annund
700	B Answer2
199	时形. [A] Question: Describe examples of cache write policies in computer architecture
008	Answer 1: In computer architecture, common write-back strategies include: 1. Full Write: Updates the entire memory chip when a memory address is modified ensuring accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy but may be inofficient 3. Partial Write: Updates and a service accuracy bu
801	reducing power consumption. 3. Overlap Writing: Combines multiple write operations into one for efficiency, especially for large files or frequently
802	accessea coae. 4. Asynchronous Write-back: Writes back aata at optimal times based on system conditions, optimizing resource use and minimizing delays. These strategies vary in their features and applications, requiring careful selection for optimal system performance.
803	Answer 2: Computer systems often use write buffer (Cache) technology to speed up memory access and reduce CPU wait times. Data updates are first stored in the write buffer, then written to memory, reducing direct memory access delays. This technique, evolving with newer, more efficient
804	schemes, enhances overall system performance and user experience. Choose the better response:
805	A. Answer 1 B. Answer 2
806	Evaluation: [A]
807	
808	Figure 11: Examples in Human Evaluation of Data and Model Performance.
000	



ModelC-EVAL PPLCMMLU PPLPPLGENPPLGENKun-9k72.9572.9673.5324.53Kun-26k73.73 <u>73.26</u> <u>75.75</u> 25.23Kun-39k <u>73.84</u> 73.38 75.6825.31Kun-52k73.3672.7275.4825.50Kun-100k72.9572.9673.5324.53Kun-200k73.3672.0774.8131.81Kun-360k73.2571.8874.14 46.58 Kun-52k + COIG-one-38k73.61 73.38 75.7034.84COIG-38k 74.09 72.95 75.92 38.40Belle-52k72.3172.5674.73 <u>44.56</u> OL-CC-10k72.0071.6274.9734.38						
Model PPL GEN PPL GEN Kun-9k 72.95 72.96 73.53 24.53 Kun-26k 73.73 73.26 75.75 25.23 Kun-39k 73.84 73.38 75.68 25.31 Kun-52k 73.36 72.72 75.48 25.50 Kun-100k 72.95 72.96 73.53 24.53 Kun-52k 73.36 72.72 75.48 25.50 Kun-100k 72.95 72.96 73.53 24.53 Kun-200k 73.36 72.07 74.81 31.81 Kun-360k 73.25 71.88 74.14 46.58 Kun-52k + 73.61 73.38 75.70 34.84 COIG-0ne-38k 74.09 72.95 75.92 38.40 Belle-52k 72.31 72.56 74.73 44.56 OL-CC-10k 72.00 71.62 74.97 34.38			C-EVAL		CMMLU	
Kun-9k72.9572.9673.5324.53Kun-26k73.7373.2675.7525.23Kun-39k73.8473.3875.6825.31Kun-52k73.3672.7275.4825.50Kun-100k72.9572.9673.5324.53Kun-200k73.3672.0774.8131.81Kun-360k73.2571.8874.1446.58Kun-52k + COIG-one-38k73.6173.3875.7034.84COIG-38k74.0972.9575.9238.40Belle-52k72.3172.5674.7344.56OL-CC-10k72.0071.6274.9734.38		Model	PPL	GEN	PPL	GEN
Kun-26k 73.73 73.26 75.75 25.23 Kun-39k 73.84 73.38 75.68 25.31 Kun-52k 73.36 72.72 75.48 25.50 Kun-100k 72.95 72.96 73.53 24.53 Kun-200k 73.36 72.07 74.81 31.81 Kun-360k 73.25 71.88 74.14 46.58 Kun-52k + COIG-one-38k 73.61 73.38 75.70 34.84 COIG-38k 74.09 72.95 75.92 38.40 Belle-52k 72.31 72.56 74.73 44.56 OL-CC-10k 72.00 71.62 74.97 34.38		Kun-9k	72.95	72.96	73.53	24.53
Kun-39k 73.8473.3875.68 25.31Kun-52k73.3672.7275.4825.50Kun-100k72.9572.9673.5324.53Kun-200k73.3672.0774.8131.81Kun-360k73.2571.8874.14 46.58 Kun-52k + COIG-one-38k73.61 73.38 75.7034.84COIG-38k 74.09 72.95 75.92 38.40Belle-52k72.3172.5674.73 <u>44.56</u> OL-CC-10k72.0071.6274.9734.38		Kun-26k	73.73	73.26	75.75	25.23
Kun-52k73.3672.7275.4825.50Kun-100k72.9572.9673.5324.53Kun-200k73.3672.0774.8131.81Kun-360k73.2571.8874.14 46.58 Kun-52k + COIG-one-38k73.61 73.38 75.7034.84COIG-38k 74.09 72.95 75.92 38.40Belle-52k72.3172.5674.73 <u>44.56</u> OL-CC-10k72.0071.6274.9734.38		Kun-39k	73.84	73.38	75.68	25.31
Kun-100k 72.95 72.96 73.53 24.53 Kun-200k 73.36 72.07 74.81 31.81 Kun-360k 73.25 71.88 74.14 46.58 Kun-52k + COIG-one-38k 73.61 73.38 75.70 34.84 COIG-38k 74.09 72.95 75.92 38.40 Belle-52k 72.31 72.56 74.73 44.56 OL-CC-10k 72.00 71.62 74.97 34.38		Kun-52k	73.36	72.72	75.48	25.50
Kun-200k73.3672.0774.8131.81Kun-360k73.2571.8874.1446.58Kun-52k + COIG-one-38k73.6173.3875.7034.84COIG-38k74.0972.9575.9238.40Belle-52k72.3172.5674.7344.56OL-CC-10k72.0071.6274.9734.38		Kun-100k	72.95	72.96	73.53	24.53
Kun-500k 73.23 71.88 74.14 40.36 Kun-52k + COIG-one-38k 73.61 73.38 75.70 34.84 COIG-38k 74.09 72.95 75.92 38.40 Belle-52k 72.31 72.56 74.73 44.56 OL-CC-10k 72.00 71.62 74.97 34.38		Kun-200k	73.30	71.00	74.81	51.81 46 59
Kun-52k + COIG-one-38k73.61 73.38 75.7034.84COIG-38k Belle-52k 74.09 72.95 75.92 38.40OL-CC-10k72.0071.6274.73 $\frac{44.56}{34.38}$	_	Kun-300k	13.25	/1.88	/4.14	40.58
COIG-38k Belle-52k 74.09 72.95 75.92 38.40OL-CC-10k72.3172.5674.73 $\frac{44.56}{34.38}$		Kun-52k + COIG-one-38k	73.61	73.38	75.70	34.84
COIG-38k 74.09 72.95 75.92 38.40 Belle-52k 72.31 72.56 74.73 <u>44.56</u> OL-CC-10k 72.00 71.62 74.97 34.38		COIC 291-	74.00	72.05	75.00	29.40
Bene-52k 72.51 72.56 74.73 44.56 OL-CC-10k 72.00 71.62 74.97 34.38		COIG-38K	74.09	12.95	13.92	38.40
OL-CC-10K /2.00 /1.62 /4.9/ 34.38		OL CC 10b	72.00	71.60	74.73	$\frac{44.50}{24.29}$
	_	OL-CC-10k	72.00	/1.62	/4.9/	54.38

A.6 STANDARD BENCHMARKS

Table 3: Performance Statistics on Standard Benchmarks. The best results in each section are **Bold**, the second-best results are <u>underlined</u>, while the results of our best model are in <u>Blue</u>. PPL:Perplexity GEN:Generation