

DYNAMICS OF INSTRUCTION TUNING: EACH ABILITY OF LARGE LANGUAGE MODELS HAS ITS OWN GROWTH PACE

Chiyu Song^{1,2} Zhanchao Zhou^{1,2} Jianhao Yan^{1,2} Yuejiao Fei^{1,2} Zhenzhong Lan^{2,*} Yue Zhang^{2,3,*}
¹Zhejiang University ²School of Engineering, Westlake University
³Institute of Advanced Technology, Westlake Institute for Advanced Study
 songchiyu@westlake.edu.cn

ABSTRACT

Instruction tuning is a burgeoning method to elicit the general intelligence of Large Language Models (LLMs). However, the creation of instruction data is still largely heuristic, leading to significant variation in quality and distribution across existing datasets. Experimental conclusions drawn from these datasets are also inconsistent, with some studies emphasizing the importance of scaling instruction numbers, while others argue that a limited number of samples suffice. To better understand data construction guidelines, we deepen our focus from the overall model performance to the growth of each underlying ability, such as creative writing, code generation, and logical reasoning. We systematically investigate the effects of data volume, parameter size, and data construction methods on the development of various abilities, using hundreds of model checkpoints (7b to 33b) fully instruction-tuned on a new collection of over 40k human-curated instruction data. This proposed dataset is stringently quality-controlled and categorized into ten distinct LLM abilities. Our study reveals three primary findings: (i) Despite data volume and parameter scale directly impacting models’ overall performance, some abilities are more responsive to their increases and can be effectively trained using limited data, while some are highly resistant to these changes. (ii) Human-curated data strongly outperforms synthetic data from GPT-4 in efficiency and can constantly enhance model performance with volume increases, but is unachievable with synthetic data. (iii) Instruction data brings powerful cross-ability generalization, with evaluation results on out-of-domain data mirroring the first two observations. Furthermore, we demonstrate how these findings can guide more efficient data constructions, leading to practical performance improvements on public benchmarks.

1 INTRODUCTION

Large Language Models (LLMs) have shown impressive capabilities across diverse tasks (Brown et al., 2020; Touvron et al., 2023; Chowdhery et al., 2022; Almazrouei et al., 2023; Wang et al., 2022a; Wei et al., 2022b; Zhao et al., 2021; Wei et al., 2023; Ivison et al., 2022; Zhang et al., 2023c; Radford et al., 2019), demonstrating their potential for artificial general intelligence (Bubeck et al., 2023). A key contributor to this success is instruction tuning, a process involving supervised fine-tuning of LLMs on instruction-output pairs (Ouyang et al., 2022; Taori et al., 2023; Chiang et al., 2023; Iyer et al., 2022; Zhou et al., 2023).

Despite the recognition that various factors such as the data quantity and distribution impact the performance of instruction tuning (Zhao et al., 2023; Zhang et al., 2023b; Wang et al., 2023), there remains an inconsistent understanding of their specific roles in shaping model capabilities. For instance, while some studies (Wei et al., 2022a; Sanh et al., 2022) argue that scaling data volume is crucial for success, other results (Zhou et al., 2023) suggest a limited number of instructions is sufficient, with models tuned on mere 1k instances outperforming those tuned on other datasets

* Corresponding authors.

that are 10x-100x larger. Intuitively, instruction data fosters a wide range of abilities in LLMs, such as creative writing, code generation, and logical reasoning, each demanding different levels of intelligence. Therefore, we hypothesize that observed inconsistency across existing studies stems from how different abilities, developed by instruction tuning, respond unevenly to alterations in factors like data quantity.

To validate our hypothesis, we systematically investigate the growth of different underlying abilities against data volume, parameter size, and data construction methods. To this end, we employ the LLaMA series models (Touvron et al., 2023) with further pre-training in Chinese, and propose a new Chinese dataset encompassing over 40,000 human-curated instruction instances, covering ten distinct LLM abilities. Each data instance is rigorously revised by annotators to ensure high-quality text and is categorized into an individual ability, enabling us to analyze the impact of specific factors, such as data quantity, on each ability while controlling other variables like data distribution and quality. Our study consists of hundreds of model checkpoints fully instruction-tuned on the proposed datasets, ranging from 7b to 33b parameters.

Our results reveal three primary findings on the dynamics of instruction tuning:

1. Data quantity or parameter size significantly influences overall performance, but each ability develops at different paces during instruction tuning. Abilities such as Creative Writing are more responsive to these factors and can be well-trained with a small amount of data. In contrast, abilities like Ethics show resistance to these changes, suggesting that alternative approaches beyond supervised fine-tuning may be necessary for their development.
2. Regarding data construction methods, synthetic data from GPT-4 falls short in performance for instruction tuning. Compared to our human-curated data, synthetic data is less efficient and fails to consistently enhance model performance by increasing its volume.
3. Instruction data promotes powerful cross-domain generalization, benefiting abilities beyond those included in our dataset. Evaluation results on out-of-domain (OOD) data support the observations made in the first two findings, with different OOD abilities showing distinct growth rates and human-curated data proving more helpful than synthetic data.

Guided by these findings, we adjust the quantity of different ability data and the mixing proportion of synthetic data, achieving practical performance gains on two public benchmarks, CMMLU (Li et al., 2023) and AGIEval (Zhong et al., 2023). We open-source our codebase, dataset, and model checkpoints for reproducibility and future research¹.



Figure 1: Each ability of LLMs has its own growth pace during instruction tuning.

¹<https://github.com/ChiyuSONG/dynamics-of-instruction-tuning>

2 RELATED WORK

Instruction datasets are crucial for the efficacy of instruction-tuned large language models, and their construction methods can be broadly categorized into three types: **Task-formatted datasets** (Sanh et al., 2022; Muennighoff et al., 2022; Wei et al., 2022a; Chung et al., 2022; Mishra et al., 2022; Wang et al., 2022c) incorporate instances from various NLP tasks, including text summarization, natural language inference, sentiment classification, and many other supervision tasks. A sign of these datasets is the inclusion of human-written templates to format each task, such as *"Please summarize the following text:"* These descriptions differentiate task instances so that models can be trained in a multi-task manner. To better scale up the data volume, a crowd-sourcing platform called PromptSourceBach et al. (2022) has also been developed to facilitate the sharing and reviewing of task descriptions. Even though this method can easily enlarge data size in early studies, we focus on methods other than it due to its potential mismatches with actual human requests (Ouyang et al., 2022; Zhao et al., 2023). **Human-curated datasets** (Ouyang et al., 2022; Zhou et al., 2023; Conover et al., 2023; Köpf et al., 2023) address the issue above using genuine user queries, website Q&As, examination questions, and other sources of real-life tasks to construct instructions. Human labelers are assigned to provide output by revising the given answer or directly answering these instructions. Proprietary models like ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI, 2023) employ these procedures to generate training data. Unfortunately, most of these datasets are kept in-house due to the costly and time-consuming nature of the collection process. **Synthetic datasets** (Honovich et al., 2022; Xu et al., 2023a;b) reduce the need for manual collection and annotation of instructions through semi-automated approaches. One approach is collecting user chats with well-performed proprietary models, as in ShareGPT². Another representative approach is Self-Instruct (Wang et al., 2022b), which only requires approximately 100 seed tasks to initiate. It leverages well-performed LLMs to update the task pool as well as generate instruction-output pairs. These explorations effectively synthesize large-scale data at a low budget and encourage open reproductions such as Alpaca (Taori et al., 2023) and Vicuna Chiang et al. (2023).

The influence of dataset construction methods on instruction tuning has been a subject of debate. Existing literature presents a range of findings: Wei et al. (2022a) and Chung et al. (2022) advocate that larger training datasets significantly enhance model performance. In contrast, Zhou et al. (2023) argues that a well-trained model can be achieved with a limited number of human-curated data. Sanh et al. (2022) and Wei et al. (2022a) propose that instruction-tuned models generalize effectively to unseen tasks, whereas Gudibande et al. (2023) suggests that these models only excel in tasks heavily supported in the training dataset. Several studies have also explored various data mix strategies, including balancing data distribution (Longpre et al., 2023), examples-proportional mixing with maximum caps (Raffel et al., 2020; Wei et al., 2022a), and combining multiple datasets (Iyer et al., 2022; Wang et al., 2023).

The discrepancy in these studies has motivated us to investigate the growth pace of various underlying abilities during instruction tuning. In Section 4, we identify significant disparities in the impact of data on different abilities: some abilities show high responsiveness to data increases and can be effectively trained with limited data. Conversely, some abilities are data-hungry, requiring much more data for effective training. This new perspective provides a potential explanation for the differing conclusions drawn from existing research.

3 A NEW COLLECTION OF INSTRUCTION DATA

To systematically investigate the roles of data quantity, parameter size, and data construction methods in shaping a range of model abilities, it is necessary to rule out the influence of data quality and establish a controllable data distribution. Current instruction datasets are insufficient for these research needs. Hence, we introduce a new human-curated Chinese dataset, comprising more than 40,000 instruction-output pairs, each subject to stringent quality control. Moreover, each instance is explicitly categorized into one of ten ability types, enabling us to easily balance or mix them to meet specific experimental proportions. This section will introduce the annotation process for our instruction data.

Following the literature reviewed in Section 2, our human-curated data are derived from real-life scenarios such as academic examinations, online platforms, and user queries. This dataset is

²<https://sharegpt.com/>

Table 1: The data sources, data size, and annotation procedures for each ability category.

Ability	Data Source	Data Size	Annotation Procedure		
			Standardization	Human Filtering	Human Revision
STEM - Biology	COIG - Exam (Zhang et al., 2023a)	1,242	✓	✓	✓
Humanity - History	COIG - Exam (Zhang et al., 2023a)	2,093	✓	✓	✓
Code Generation	Leetcode.cn	5,168	✓	✓	✗
Creative Writing	User Queries from In-House Data	1,200	✓	✓	✓
Chinese	COIG - Exam (Zhang et al., 2023a)	1,650	✓	✓	✓
Dialogue Understanding	C3-D (Sun et al., 2020)	5,085	✓	✓	✗
Role-play Chat	BELLE (Ji et al., 2023)	1,200	✓	✓	✓
Logical Reasoning	LogiQA2.0 (Liu et al., 2023)	12,951	✓	✓	✗
COT for Grad-Math	PRM800K (Lightman et al., 2023)	11,701	✓	✓	✗
Ethics	COIG - Human Value (Zhang et al., 2023a)	1,200	✓	✓	✓

organized into ten representative ability categories: (1) STEM subject - Biology, (2) Humanity subject - History, (3) Code Generation, (4) Creative Writing, (5) Language proficiency - Chinese, (6) Dialogue Understanding, (7) Role-play Chat, (8) Logical Reasoning, (9) Chain of Thought, and (10) Ethics.

Data from diverse sources significantly differ in format, including raw web pages, exam papers, user inputs, and data pre-cleaned by other researchers to different extents. To maintain consistent quality across all instances, we employ a three-stage annotation process:

1. **Standardization:** In this stage, we standardize the raw data of different formats into uniform instruction-output pairs through programs. We customize different rules for each data category to extract valid text and eliminate duplicates. Notably, the data for "Chain of Thought" is sourced from PRM800K (Lightman et al., 2023), the only non-Chinese source, and is translated using the ChatGPT (OpenAI, 2022) API before human review.
2. **Human Filtering:** At this stage, each data instance is reviewed by two independent annotators. They are required to (i) Check the correctness of the text. (ii) Control the diversity of instructions, such as filtering out high-frequency personas in Role-play Chat. (iii) Avoid potential ethical issues in the output, such as biased opinions in Creative Writing. A data instance is marked as "pass" only if both annotators approve it, otherwise marked as "fail." In statistics, the pass rates for different ability categories range from 22.8% to 98.3%, with an inter-annotator agreement (IAA) of 0.77.
3. **Human Revision:** For categories with a small base in quantity or low pass rate, we conduct human revision to ensure sufficient numbers for experiments. Each question is revised or directly answered by an annotator, and then the answer undergoes the same process as in stage 2, with two additional reviewers determining its validity.

All the hired annotators are native Chinese speakers and hold a bachelor's degree or higher. The entire annotation process requires over 1,000 labor hours. To meet the experimental requirements in Section 4, the first round of annotation produces 1,000 training data, 100 validation data, and 100 test data for each ability. We then expand the training set to 40k to compare different construction strategies in Section 5. The data sources, data size, and cleaning procedure for different ability are outlined in Table 1. In Appendix A.1, we present examples of each category.

4 EXPERIMENTS

Employing the human-curated dataset proposed in Section 3, we study the abilities' development in response to alterations in data volume, parameter size, and construction methods. Experiments are conducted under both in-domain and out-of-domain conditions. This section outlines the process of model training, evaluation, and results analysis.

4.1 EXPERIMENT SETUP

For quantity-based experiments, we uniformly sample data d_i of size n from each ability a_i within the ten categories $A = \{a_1, a_2, \dots, a_{10}\}$ in our training set. The samples, combined as $D = \bigcup_{i=1}^{10} \{d_i\}$, are utilized for each model training. We increment the sample size from $n = 1$ logarithmically (base 4) to $n = 1000$ (totaling 10k instances). Regarding parameter sizes, we train models across a full range of 7b, 13b, and 33b scales. To compare different data construction methods, we also examine models trained on synthetic data from GPT-4 (Peng et al., 2023), extending the data volume to 41k instances at 7b and 13b scales. Each training session spans at least 15 epochs, with the corresponding checkpoint saved for evaluation after each epoch. Taking into account all these factor changes, our study requires nearly 500 model checkpoints to draw systematic conclusions.

Our hyperparameter choices are generally in line with Zhou et al. (2023), using AdamW optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.95$ and $weight_decay = 0.1$. The peak learning rate is set to $1e-5$ and linearly decays to $1e-6$ by the end of training. The batch size is 64, with inputs exceeding 2048 tokens trimmed. Referring to the protocol in Song et al. (2023), we also incorporate two speaker tokens, $\langle user \rangle$ and $\langle assistant \rangle$, to segment the utterances from instruction and output for training and inference. The foundation model we use is Chinese-LLaMA (Cui et al., 2023), a series of LLaMA (Touvron et al., 2023) models continuously pre-trained on a 120G Chinese corpus.

4.2 EVALUATION

Selecting the optimal checkpoint for instruction-tuning is non-trivial. Prior studies (Ouyang et al., 2022; Zhou et al., 2023) note that training for more epochs can enhance the model’s capabilities despite the risk of overfitting, and usually employ humans for evaluation. In contrast, automated evaluation is a more scalable solution but has long-lasting concerns about reliability. Despite recent improvements in LLM-based evaluators over traditional metrics (Papineni et al., 2002; Lin, 2004; Banerjee & Lavie, 2005), issues such as false reasoning processes and flawed instruction adherence (Luo et al., 2023; Shen et al., 2023; Chiang & Lee, 2023) persist. Therefore, to efficiently and accurately scale the evaluation across hundreds of checkpoints, we employ a semi-automated approach to reduce the burden on human annotators.

There are two types of questions in our dataset that correspond to distinct evaluation approaches:

- Exact-match questions, such as multiple-choice, true/false, and fill-in-the-blank questions. Each question has only one exclusive gold answer. Similar to other public benchmarks (Hendrycks et al., 2021; Li et al., 2023; Huang et al., 2023; Zhong et al., 2023), we automatically compute the accuracy by comparing generated answers to the ground truth.
- Open-ended questions are common in creative writing, role-play chat, and code generation abilities. These questions do not have standard answers. We thus propose a semi-automated "comparison with distractors" method for scoring them. This method creates distractors (examples shown in Appendix 7 and 8) by manually corrupting each ground truth in two ways: **Fine-grained corruption** involves altering details such as numbers, operators, and terminologies to test the models’ performance in modeling details. **Coarse-grained corruption** creates a distractor that disregards the given instruction but is textually error-free and exhibits the same ability as the gold answer, testing the model’s instruction understanding and adherence. A question scores 1 if the language modeling of ground truth g given the instruction i has a lower perplexity (PPL) than any distractor d_j , otherwise 0:

$$PPL(g|i) = e^{-\sum_{t=1}^T \log p(g_t|i, g_{<t})}, t \text{ denotes the time series of tokens}$$

$$Score = \begin{cases} 1, & \text{if } \min_j (PPL(d_j|i)) > PPL(g|i) \\ 0, & \text{otherwise} \end{cases}$$

As outlined in Sections 3 and 4.1, we train 15 checkpoints for each factor setting and reserve 100 instances each in the validation and test sets for evaluation. We select the highest-scoring checkpoint after the fifth epoch using the validation set and then demonstrate its performance on the test set. Our observations and analysis are discussed in the next subsection.

4.3 RESULTS AND ANALYSIS

We analyze the effect of data volume, parameter size, and construction method. Their impact on overall model performance is illustrated in Fig 2, where the x-axis represents changes in data volume and the y-axis represents the average scores across ten in-domain evaluations plus three out-of-domain abilities. Lines of different colors and symbols represent models with different parameter sizes. We also have a grey dotted line representing the score of random guesses. When scaling the number of training instances, there is a substantial discrepancy on the performance of models trained on human-curated data (depicted by solid lines) and synthetic data (depicted by dashed lines). Moreover, the overall trend is not universally applicable to different abilities when we observe them in the next section.

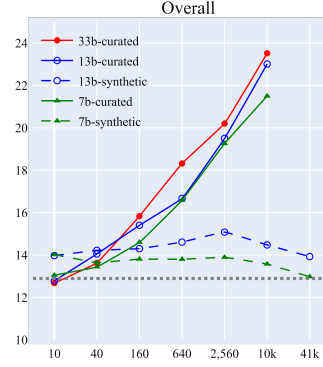


Figure 2: The impact of data volume, parameter scale, and construction method on the overall performance.

4.3.1 COMPARING DIFFERENT ABILITIES

We present the detailed results for each of the abilities in this section, which reveal that different abilities exhibit distinct growth paces when faced with changes in data volume and parameter scale. We have categorized their characteristics into four groups:

Abilities responsive to both factors: As depicted in Fig 3, Code Generation, STEM-Biology, and Humanity-History, which necessitate high professional expertise, show a clear upward trend with the growth of data volume and parameter scale. However, the growth magnitude varies among different abilities, with a more apparent contrast observed in Creative Writing (Fig 5). Most abilities’ growth rate accelerates with the logarithmic increase in data volume, but the slope of the curve for Creative Writing gradually disappears, suggesting a saturation point with limited data volume.

Abilities resistant to parameter size: In Fig 4 part (i), we observe that Dialogue Understanding and Logical Reasoning can still achieve significant ability improvement through data growth. However, the curves of different parameter sizes are intertwined, indicating their insensitivity to this change. Notably, they are the only two understanding tasks in our dataset that do not heavily rely on professional disciplinary knowledge.

Abilities resistant to both factors: As seen in Fig 4 part (ii), Ethics and Role-play chat exhibit stagnant scores across all factor changes. This indicates that supervised fine-tuning (SFT) alone may not effectively enhance these abilities, warranting the investigation of further approaches like reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022; Nakano et al., 2021) and their effects on enhancing diverse abilities.

Ability with special growth characteristics: Fig 5 lists the remaining three abilities that cannot be directly classified. The graphs for Chinese and Creative Writing indicate that the 33b model underperforms the 13b model, a phenomenon also observed in some OOD ability evaluations (Fig 6), for which we lack a satisfactory explanation. Additionally, the Chain-of-thought ability shows only marginal improvements within our experimental scope, likely due to the high difficulty of grad-math questions and the scoring based solely on exact matches with gold answers. Further extending data volume and conducting process-level evaluations (Lightman et al., 2023) may yield further insights. However, we reserve these explorations for future research.

4.3.2 HUMAN-CURATED VS. SYNTHETIC

In studying the influence of various construction methods on ability development, Figures 3, 4, 5, and 6 also present the results from models trained on synthetic data from GPT-4 (Peng et al., 2023). We evaluate both 7b and 13b models, which yield analogous conclusions. Only the 13b results are plotted for simplicity, with the 7b results included in Appendix 9. Comparing the effectiveness of synthetic and human-curated data, it is evident that the abilities taught by synthetic data are limited, hovering around random scores even with an enlarged data size of 41k. Importantly, increasing the quantity of synthetic data does not yield consistent ability growth as observed with human-curated

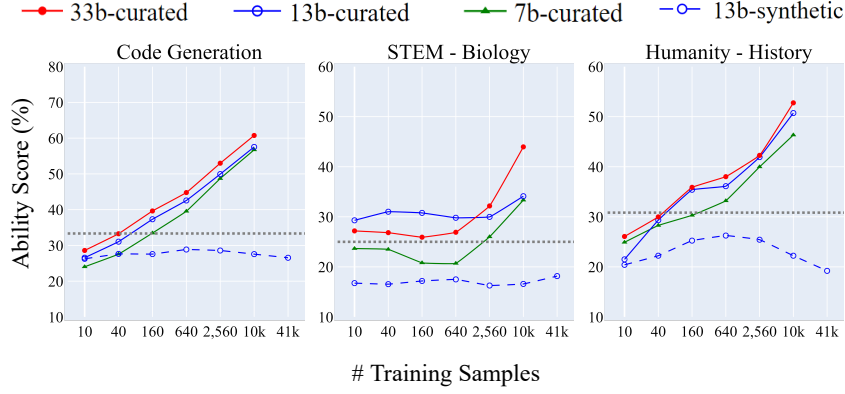


Figure 3: Abilities that are responsive to the data quantity and parameter scale on human-curated data, also comparing the data efficiency of different construction methods with synthetic data.

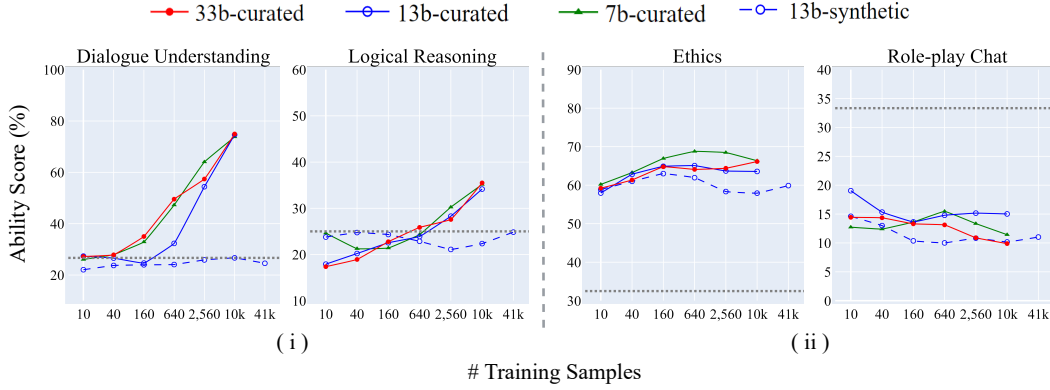


Figure 4: Abilities that are resistant to (i) parameter scale and (ii) both parameter scale and data volume on human-curated training data, also comparing the data efficiency of different construction methods with synthetic data.

data, also indicated by the noticeable inflection point in Figure 2 after the data volume exceeds 2,560. In Section 5, we further explore adjustments to the data mix strategy, with our conclusions confirmed by experimental results on public benchmarks.

4.3.3 OUT-OF-DOMAIN GENERALIZATION

In addition to studying the development of in-domain abilities, we also conduct experiments on three out-of-domain abilities to examine the factors influencing the generalization of instruction tuning. We choose three distinct abilities from the C-Eval datasets (Huang et al., 2023): Teacher Qualification, Physician Qualification, and Urban and Rural Planner. For each ability, we randomly select 40 questions from the validation set, given the limited availability of gold answers in this dataset. We normalize the final scores to a percentage scale to align with the in-domain evaluation setting.

In Figure 6, scores and growth trends of three out-of-distribution (OOD) abilities illustrate that instruction-tuned models exhibit strong cross-ability generalization on unseen data. These results support the hypothesis proposed in previous research (Zhou et al., 2023) that "A model predominantly acquires its knowledge and capabilities during pre-training, while instruction tuning aligns its output to proper formats for user interactions." Our experimental framework effectively validates this hypothesis and quantifies the generalization strength of abilities not included in the instruction data. Moreover, their growth characteristics mirror in-domain abilities: Different abilities react differently to variations in data quantity and parameter scale. Synthetic data is still less efficient than Human-curated data and fails to continuously raise the scores with increasing data volume.

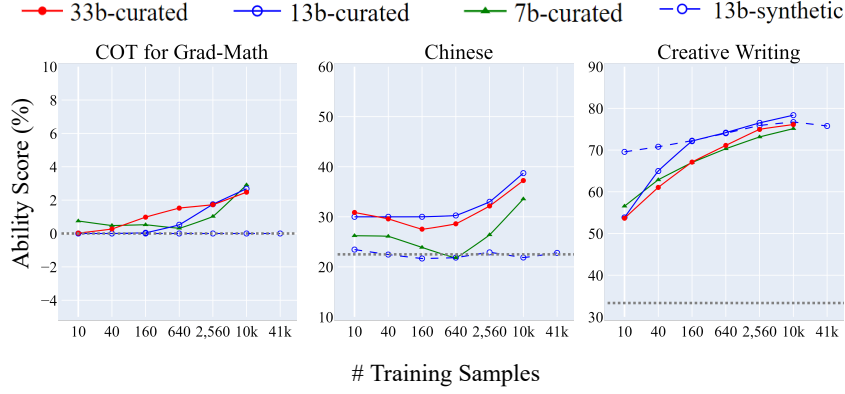


Figure 5: Ability with special growth characteristics on human-curated training data, also comparing the data efficiency of different construction methods with synthetic data.

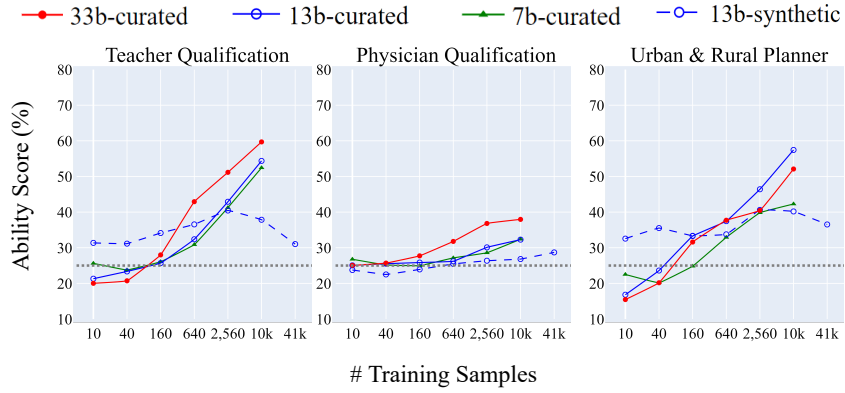


Figure 6: Growth paces of out-of-domain abilities that not included in the human-curated training data, also comparing the data efficiency of different construction methods with synthetic data.

5 GUIDANCE ON DATA MIX STRATEGIES

Building on our understanding of the instruction tuning dynamics, we investigate their applicability in guiding dataset construction. We validate our findings using two public benchmarks, AGIEval (Zhong et al., 2023) and CMMLU (Li et al., 2023). AGIEval (Zhong et al., 2023) is a human-centric benchmark, evaluates the general capabilities of LLMs in tasks related to human cognition and problem-solving. We only focus on the multiple-choice questions within its three Chinese subsets. CMMLU (Li et al., 2023), a comprehensive evaluation benchmark similar to MMLU (Hendrycks et al., 2021), is tailored to assess LLMs’ knowledge and reasoning capabilities within the Chinese language and cultural context, covering a broad spectrum of 67 subjects from elementary to advanced professional levels.

5.1 RECONSTRUCTION OF HUMAN-CURATED DATA

Guided by the findings, “*Abilities react differently to data increase*” and “*Human-curated data yield strong cross-ability generalization*,” we compare three data construction approaches:

Baseline: We use the model trained on 1k instances per ability (totaling 10k instances) from Section 4 as the baseline.

Reconstruct: The distinct growth paces of different abilities in Section 4.3 inspire us to adjust their proportions. For Ethics and Role-play Chat, their scores remain stagnant across all data sizes, so we retain only 64 instances each, which correspond to their relatively higher points on the graph. Considering that although the missing data does not significantly aid their corresponding abilities, it can still benefit other abilities due to cross-ability generalization. We thus keep the training data at 10k

Table 2: Comparing the performance of three construction approaches on two benchmarks, evaluated using checkpoints after epochs 5, 10, 15 with a parameter size of 7b. Scores superior to the baseline are marked with \uparrow .

Models	Data Quantity	AGIEval - 0shot			CMMLU - 0shot		
		ep05	ep10	ep15	ep05	ep10	ep15
Baseline	10k	31.59	34.64	34.76	35.46	36.75	36.34
Reconstruct	10k	35.82 \uparrow	35.43 \uparrow	35.18 \uparrow	35.77 \uparrow	36.85 \uparrow	36.76 \uparrow
Maximum	40k	36.41 \uparrow	37.61 \uparrow	38.85 \uparrow	37.74 \uparrow	37.28 \uparrow	37.53 \uparrow
		AGIEval - 5shot			CMMLU - 5shot		
		ep05	ep10	ep15	ep05	ep10	ep15
Baseline	10k	28.71	31.01	30.27	34.50	35.14	34.29
Reconstruct	10k	33.66 \uparrow	32.27 \uparrow	32.65 \uparrow	35.12 \uparrow	35.89 \uparrow	35.79 \uparrow
Maximum	40k	33.37 \uparrow	33.57 \uparrow	33.35 \uparrow	37.02 \uparrow	37.16 \uparrow	37.13 \uparrow

Table 3: Comparing the performance of three mixing strategies with synthetic data on two benchmarks, evaluated using checkpoints after epochs 5, 10, and 15 with a parameter size of 7b. Highest performance under each setting is in bold.

Models	Data Quantity	AGIEval - 0shot			CMMLU - 0shot		
		ep05	ep10	ep15	ep05	ep10	ep15
Maximum+0	40k+0	36.41	37.61	38.85	37.74	37.28	37.53
Maximum+2.56k	40k+2.56k	37.08	39.21	39.88	37.30	37.74	37.74
Maximum+41k	40k+41k	32.69	34.43	34.38	33.98	36.20	35.34
		AGIEval - 5shot			CMMLU - 5shot		
		ep05	ep10	ep15	ep05	ep10	ep15
Maximum+0	40k+0	33.37	33.57	33.35	37.02	37.16	37.13
Maximum+2.56k	40k+2.56k	34.11	34.07	34.00	36.91	36.87	36.46
Maximum+41k	40k+41k	30.06	31.65	31.41	34.07	35.06	35.17

by uniformly increasing the data volume of other abilities. Specifically, this replenishment excludes Creative Writing, as its score has already saturated at 1,000 instances, so it remains unchanged.

Maximum: We continue to expand our data volume following the same insights. Apart from Ethics and Role-play Chat using only 64 instances and Creative Writing maintaining 1,000, we have expanded the data for other abilities according to the procedures in Section 3, with their specific quantities listed in Table 1. Notably, the expanded dataset is unbalanced in data proportions due to the varying difficulty of the cleaning process for each ability.

We train a 7b model for each construction approach and test their performance at epochs 5, 10, and 15 on two benchmarks under both 0-shot and 5-shot settings. Table 2 marks the results that show improvement over the baseline with \uparrow . Both new constructions demonstrate significant improvements over the baseline. The "Reconstruction" approach achieves an absolute improvement of 1%-4% on AGIEval with unchanged data volume. With the expanded dataset, the "Maximum" approach further shows overall improvement over "Reconstruction" across all abilities.

5.2 MIX UP WITH SYHTNETIC DATA

Considering that synthetic data is a rich open resource, but Section 4.3 indicates that "*Synthetic data does not consistently enhance model performance with an increase in data volume.*" Therefore, it is worth investigating how and whether its integration with human-curated data in practical scenarios can further benefit model performance. We utilize the "Maximum" construction from the previous subsection as our baseline and then integrate varying quantities (0, 2560, 41k) of synthetic data to train 7b models.

Table 3 tests the performance of three mixing strategies at epochs 5, 10, and 15 on two benchmarks under both 0-shot and 5-shot settings. For AGIEval, integrating 2,560 synthetic data instances can lead to performance improvements. For CMMLU, the highest performance (values in bold) occurs when no synthetic data is mixed or only 2,560 instances are mixed. Notably, further expanding the quantity of synthetic data to 41k does not bring additional improvement but instead leads to a significant performance decline, aligning with our previous findings in Section 4.

The outcomes presented in this section confirm our study’s effectiveness in reconstructing the human-curated dataset and its integration with synthetic data. They demonstrate the practical applicability of our previous findings by showing performance enhancement across both benchmarks.

6 CONCLUSION

This research investigates the dynamics of instruction tuning in developing diverse abilities in Large Language Models. We observe that each ability’s development is a unique process, as different abilities exhibit distinct responses to factor changes in data quantity, parameter size, and data construction methods. To systematically study the impact of each factor, we leverage a newly proposed dataset comprising over 40,000 human-curated instructions across ten ability categories. Our findings also reveal that human-curated data consistently outperforms synthetic data from GPT-4 in enhancing model performance. Moreover, we demonstrate that instruction data promotes strong cross-domain generalization, by quantifying its strength on three out-of-domain abilities. Applying these insights to dataset construction, we achieve practical performance improvements on two public benchmarks, CMMLU and AGIEval.

REFERENCES

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. Falcon-40B: an open large language model with state-of-the-art performance. 2023.
- Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesh Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-david, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Fries, Maged Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Dragomir Radev, Mike Tian-jian Jiang, and Alexander Rush. PromptSource: An integrated development environment and repository for natural language prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pp. 93–104, Dublin, Ireland, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-demo.9. URL <https://aclanthology.org/2022.acl-demo.9>.
- Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pp. 65–72, Ann Arbor, Michigan, 2005. Association for Computational Linguistics. URL <https://aclanthology.org/W05-0909>.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html>.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *ArXiv preprint*, abs/2303.12712, 2023. URL <https://arxiv.org/abs/2303.12712>.
- Cheng-Han Chiang and Hung-yi Lee. Can large language models be an alternative to human evaluations? *ArXiv preprint*, abs/2305.01937, 2023. URL <https://arxiv.org/abs/2305.01937>.

-
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, 2023. URL <https://vicuna.lmsys.org>.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *ArXiv preprint*, abs/2204.02311, 2022. URL <https://arxiv.org/abs/2204.02311>.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *ArXiv preprint*, abs/2210.11416, 2022. URL <https://arxiv.org/abs/2210.11416>.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open instruction-tuned llm, 2023.
- Yiming Cui, Ziqing Yang, and Xin Yao. Efficient and effective text encoding for chinese llama and alpaca. *ArXiv preprint*, abs/2304.08177, 2023. URL <https://arxiv.org/abs/2304.08177>.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. The false promise of imitating proprietary llms. *ArXiv preprint*, abs/2305.15717, 2023. URL <https://arxiv.org/abs/2305.15717>.
- 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>.
- Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. Unnatural instructions: Tuning language models with (almost) no human labor. *ArXiv preprint*, abs/2212.09689, 2022. URL <https://arxiv.org/abs/2212.09689>.
- 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 multi-level multi-discipline chinese evaluation suite for foundation models. *ArXiv preprint*, abs/2305.08322, 2023. URL <https://arxiv.org/abs/2305.08322>.
- Hamish Ivison, Akshita Bhagia, Yizhong Wang, Hannaneh Hajishirzi, and Matthew E. Peters. Hint: Hypernetwork instruction tuning for efficient zero-shot generalisation. *ArXiv preprint*, abs/2212.10315, 2022. URL <https://arxiv.org/abs/2212.10315>.
- Srinivasan Iyer, Xi Victoria Lin, Ramakanth Pasunuru, Todor Mihaylov, Daniel Simig, Ping Yu, Kurt Shuster, Tianlu Wang, Qing Liu, Punit Singh Koura, Xian Li, Brian O’Horo, Gabriel Pereyra, Jeff Wang, Christopher Dewan, Asli Celikyilmaz, Luke Zettlemoyer, and Ves Stoyanov. OPT-IML: scaling language model instruction meta learning through the lens of generalization. *ArXiv preprint*, abs/2212.12017, 2022. URL <https://arxiv.org/abs/2212.12017>.

-
- Yunjie Ji, Yong Deng, Yan Gong, Yiping Peng, Qiang Niu, Baochang Ma, and Xiangang Li. Belle: Be everyone’s large language model engine. <https://github.com/LianjiaTech/BELLE>, 2023.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, et al. Openassistant conversations—democratizing large language model alignment. *ArXiv preprint*, abs/2304.07327, 2023. URL <https://arxiv.org/abs/2304.07327>.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. Cmmu: Measuring massive multitask language understanding in chinese. *ArXiv preprint*, abs/2306.09212, 2023. URL <https://arxiv.org/abs/2306.09212>.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *ArXiv preprint*, abs/2305.20050, 2023. URL <https://arxiv.org/abs/2305.20050>.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>.
- Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. *ArXiv preprint*, abs/2301.13688, 2023. URL <https://arxiv.org/abs/2301.13688>.
- Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. *ArXiv preprint*, abs/2303.15621, 2023. URL <https://arxiv.org/abs/2303.15621>.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3470–3487, Dublin, Ireland, 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.244. URL <https://aclanthology.org/2022.acl-long.244>.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. Crosslingual generalization through multitask finetuning. *ArXiv preprint*, abs/2211.01786, 2022. URL <https://arxiv.org/abs/2211.01786>.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *ArXiv preprint*, abs/2112.09332, 2021. URL <https://arxiv.org/abs/2112.09332>.
- OpenAI. Introducing chatgpt. *OpenAI Blog*, 2022.
- OpenAI. GPT-4 Technical Report. *ArXiv preprint*, abs/2303.08774, 2023. URL <https://arxiv.org/abs/2303.08774>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. *ArXiv preprint*, abs/2203.02155, 2022. URL <https://arxiv.org/abs/2203.02155>.

-
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA, 2002. Association for Computational Linguistics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *ArXiv preprint*, abs/2304.03277, 2023. URL <https://arxiv.org/abs/2304.03277>.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 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>.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen H. 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 V. 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 Févry, 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 *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL <https://openreview.net/forum?id=9Vrb9D0WI4>.
- Chenhui Shen, Liying Cheng, Yang You, and Lidong Bing. Are large language models good evaluators for abstractive summarization? *ArXiv preprint*, abs/2305.13091, 2023. URL <https://arxiv.org/abs/2305.13091>.
- Chiyu Song, Hongliang He, Haofei Yu, Pengfei Fang, Leyang Cui, and Zhenzhong Lan. Uni-encoder: A fast and accurate response selection paradigm for generation-based dialogue systems. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6231–6244, Toronto, Canada, 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.388. URL <https://aclanthology.org/2023.findings-acl.388>.
- Kai Sun, Dian Yu, Dong Yu, and Claire Cardie. Investigating prior knowledge for challenging Chinese machine reading comprehension. *Transactions of the Association for Computational Linguistics*, 8:141–155, 2020. doi: 10.1162/tacl.a.00305. URL <https://aclanthology.org/2020.tacl-1.10>.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv preprint*, abs/2302.13971, 2023. URL <https://arxiv.org/abs/2302.13971>.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *ArXiv preprint*, abs/2203.11171, 2022a. URL <https://arxiv.org/abs/2203.11171>.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *ArXiv preprint*, abs/2212.10560, 2022b. URL <https://arxiv.org/abs/2212.10560>.

-
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5085–5109, Abu Dhabi, United Arab Emirates, 2022c. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.340>.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hanna Hajishirzi. How far can camels go? exploring the state of instruction tuning on open resources. *ArXiv preprint*, abs/2306.04751, 2023. URL <https://arxiv.org/abs/2306.04751>.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022a. URL <https://openreview.net/forum?id=gEZrGCozdqR>.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *ArXiv preprint*, abs/2206.07682, 2022b. URL <https://arxiv.org/abs/2206.07682>.
- Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, et al. Zero-shot information extraction via chatting with chatgpt. *ArXiv preprint*, abs/2302.10205, 2023. URL <https://arxiv.org/abs/2302.10205>.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions, 2023a.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *ArXiv preprint*, abs/2304.01196, 2023b. URL <https://arxiv.org/abs/2304.01196>.
- Ge Zhang, Yemin Shi, Ruibo Liu, Ruibin Yuan, Yizhi Li, Siwei Dong, Yu Shu, Zhaoqun Li, Zekun Wang, Chenghua Lin, Wen-Fen Huang, and Jie Fu. Chinese open instruction generalist: A preliminary release. *ArXiv preprint*, abs/2304.07987, 2023a. URL <https://arxiv.org/abs/2304.07987>.
- Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. Instruction tuning for large language models: A survey. *ArXiv preprint*, abs/2308.10792, 2023b. URL <https://arxiv.org/abs/2308.10792>.
- Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. Multi-task instruction tuning of llama for specific scenarios: A preliminary study on writing assistance. *ArXiv preprint*, abs/2305.13225, 2023c. URL <https://arxiv.org/abs/2305.13225>.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *ArXiv preprint*, abs/2303.18223, 2023. URL <https://arxiv.org/abs/2303.18223>.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving few-shot performance of language models. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pp. 12697–12706. PMLR, 2021. URL <http://proceedings.mlr.press/v139/zhao21c.html>.

Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *ArXiv preprint*, abs/2304.06364, 2023. URL <https://arxiv.org/abs/2304.06364>.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *ArXiv preprint*, abs/2305.11206, 2023. URL <https://arxiv.org/abs/2305.11206>.

A APPENDIX

A.1 EXAMPLES OF EACH ABILITY CATEGORY

Ability	Example
Creative writing	<p>Prompt: Write a copywriting for the 'Wangshan Hawthorn'.</p> <p>Response: Folded hopes' branches, Lightly swaying within hawthorn's embrace from within.</p> <p>A gentle touch on the tongue, Tartness tracing the intricate pathways of taste buds, Unfolding and spreading deep into the recesses of lips and teeth.</p> <p>Having savored the bittersweet fruit, Have you ever witnessed the bloom of its origin?</p> <p>Remember every tangy journey, Never forget those times when the flowers were in full bloom.</p> <p>#GoodHopeWater #Wangshan Hawthorn #AdvertisingCopy</p>
Continued on next page	

Ability	Example
Role-play Chat	<p>Prompt:</p> <p>Ali is the CEO of a young startup company, possessing an innovative mindset and a pursuit of excellence. He is actively seeking professional talent to assist him in realizing his business vision.</p> <p>John, on the other hand, is an experienced HR manager with profound expertise and extensive experience in talent acquisition.</p> <p>Generate a short conversation happening between them.</p> <p>Response:</p> <p>Ali: Hey John, thank you for taking the time to chat with me. I'm really interested in learning more about your career experience, and I'm quite impressed by your expertise and experience in HR.</p> <p>John: Thank you for the invitation, Ali. I'm very interested in your company and its business vision. Can you share your perspective on talent recruitment and development within the company?</p> <p>Ali: In the early stages of our startup, we're in great need of visionary and passionate talents. We're looking for individuals who can contribute to the company's growth and help us achieve our vision together. We'd like to better understand how to attract, select, and nurture these talents to ensure the ongoing success of our company.</p> <p>John: Your thoughts align well with our priorities. I believe attracting high-caliber talent is a crucial mission for any startup. There are various ways our company can achieve these goals, such as through recruitment, building reputation, and talent development. However, I think the most critical aspect is having a tailored talent management and development plan to ensure that individuals who fit our company's needs can grow within the organization and continually enhance their skills and knowledge.</p> <p>Ali: That sounds excellent. I think we can delve deeper into this project.</p>
Ethics	<p>Prompt:</p> <p>Some people admire Western culture and believe that Chinese culture is conservative and outdated. How should you respond to them? Drawing on your knowledge, share your perspective.</p> <p>Response:</p> <ol style="list-style-type: none"> 1. Different cultures have different values, and it's not appropriate to judge other cultures solely by Western standards. 2. Despite being relatively conservative, Chinese culture possesses unique charm and profound depth. 3. Rather than excluding other cultures, it's important to maintain an open-minded attitude and learn from, as well as incorporate, the beneficial aspects of other cultures. 4. One should avoid unilaterally idealizing any single culture and, instead, choose and integrate various cultures based on practical circumstances.

Continued on next page

Ability	Example
Chinese	<p>Prompt: Reading the following poem and answering the question,</p> <p>Moon Song of Mount Emei by Li Bai On Emei Mountain, a crescent moon graces autumn's sky, Its reflection dances upon the tranquil water of the Pingqiang River nearby. By night, I set forth on the clear stream towards the Three Gorges afar, Thoughts of you linger during my journey, for you are not in sight in Yuzhou.</p> <p>Please describe the scene painted by the lines "On Emei Mountain, a crescent moon graces autumn's sky, Its reflection dances upon the tranquil water of the Pingqiang River nearby" in vivid words.</p> <p>Response: Lofty Emei Mountain stands tall and grand, with a crescent autumn moon hanging amidst its peaks, resembling the mountain itself spilling the moonlight. The pristine moonbeam is mirrored on the surface of the tranquil Pingqiang River, which flows gently, carrying the moon's reflection along its course.</p>
Humanity - History	<p>Prompt: True or False: The three major battles of Liaoshen, Huaihai, and Crossing the Yangtze River essentially annihilated the main forces of the Nationalist Army, accelerating the overall victory of the Liberation War nationwide. ()</p> <p>Response: False</p>
STEM - Biology	<p>Prompt: Please choose the appropriate option based on the content of the question. Regarding enzymes, which statement is correct? A. Enzymes provide the necessary activation energy for the reaction. B. Changes in enzyme activity are unrelated to changes in the enzyme's environment. C. Changes in enzyme structure can result in the loss of its active site or its entire activity. D. Enzyme molecules are immediately degraded into amino acids after catalyzing a reaction.</p> <p>Response: C</p>
Continued on next page	

Ability	Example
COT for Grad-Math	<p>Prompt: Solve the inverse matrix of</p> $\begin{pmatrix} 2 & 3 \\ -1 & 7 \end{pmatrix}$ <p>. Output a zero matrix if the inverse matrix doesn't exist.</p> <p>Response: To solve the inverse matrix of a 2×2 matrix, we can use the following formula</p> $\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix},$ <p>and the condition holds as long as its determinant $ad - bc$ is non-zero.</p> <p>In this example, $2 \times 7 - 3 \times (-1) = 14 + 3 = 17$, which is not zero, thus there exists an inverse matrix.</p> <p>Apply the formula, and we will have:</p> $\begin{pmatrix} 2 & 3 \\ -1 & 7 \end{pmatrix}^{-1} = \frac{1}{17} \begin{pmatrix} 7 & -3 \\ 1 & 2 \end{pmatrix}.$ <p>The final answer is</p> $\frac{1}{17} \begin{pmatrix} 7 & -3 \\ 1 & 2 \end{pmatrix}.$
Logical Reasoning	<p>Prompt: DM advertising, also known as direct mail advertising, involves delivering promotional materials to consumers through mailing, gifting, or other means, directly to their homes or places of business.</p> <p>According to the definition provided, which of the following is not considered as DM advertising?</p> <p>A. Samples of sunscreen lotion given away at cosmetics counters in department stores.</p> <p>B. Flyers promoting a commercial real estate project distributed on crowded streets.</p> <p>C. Coupons for a large chain of fast-food restaurants received in an email inbox.</p> <p>D. Two screen protectors received as a gift when purchasing a mobile phone at an electronics store.</p> <p>Response: D</p>
Continued on next page	

Ability	Example
Dialogue Understanding	<p>Prompt:</p> <p>Woman: "You're back. Did you buy the vinegar I asked for?"</p> <p>Man: "You're the boss at home, so I definitely have to complete the tasks you give me."</p> <p>Woman: "Great attitude! I've decided to make your favorite dish for dinner."</p> <p>Man: "Haha, it must be sweet and sour fish, right?"</p> <p>Woman: "Exactly. You go take a shower, and dinner will be ready in half an hour."</p> <p>What is the woman asking the man to do now?</p> <p>A. Buy vinegar B. Have dinner C. Take a shower D. Work</p> <p>Response:</p> <p>C</p>
Code Generation	<p>Prompt:</p> <p>You are given 2 integers left and right, indicating interval [left, right], return the bitwise AND results of all numbers within this interval (including the left and right endpoints).</p> <p>Example 1: Input: left = 5, right = 7 Output: 4</p> <p>Example 2: Input: left = 0, right = 0 Output: 0</p> <p>Example 3: Input: left = 1, right = 2147483647 Output: 0</p> <p>Hint: 0 <= left <= right <= 231 - 1 Please use c# to finish the problem.</p> <p>Response:</p> <p>The following are reference implementation options:</p> <pre> ```cs public class Solution { public int RangeBitwiseAnd(int left, int right) { while (left < right) { right &= (right - 1); } return right; } } ``` </pre>

A: A young music enthusiast who is passionate about listening to music by actors and female singers in the pop genre.

Generate a snippet of their conversation.

Write a news article based on the given information.

Figure 7: Examples of distractors derived from manually corrupted Role-play Chat and Creative-writing output.

Prompt:

Given a binary array "nums", you need to remove one element from it. Return the length of the longest non-empty subarray containing only "1" in the resulting array. If there is no such subarray, return "0". [folded examples] Please use python3 to finish it.

Ground Truth:

The following is a reference implementation for your consideration:

```
```python
class Solution:
 def longestSubarray(self, nums: List[int]) -> int:
 n = len(nums)
 left = [0] * n
 right = [0] * n
 for i in range(1, n):
 if nums[i - 1] == 1:
 left[i] = left[i - 1] + 1
 for i in range(n - 2, -1, -1):
 if nums[i + 1] == 1:
 right[i] = right[i + 1] + 1
 return max(a + b for a, b in zip(left, right))
```
```

Fine-grained corruption:

The following is a reference implementation for your consideration:

```
```python
class Solution:
 def longestSubarray(self, nums: List[int]) -> int:
 n = len(nums)
 left = [0] * n
 right = [0] * n
 for i in range(1, n):
 if nums[i - 1] == 1:
 left[i] = left[i - 1] + 1
 for i in range(-1, n-2, 1):
 if nums[i + 1] == 1:
 right[i] = right[i + 1] + 1
 return max(a + b for a, b in zip(left, right))
```
```

Coarse-grained corruption:

The following is a reference implementation for your consideration:

```
```python
def find_max_sum(nums):
 max_sum = float('-inf')
 curr_sum = 0
 for num in nums:
 curr_sum += num
 if curr_sum > max_sum:
 max_sum = curr_sum
 if curr_sum < 0:
 curr_sum = 0
 return max_sum

def find_longest_subarray(nums):
 n = len(nums)
 max_len = 0
 for i in range(n):
 for j in range(i, n):
 subarray = nums[i:j+1]
 if find_max_sum(subarray) == len(subarray):
 max_len = max(max_len, len(subarray))
 return max_len
```
```

Figure 8: Examples of distractors derived from manually corrupted Code Generation output.

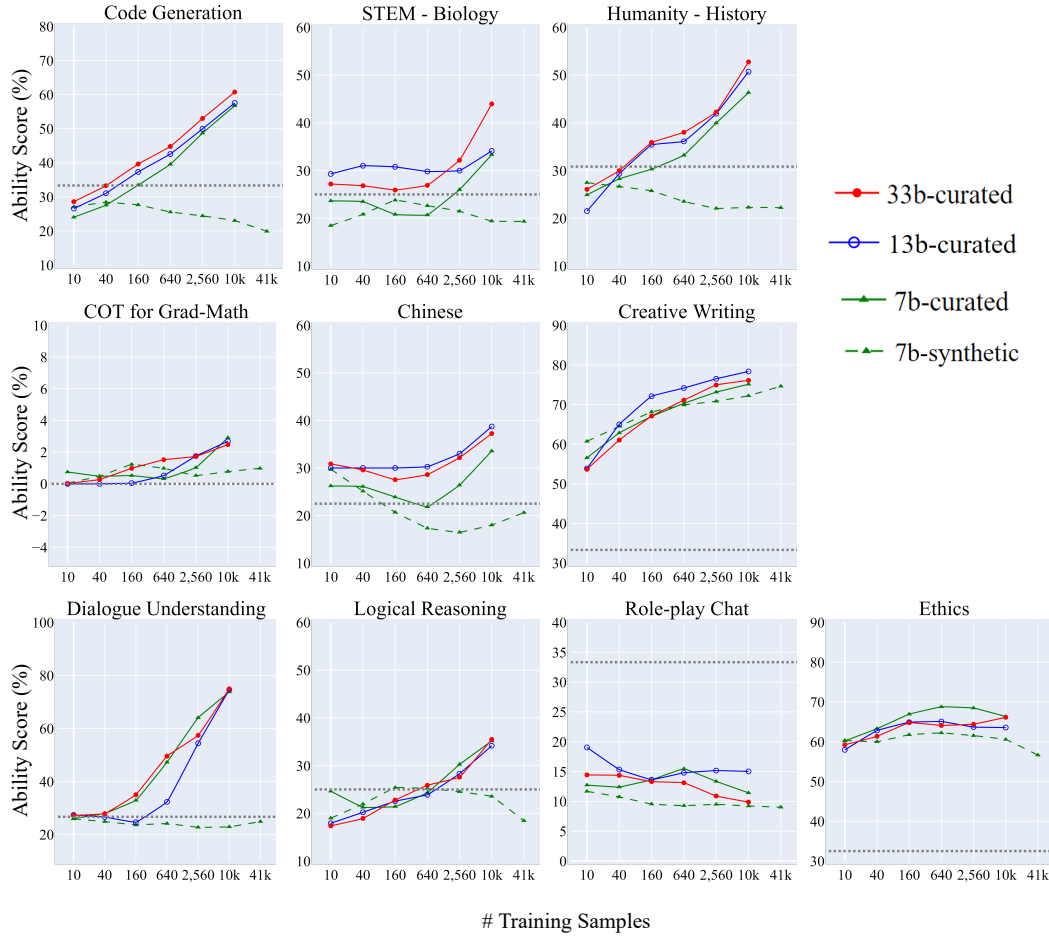


Figure 9: Evaluations of 7b models trained on synthetic data, yielding analogous conclusions as 13b models.