

Dynamics of Instruction Tuning: Each Ability of Large Language Models Has Its Own Growth Pace

Anonymous ACL submission

Abstract

Instruction tuning is a burgeoning method to elicit the general intelligence of Large Language Models (LLMs). However, the understanding of its scaling properties remains underexplored. While some research advocates for expanding the number of instructions, others suggest that a small set of well-chosen examples is adequate. To understand such discrepancy, our work systematically studies the effectiveness of data volume, parameter size, and data construction methods on the development of each underlying ability of LLM, such as creative writing, code generation, and logical reasoning. Our study reveals three primary findings: (i) Despite these factors significantly influencing overall model performance, some abilities are more responsive to scaling, while others show high resistance. (ii) The sensitivity of different abilities to these factors can be explained by their Complexity and Transference, which indicate the relative importance of each factor in learning specific abilities. (iii) Tailoring data construction based on these sensitivities results in performance gains on two public benchmarks. Additionally, we curate a comprehensive dataset containing over 40k instances across ten abilities for our experiments.

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, distribution, and construction method, directly 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., 2021) argue that scaling data volume is crucial for the success of some tasks, other results (Zhou et al., 2023) suggest a limited number of instructions is sufficient. This discrepancy indicates that different abilities of LLMs may respond unevenly to changes in these factors, yet comparative analyses are lacking (Hestness et al., 2017; Zhang et al., 2024). Moreover, current investigation and understanding of instruction tuning are predominantly based on English datasets, with explorations in other languages remaining limited.

To bridge this research gap, we systematically examine the development of various underlying capabilities in relation to data volume, parameter size, and data construction methods. For this purpose, we employ pre-trained Chinese models such as Chinese-LLaMA(Cui et al., 2023), Baichuan2(Yang et al., 2023), and Qwen1.5(Bai et al., 2023). Additionally, we propose *DoIT*, a new Chinese dataset encompassing over 40,000 human-curated instruction instances that span ten distinct LLM abilities. Each data instance is rigorously revised by annotators to ensure high-quality text and is categorized according to its specific ability.

Our analysis disentangles the effects of each factor by maintaining control over others, resulting in a comprehensive set of instruction-tuned models with sizes ranging from 7 billion to 33 billion parameters. The results reveal three primary findings:

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

082 are more responsive to these factors and can 129
083 be well-trained with a small amount of data. 130
084 In contrast, abilities like Ethics show resis- 131
085 tance to these changes, suggesting that alterna- 132
086 tive approaches beyond supervised fine-tuning 133
087 may be necessary for their development. 134

088 2. We investigate the reasons behind these dis- 136
089 crepancies and identify two features, Com- 137
090 plexity and Transference, which can be cal- 138
091 culated in low-resource scenarios. These fea- 139
092 tures help infer the potential for ability growth 140
093 when scaling up data or model parameters. 141

094 3. By adjusting the mixing strategies of differ- 142
095 ent ability data based on their sensitivities 143
096 to scaling, we achieve practical performance 144
097 improvements on two comprehensive bench- 145
098 marks, CMMLU (Li et al., 2023) and AGIEval 146
099 (Zhong et al., 2023). 147

100 We open-source our codebase, dataset, and 148
101 model checkpoints for reproducibility and future 149
102 research¹. 150

103 2 Related Work 151

104 Instruction datasets are crucial for the efficacy of 152
105 instruction-tuned LLMs, and their construction 153
106 methods can be broadly categorized into three 154
107 types: Task-formatted datasets (Sanh et al., 2021; 155
108 Muennighoff et al., 2022; Wei et al., 2022a; Chung 156
109 et al., 2022; Mishra et al., 2021; Wang et al., 2022c) 157
110 incorporate instances from diverse NLP tasks us- 158
111 ing human-crafted templates to enable multi-task 159
112 training. While platforms like PromptSource(Bach 160
113 et al., 2022) have been developed to expand these 161
114 datasets, concerns about their alignment with real 162
115 user requests (Ouyang et al., 2022; Zhao et al., 163
116 2023) have led to the exploration of alternative 164
117 methods. Human-curated datasets (Ouyang et al., 165
118 2022; Zhou et al., 2023; Conover et al., 2023; Köpf 166
119 et al., 2023) address the issue above using real- 167
120 life tasks with human labeling, such as genuine 168
121 user queries or examination questions. Proprietary 169
122 models like ChatGPT (OpenAI, 2022) and GPT-4 170
123 (OpenAI, 2023) employ this data source for train- 171
124 ing. Unfortunately, these datasets are often not 172
125 publicly available due to the high cost and effort 173
126 required. Synthetic datasets (Honovich et al., 2022; 174
127 Xu et al., 2023a,b) offer a cost-effective solution 175
128 by semi-automating instruction generation. One 176
177

approach is collecting user chats with proprietary 129
models as in ShareGPT². Self-Instruct (Wang et al., 130
2022b) is another representative approach, which 131
bootstraps datasets from a small set of seed tasks. 132
This approach has inspired open-source projects 133
like Alpaca (Taori et al., 2023) and Vicuna (Chiang 134
et al., 2023). 135

The influence of data factors on instruction tun- 136
ing has been a subject of debate. Some studies 137
(Wei et al., 2022a; Chung et al., 2022) suggest that 138
larger datasets improve model performance, while 139
others (Zhou et al., 2023) show that a smaller, high- 140
quality dataset can suffice. While there is evidence 141
that instruction-tuned models generalize well (Sanh 142
et al., 2021; Wei et al., 2022a), some argue this 143
is limited to tasks heavily supported in the train- 144
ing data (Gudibande et al., 2023). Synthetic data 145
has shown promise (Wang et al., 2022b; Yin et al., 146
2023), but the model’s capability proved limited 147
by imitating proprietary systems (Gudibande et al., 148
2023). 149

The discrepancy in these studies has motivated 150
us to investigate how various abilities develop dur- 151
ing instruction tuning. Our research, detailed in 152
Section 4, identifies significant disparities in the 153
data impact on different abilities. This insight may 154
reconcile the differing conclusions drawn from 155
prior studies. 156

157 3 DoIT: A New Instruction Dataset 157

To systematically investigate the roles of data quan- 158
tity, parameter size, and data construction methods 159
in shaping a range of model abilities, it is neces- 160
sary to rule out the influence of data quality and 161
establish a controllable data distribution among dif- 162
ferent abilities. To fulfill these research needs, we 163
introduce *DoIT*, a new human-curated dataset. This 164
dataset contains over 40,000 quality-controlled Chi- 165
nese instances, categorized into ten distinct ability 166
classes, allowing for tailored experimental setups. 167

Following the literature reviewed in Section 2, 168
our human-curated data are derived from real-life 169
contexts such as academic examinations, online 170
platforms, and user queries. This dataset is orga- 171
nized into ten representative ability categories: (1) 172
STEM subject - Biology, (2) Humanity subject - 173
History, (3) Code Generation, (4) Creative Writing, 174
(5) Language proficiency - Chinese, (6) Dialogue 175
Understanding, (7) Role-play Chat, (8) Logical 176
Reasoning, (9) Chain of Thought, and (10) Ethics. 177

¹The link to be added.

²<https://sharegpt.com/>

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

Table 1: The data sources and data size after two rounds of human annotation for each ability category.

To maintain consistent quality across all instances, we employ a three-stage annotation process:

- 1. Standardization:** 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. In this stage, we convert them into consistent instruction-output pairs, applying tailored rules for each category to extract relevant text and eliminate duplicates. Notably, the "Chain of Thought" data originated from PRM800K (Lightman et al., 2023) is the only non-Chinese source and is translated using the ChatGPT (OpenAI, 2022) API before human review.
- 2. Human Filtering:** Each item is then 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. Only items approved by both annotators are accepted, with pass rates ranging from 22.8% to 98.3% across different categories and an inter-annotator agreement (IAA) of 0.77.
- 3. Human Revision:** To ensure adequate data for underrepresented or low-approval categories, we conduct human revision to ensure sufficient numbers for experiments. In this stage, each question is revised or answered by an annotator. 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, hold a bachelor’s degree or higher, and dedicate

over 1,000 labor hours to annotation. 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 and sizes for each ability category are outlined in Table 1, with examples provided in Appendix A.2.

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. Each training session spans at least 15 epochs, with the corresponding checkpoint saved for evaluation after each epoch. We also compare our human-curated dataset, *DoIT*, with a synthetic dataset proposed by Peng et al. (2023) for instruction tuning. The synthetic dataset utilizes the Alpaca (Taori et al., 2023) instruction pool, created through the Self-Instruct (Wang et al., 2022b) framework, with responses generated by GPT-4 (OpenAI, 2023). By leveraging the cost-effectiveness of synthetic data to acquire a large and diverse set of instances, we can expand our experimental data volume to 41k on this dataset.

Taking into account all these factors, our study requires nearly 500 model checkpoints to draw systematic conclusions. To ensure the generalizability of our findings, we first analyze the scaling properties of different capabilities using the Chinese-LLaMA model (Cui et al., 2023), which maintains the straightforward architecture of LLaMA (Touvron et al., 2023) without any modifications. Subsequently, we employ more sophisticated foundation models such as Qwen1.5 (Bai et al., 2023) and

Baichuan2 (Yang et al., 2023) to further validate how our insights can enhance model performance. Detailed hyperparameter choices and training procedures are provided in Appendix A.1.

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 in both statistical (Papineni et al., 2002; Lin, 2004; Banerjee and Lavie, 2005) and LLM-based (Luo et al., 2023; Shen et al., 2023; Chiang and Lee, 2023) metrics. 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 (e.g., multiple-choice, true/false, fill-in-the-blank) have one exclusive gold answer. Similar to other public benchmarks (Hendrycks et al., 2020; 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, common in creative writing, role-play chat, and code generation abilities, lack standard answers. We introduce a semi-automated *comparison with distractors* method for these. This method creates distractors (examples shown in Appendix 9 and 10) by manually corrupting each ground truth in two ways: **Fine-grained corruption** subtly alter some 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, 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:

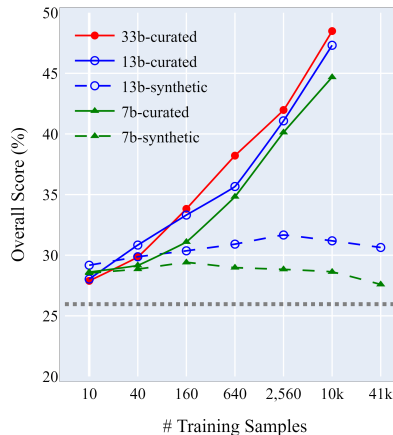


Figure 1: The impact of data volume, parameter scale, and construction method on the overall performance.

$$PPL(g|i) = e^{-\sum_{t=1}^T \log p(g_t|i, g_{<t})} \quad 312$$

, t denotes the time series of tokens 313

$$Score = \begin{cases} 1, & \text{if } \min_j (PPL(d_j|i)) > PPL(g|i) \\ 0, & \text{otherwise} \end{cases} \quad 314$$

As outlined in Sections 3 and 4.1, we train 15 checkpoints for each factor setting and reserve 100 instances in both validation and test sets for evaluation. We select the highest-scoring checkpoint after 5 epochs using the validation set and then demonstrate its performance on the test set. Our observations and analysis are discussed in the next subsection. 315-322

4.3 Results and Analysis 323

We analyze the effect of data volume, parameter size, and construction method. Their impact on overall model performance is illustrated in Fig 1, 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). 324-337

Moreover, the overall trend is not universally applicable to different abilities when we observe them in the next section. Subsequently, we quantify the scaling sensitivity of each ability by investigating the relationship between its task accuracy and the 338-342

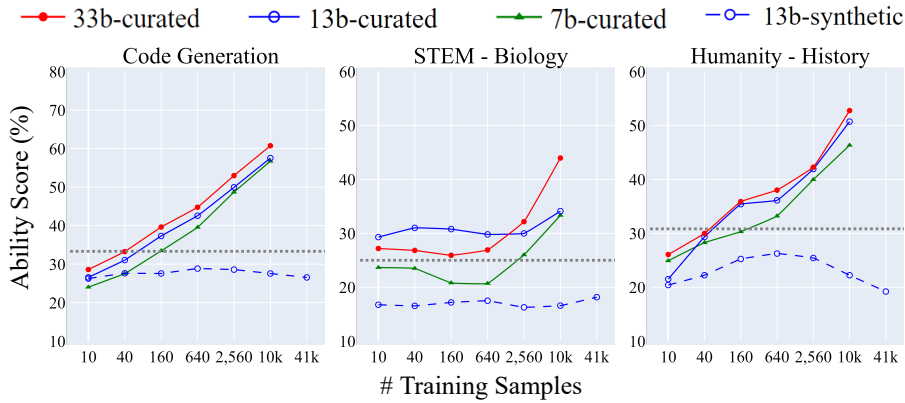


Figure 2: Abilities that are responsive to the data quantity and parameter scale in the human-curated dataset, also comparing the data efficiency of different construction methods with synthetic data.

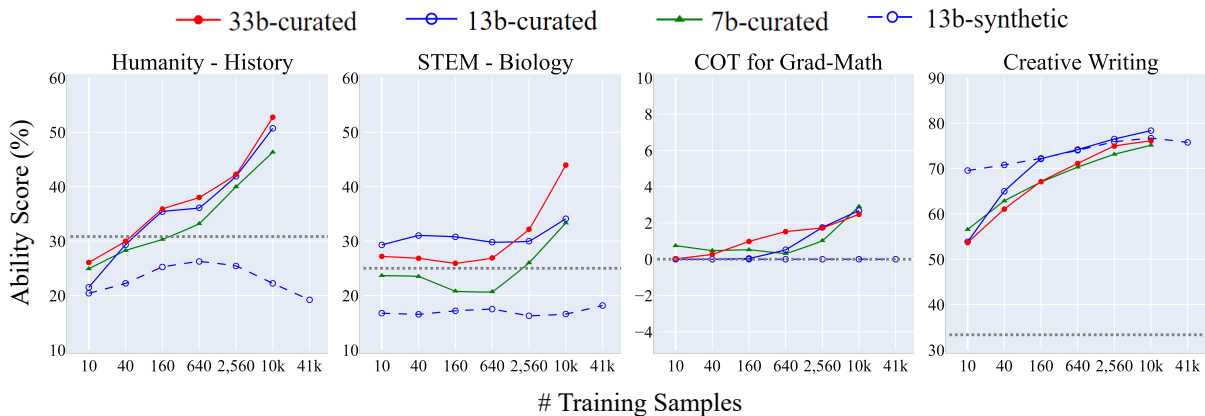


Figure 3: Comparison of abilities with varying sensitivities to data scaling in the human-curated dataset, also comparing the data efficiency of different construction methods with synthetic data.

factors above. We further analyze two interpretable features that potentially forge different scaling sensitivities across these abilities.

4.3.1 Disparities in ability growth trajectories

We present the empirical results for each ability in this section, exhibiting their distinct growth paces when facing factor changes:

Abilities responsive to scaling: Some abilities such as Code Generation, STEM-Biology, and Humanity-History are responsive to factor changes. As illustrated in Fig 2, they show clear upward trends with the growth of data volume and parameter scale.

Varying sensitivities to data scaling: As depicted in Fig 3, the rate of improvement is not uniform across abilities. This figure reveals varying degrees of data scaling sensitivity, with Creative Writing being a notable case. The slope of its curve gradually disappears, indicating a plateau with limited data expansion.

Varying sensitivities to parameter scaling: Fig 4 demonstrates that the impact of parameter size

scaling also varies among abilities. From left to right in the figure, their curves for different model sizes become increasingly intertwined, indicating the insensitivity to this change.

Abilities resistant to scaling: As shown in Fig 5, certain abilities like Ethics and Role-play Chat appear to resist both factors and maintain stagnant scores across all changes. This lack of progress implies that supervised fine-tuning (SFT) alone may not effectively advance these abilities, warranting the investigation of approaches beyond it, such as reinforcement learning from human feedback (RLHF) (Ouyang et al., 2022; Nakano et al., 2021).

Out-of-domain evaluation: Beyond in-domain abilities, Figure 6 evaluates model performance on three out-of-distribution (OOD) tasks from the C-Eval datasets (Huang et al., 2023): Teacher Qualification, Physician Qualification, and Urban and Rural Planning. The observed growth trends suggest robust cross-ability generalization. Similar to in-domain evaluations, these OOD tasks exhibit diverse responses to variations in data quantity and parameter scale.

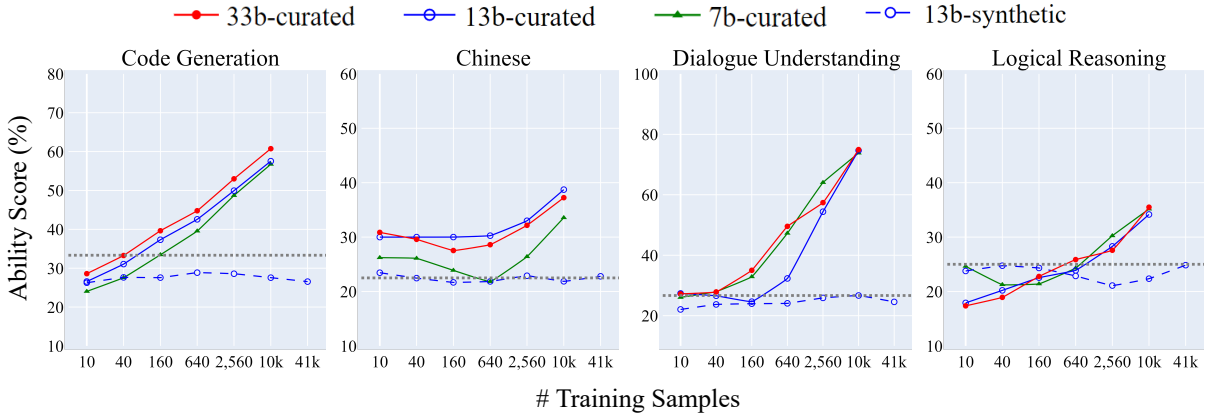


Figure 4: Comparison of abilities with varying sensitivities to parameter size scaling in the human-curated dataset, also comparing the data efficiency of different construction methods with synthetic data.

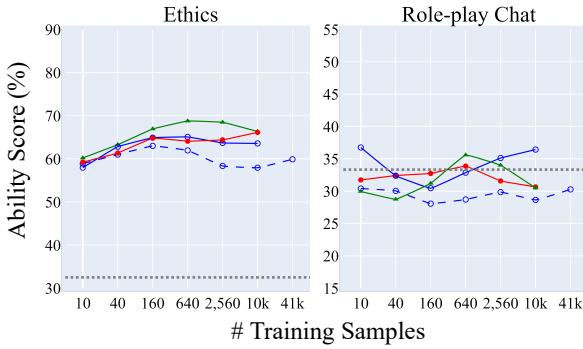


Figure 5: Abilities that are resistant to the data quantity and parameter scale in the human-curated dataset, also comparing the data efficiency of different construction methods with synthetic data.

Human-curated vs. Synthetic: Figures 2-6 also present the results from models trained on synthetic data generated by GPT-4 (Peng et al., 2023). We evaluate both 7B and 13B models, which yield analogous conclusions. For simplicity, only the results for the 13B model are plotted, with the 7B results included in Appendix 11. Our findings align with previous studies (Gudibande et al., 2023), indicating that synthetic data is only effective for learning partial abilities. Additionally, Figure 1 demonstrates that increasing the volume of synthetic data does not continuously improve model performance. We further empirically demonstrate in Appendix A.3 that even when synthetic data is combined with human-curated data, its effectiveness still has an upper limit. Consequently, in subsequent experiments, we focus solely on exploring the scaling properties of human-curated data.

4.3.2 Understanding Diverse Scaling Behaviors

To understand the varying scaling properties of abilities, a notable observation from Section 4.3.1 is

that abilities tied to professional (academic) knowledge are more sensitive to parameter scaling. We define such common feature as **Complexity**, indicating these abilities are inherently "more challenging for language modeling" and "benefit less from the training of other abilities". We hypothesize that Complexity is associated with how different abilities respond to model size changes.

To examine the relationship between Complexity and parameter sensitivity, we first quantify the sensitivities of individual abilities. Adopting the scaling law function similar to Kaplan et al. (2020), we model the task score ACC (averaged across varying data volume) as a function of model size N for each ability i :

$$\exp(ACC_i) = (\exp(c_i) \cdot N)^{\alpha_i}, c_i \text{ is constant} \quad (1)$$

Here, the exponent α_i represents the rate of accuracy improvement with increasing model size, indicating the *scaling sensitivity*.

$$ACC_i = \alpha_i \cdot \log(N) + c_i \quad (2)$$

We further demonstrate that Complexity can be measured in a low-resource setting by fine-tuning separate 7b models with only 64 data points per ability. According to its definition, Complexity is calculated as a weighted sum of the test loss (L) for each model trained individually on ability i , along with its accuracy achieved by training on ability j data, $Acc(j, i)$, improved over the foundation model's performance $Acc(f, i)$:

$$\text{Complexity}_i = w_1 \cdot L_i - w_2 \cdot \sum_{j \neq i} (Acc(j, i) - Acc(f, i))$$

Results depicted in Figure 7 (i) show a clear linear relationship between Complexity and the sensitivity to parameter scale, indicating that even with minimal resources (7b model size, 64 data

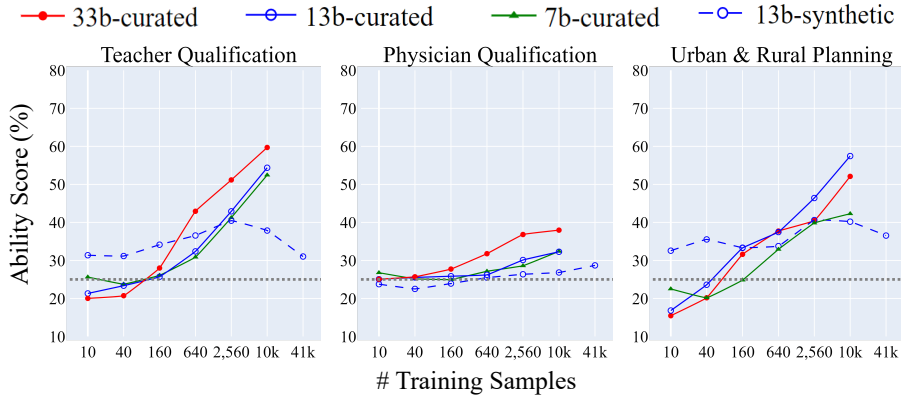


Figure 6: Growth paces of out-of-domain abilities that not included in the human-curated dataset, also comparing the data efficiency of different construction methods with synthetic data.

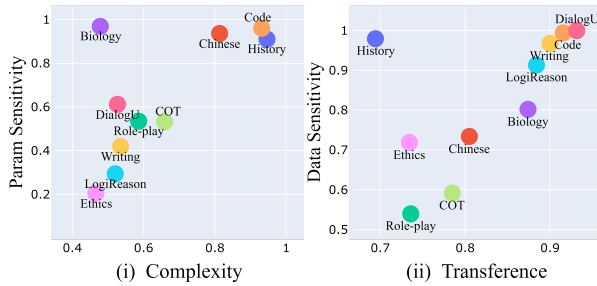


Figure 7: Two interpretable features, (i) Complexity and (ii) Transference, of different abilities demonstrate linear relationships with their sensitivities to scaling in parameter size and data volume. These features can infer the growth of abilities after scale-up, as discussed in Section 5. Sensitivity values are normalized to (0, 1) for visualization.

points), we can forecast how abilities will develop with increased model size.

Correspondingly, we have also computed the **Transference** for each ability, which reflects "how well ability i data enhances other abilities j " via this formula:

$$\text{Transference}_i = w \cdot \sum_{j \neq i} (\text{Acc}(i, j) - \text{Acc}(j, j))$$

By substituting model size (N) with data volume (D) in Equation (1), we can also evaluate each ability's sensitivity to data volume. Figure 7 (ii) shows that Transference is linearly related to data scaling sensitivity. This confirms that using checkpoints trained with limited data can also infer the impact of enlarged data volume on ability development.

5 Guidance on Model Training

After understanding that "abilities react differently to factor changes" and learning "how to estimate the sensitivities of different abilities to scaling," we further investigate how to leverage their varying sensitivities to enhance specific or overall model abilities more effectively. At this stage,

we employ two more advanced foundation models, Qwen1.5(Bai et al., 2023) and Baichuan2(Yang et al., 2023), to conduct our experiments.

5.1 Learning a Specific Ability

We select two distinct abilities to test the effectiveness of model training based on their sensitivities:

Logical Reasoning: As demonstrated in Figure 7, this ability is data-biased, meaning it is more sensitive to data scaling than to parameter scaling.

Novel Generation³: This OOD task involves continuing a novel from a given starting point. According to the calculations described in Section 4.3.2, Novel Generation is parameter-biased with a parameter sensitivity of 0.92 and a data sensitivity of 0.54. This indicates that increasing the model's parameter size is more effective for enhancing this ability than increasing the data volume.

Table 2 illustrates the performance of Qwen1.5 and Baichuan2 when trained with varying amounts of data and parameters, following the evaluation methodology outlined in Section 4. The results show that for the data-biased ability, Logical Reasoning, increasing the data volume from 2000 to 4000 can yield performance that matches or exceeds that of models with 13/14B parameters. This suggests that annotating more data can effectively reduce computational resource requirements. Conversely, for the parameter-biased ability, Novel Generation, even a fivefold increase in data volume fails to match the performance of larger models. Therefore, for such abilities, increasing the model's parameter size is the more effective strategy.

5.2 Learning Comprehensive Abilities

Guided by the sensitivities of various abilities, we explore their effectiveness in guiding dataset construction through two comprehensive benchmarks: AGIEval (Zhong et al., 2023) and CMMLU (Li

³https://huggingface.co/datasets/zxbsmk/webnovel_cn

	Parameter Scale	Logical Reasoning		Novel Generation	
		Data Size	Score	Data Size	Score
Qwen1.5	7B	2,000	26.0	2,000	53.0
	7B	4,000	37.0	10,000	52.0
	14B	2,000	34.0	2,000	62.0
Baichuan2	7B	2,000	25.0	2,000	37.0
	7B	4,000	35.0	10,000	38.0
	13B	2,000	36.0	2,000	57.0

Table 2: Performance comparison of Qwen1.5 and Baichuan2 models on Logical Reasoning and Novel Generation tasks with varying data volumes and parameter sizes.

et al., 2023). AGIEval (Zhong et al., 2023) 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), similar to MMLU (Hendrycks et al., 2020), assesses LLMs’ knowledge and reasoning capabilities within the Chinese cultural context, covering 67 diverse subjects from elementary to advanced professional levels.

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 different growth paces across abilities inspire us to adjust their proportions. Ethics and Role-play Chat, showing low sensitivity to all factor changes in Section 4.3, are reduced to 64 instances corresponding to their relatively higher points in Figure 5. Considering that the missing data can still benefit other abilities due to generalization, we maintain a total of 10k instances by uniformly increasing the data volume of other categories. Creative Writing is excluded from this replenishment, which remains at 1k due to its performance plateau.

Maximum: We continue to expand our data volume following the same insights, keeping Ethics and Role-play Chat at 64 and Creative Writing at 1k instances. Other abilities are expanded 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 each ability’s annotation.

We train 7b models for each construction approach and test their performance at epochs 10 on two benchmarks under both 0-shot and 5-shot settings. Table 3 marks the results that show improvement over the baseline with \uparrow . In our experiments with three distinct foundation models, both new

Models	Data Quantity	AGIEval - 0shot			CMMLU - 0shot		
		CLLaMA	Baichuan2	Qwen1.5	CLLaMA	Baichuan2	Qwen1.5
Baseline	10k	34.64	42.15	69.08	36.75	52.60	72.71
Reconstruct	10k	35.43 \uparrow	45.09	69.56 \uparrow	36.85 \uparrow	53.00 \uparrow	73.07 \uparrow
Maximum	40k	37.61 \uparrow	46.59 \uparrow	69.21 \uparrow	37.28 \uparrow	56.50 \uparrow	72.33
		AGIEval - 5shot			CMMLU - 5shot		
		CLLaMA	Baichuan2	Qwen1.5	CLLaMA	Baichuan2	Qwen1.5
Baseline	10k	31.01	47.03	70.12	35.14	54.87	71.97
Reconstruct	10k	32.37 \uparrow	48.46 \uparrow	70.97 \uparrow	35.89 \uparrow	55.00 \uparrow	71.82
Maximum	40k	33.57 \uparrow	53.12 \uparrow	70.96 \uparrow	37.16 \uparrow	58.02 \uparrow	71.18

Table 3: Comparing the performance of three construction approaches on two benchmarks, evaluated using checkpoints at epoch 10 with a parameter size of 7b. Scores superior to the baseline are marked with \uparrow .

strategies generally outperform the "Baseline" approach. Specifically, the "Reconstruct" strategy improves AGIEval scores by 1%-3% without increasing the data volume. Moreover, the "Maximum" strategy further enhances performance across all evaluated abilities.

6 Conclusion

This research introduces a novel, human-curated Chinese dataset comprising over 40,000 instruction instances across ten ability categories. This dataset facilitates the investigation of how the growth of large language model (LLM) abilities is influenced by data volume, parameter size, and data construction methods during instruction tuning. We are the first to disentangle the effects of these factors by examining a comprehensive set of over 500 model checkpoints, ranging from 7 billion to 33 billion parameters. Our findings reveal that the impact of increasing data volume and model size varies across different abilities. We identify two key features, Complexity and Transference, which can predict ability growth in low-resource scenarios. Guided by these findings, we enhance the effectiveness of learning specific tasks and achieve better comprehensive abilities on benchmarks such as CMMLU and AGIEval.

7 Limitations

One limitation of our study is that when calculating Complexity and Transference in a low-resource setting (7b foundation models, 64 data points per ability), most abilities display a clear linear relationship with their sensitivities, but a few outliers need to be more accurately predicted. We then try to expand the training data to 1000 instances. As illustrated in Appendix 8, allocating more resources for calculation helps mitigate this issue. Moreover, the quantity of human-curated data can be continuously expanded to further explore performance plateaus of other abilities beyond Creative Writing.

585

References

586

Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Hestlow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.

592

593

Stephen H. Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Févry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-David, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Alan Fries, Maged S. Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Mike Tian-Jian Jiang, and Alexander M. Rush. 2022. Promptsource: An integrated development environment and repository for natural language prompts. *ArXiv*, abs/2202.01279.

604

605

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

608

609

Satanjeev Banerjee and Alon Lavie. 2005. 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*, pages 65–72.

614

615

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

620

621

Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.

626

627

Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937*.

629

630

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. 2023. *Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality*.

635

636

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,

640

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. 2022. Palm: Scaling language modeling with pathways. *ArXiv*, abs/2204.02311.

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

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. 2022. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416.

659

660

661

662

663

664

665

666

667

668

669

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world’s first truly open instruction-tuned llm.

670

671

672

673

674

Yiming Cui, Ziqing Yang, and Xin Yao. 2023. *Efficient and effective text encoding for chinese llama and alpaca*. *arXiv preprint arXiv:2304.08177*.

675

676

677

Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*.

678

679

680

681

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.

682

683

684

685

Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. 2017. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*.

686

687

688

689

690

Or Honovich, Thomas Scialom, Omer Levy, and Timo Schick. 2022. Unnatural instructions: Tuning language models with (almost) no human labor. *arXiv preprint arXiv:2212.09689*.

691

692

693

694

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao

695

696

697

698	Fu, Maosong Sun, and Junxian He. 2023. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. <i>arXiv preprint arXiv:2305.08322</i> .	751
699		752
700		753
701		754
702	Hamish Ivison, Akshita Bhagia, Yizhong Wang, Hannaneh Hajishirzi, and Matthew E. Peters. 2022. Hint: Hypernetwork instruction tuning for efficient zero-shot generalisation. <i>ArXiv</i> , abs/2212.10315.	755
703		756
704		757
705		758
706	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. 2022. OPT-IML: scaling language model instruction meta learning through the lens of generalization. <i>CoRR</i> , abs/2212.12017.	759
707		760
708		761
709		762
710		763
711		764
712		765
713		766
714	Yunjie Ji, Yong Deng, Yan Gong, Yiping Peng, Qiang Niu, Baochang Ma, and Xiangang Li. 2023. Belle: Be everyone’s large language model engine. https://github.com/LianjiaTech/BELLE .	767
715		768
716		769
717		
718	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .	770
719		771
720		772
721		773
722		774
723	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. 2023. Openassistant conversations—democratizing large language model alignment. <i>arXiv preprint arXiv:2304.07327</i> .	775
724		776
725		777
726		778
727		779
728		780
729	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. Cmm1u: Measuring massive multitask language understanding in chinese. <i>arXiv preprint arXiv:2306.09212</i> .	781
730		782
731		
732		783
733		784
734	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step. <i>arXiv preprint arXiv:2305.20050</i> .	785
735		786
736		787
737		788
738		789
739	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text summarization branches out</i> , pages 74–81.	790
740		791
741		792
742	Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. 2023. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. <i>IEEE/ACM Transactions on Audio, Speech, and Language Processing</i> .	793
743		794
744		795
745		796
746		797
747	Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. <i>arXiv preprint arXiv:2303.15621</i> .	798
748		799
749		800
750		801
		802
		803
		804
		805
		806
	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-task generalization via natural language crowdsourcing instructions. <i>arXiv preprint arXiv:2104.08773</i> .	
	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. 2022. Crosslingual generalization through multitask finetuning. <i>arXiv preprint arXiv:2211.01786</i> .	
	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback . <i>ArXiv preprint</i> , abs/2112.09332.	
	OpenAI. 2022. Introducing chatgpt. <i>OpenAI Blog</i> .	
	OpenAI. 2023. GPT-4 Technical Report . <i>arXiv e-prints</i> , page arXiv: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. 2022. Training language models to follow instructions with human feedback. <i>CoRR</i> , abs/2203.02155.	
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	
	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. <i>arXiv preprint arXiv:2304.03277</i> .	
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	
	Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. 2021. Multitask prompted training enables zero-shot task generalization. <i>arXiv preprint arXiv:2110.08207</i> .	
	Chenhui Shen, Liying Cheng, Yang You, and Lidong Bing. 2023. Are large language models good evaluators for abstractive summarization? <i>arXiv preprint arXiv:2305.13091</i> .	
	Chiyu Song, Hongliang He, Haofei Yu, Pengfei Fang, Leyang Cui, and Zhenzhong Lan. 2023. Uni-encoder: A fast and accurate response selection paradigm for generation-based dialogue systems . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 6231–6244, Toronto, Canada. Association for Computational Linguistics.	

807	Kai Sun, Dian Yu, Dong Yu, and Claire Cardie. 2020. Investigating prior knowledge for challenging chinese machine reading comprehension. <i>Transactions of the Association for Computational Linguistics</i> , 8:141–155.	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022b. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> .	864 865 866 867 868
812	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, et al. 2023. Zero-shot information extraction via chatting with chatgpt. <i>arXiv preprint arXiv:2302.10205</i> .	869 870 871 872 873
817	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. 2023. Llama: Open and efficient foundation language models. <i>ArXiv</i> , abs/2302.13971.	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions.	874 875 876 877
824	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022a. Self-consistency improves chain of thought reasoning in language models. <i>ArXiv</i> , abs/2203.11171.	Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023b. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. <i>arXiv preprint arXiv:2304.01196</i> .	878 879 880 881
828	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. 2023. How far can camels go? exploring the state of instruction tuning on open resources. <i>ArXiv</i> , abs/2306.04751.	Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. 2023. Baichuan 2: Open large-scale language models. <i>arXiv preprint arXiv:2309.10305</i> .	882 883 884 885
834	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022b. Self-instruct: Aligning language model with self generated instructions. <i>arXiv preprint arXiv:2212.10560</i> .	Da Yin, Xiao Liu, Fan Yin, Ming Zhong, Hritik Bansal, Jiawei Han, and Kai-Wei Chang. 2023. Dynosaur: A dynamic growth paradigm for instruction-tuning data curation. <i>arXiv preprint arXiv:2305.14327</i> .	886 887 888 889
839	Yizhong Wang, Swaroop Mishra, Pegah Alipoor-molabashi, Yeganeh Kordi, Amirreza Mirzaei, Anjana Arunkumar, Arjun Ashok, Arut Selvan Dhanasekaran, Atharva Naik, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Gary Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krma Doshi, Maitreya Patel, Kuntal Kumar Pal, M. Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Shailaja Keyur Sampat, Savan Doshi, Siddharth Deepak Mishra, Sujan Reddy, Sumanta Patro, Tanay Dixit, Xudong Shen, Chitta Baral, Yejin Choi, Noah A. Smith, Hanna Hajishirzi, and Daniel Khashabi. 2022c. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Biao Zhang, Zhongtao Liu, Colin Cherry, and Orhan Firat. 2024. When scaling meets llm finetuning: The effect of data, model and finetuning method. <i>arXiv preprint arXiv:2402.17193</i> .	890 891 892 893
857	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022a. Finetuned language models are zero-shot learners. In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	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. 2023a. Chinese open instruction generalist: A preliminary release. <i>ArXiv</i> , abs/2304.07987.	894 895 896 897 898
		Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023b. Instruction tuning for large language models: A survey. <i>arXiv preprint arXiv:2308.10792</i> .	899 900 901 902 903
		Yue Zhang, Leyang Cui, Deng Cai, Xinting Huang, Tao Fang, and Wei Bi. 2023c. Multi-task instruction tuning of llama for specific scenarios: A preliminary study on writing assistance. <i>ArXiv</i> , abs/2305.13225.	904 905 906 907
		Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In <i>International Conference on Machine Learning</i> .	908 909 910 911
		Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> .	912 913 914 915 916

917 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang,
918 Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen,
919 and Nan Duan. 2023. Agieval: A human-centric
920 benchmark for evaluating foundation models. *arXiv*
921 *preprint arXiv:2304.06364*.

922 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao
923 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
924 Lili Yu, et al. 2023. Lima: Less is more for alignment.
925 *arXiv preprint arXiv:2305.11206*.

926 A Appendix

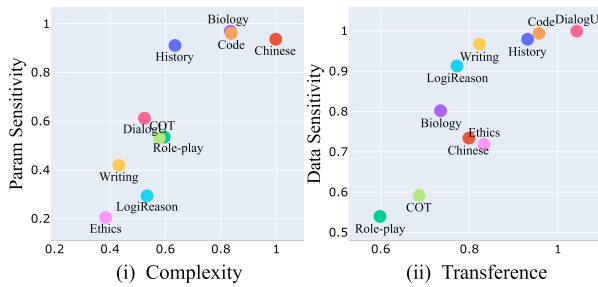


Figure 8: Two interpretable features, (i) Complexity and (ii) Transference, of different abilities demonstrate linear relationships with their sensitivities to scaling in parameter size and data volume. These features are calculated using 1000 data points instead of 64 in Section 5. Sensitivity values are normalized to (0, 1) for visualization.

927 A.1 Training Procedures

928 Our hyperparameter choices are generally in line
929 with [Zhou et al. \(2023\)](#), using AdamW optimizer
930 with $\beta_1 = 0.9$, $\beta_2 = 0.95$ and $weight_decay =$
931 0.1 . The peak learning rate is set to $1e-5$ and lin-
932 earlyly decays to $1e-6$ by the end of training. The
933 batch size is 64, with inputs exceeding 2048 tokens
934 trimmed. Referring to the protocol in [Song et al.](#)
935 [\(2023\)](#), we also incorporate two speaker tokens,
936 $\langle user \rangle$ and $\langle assistant \rangle$, to segment the utter-
937 ances from instruction and output for training and
938 inference.

939 A.2 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>
Dialogue Understanding	<p>Prompt: Man: What kind of animal do you like?</p> <p>Woman: Elephants. They have long trunks, big ears, and they look so clumsy, which makes them adorable.</p> <p>Man: Elephants aren't clumsy, though. They often help humans with work.</p> <p>Woman: I said they "look" clumsy, you know.</p> <p>Why does the woman like elephants? A. Elephants are big. B. Elephants are clumsy. C. Elephants are adorable. D. Elephants are smart.</p> <p>Response: C</p>

Continued on next page

Ability	Example
Role-play Chat	<p>Prompt: 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. John, on the other hand, is an experienced HR manager with profound expertise and extensive experience in talent acquisition. Generate a short conversation happening between them.</p> <p>Response: 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>
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
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>
Ethics	<p>Prompt: 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
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. B. Flyers promoting a commercial real estate project distributed on crowded streets. C. Coupons for a large chain of fast-food restaurants received in an email inbox. D. Two screen protectors received as a gift when purchasing a mobile phone at an electronics store.</p> <p>Response: D</p>
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> $\frac{1}{17} \begin{pmatrix} 7 & -3 \\ 1 & 2 \end{pmatrix}.$

Continued on next page

Ability	Example
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>
Code Generation	<p>Prompt: 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: $0 \leq \text{left} \leq \text{right} \leq 231 - 1$ Please use c# to finish the problem.</p> <p>Response: The following are reference implementation options: ““cs public class Solution { public int RangeBitwiseAnd(int left, int right) { while (left < right) { right &= (right - 1); } return right; } }</p>

Models	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

Table 4: Comparing the performance of three mixing strategies with synthetic data on two benchmarks, evaluated using checkpoints at epochs 5, 10, and 15 with a parameter size of 7b. Highest performance under each setting is in bold.

940 A.3 Mix up with Syhtnetic Data

941 Synthetic data is a rich open resource, but Section 4.3 indicates that "*Synthetic data does not consistently*
942 *enhance model performance with increased volume.*" Investigating the optimal use of synthetic data
943 alongside human-curated data is crucial for practical applications. We thus utilize the "Maximum"
944 construction from Section 5.2 as our baseline and then integrate varying quantities (0, 2560, 41k) of
945 synthetic data to train 7b models.

946 Table 4 presents the efficacy of three mixing strategies at epochs 5, 10, and 15, evaluated on two
947 benchmarks in both 0-shot and 5-shot scenarios. For AGIEval, a modest addition of 2,560 synthetic
948 instances shows performance gains. In contrast, for CMMLU, peak performance is achieved with no
949 synthetic data or a similar modest addition. Notably, further incorporating 41k synthetic instances does
950 not yield further benefits and instead results in a marked decrease in performance, consistent with our
951 earlier observations in Section 4.

Prompt:

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

B: **An older music enthusiast who enjoys listening to classical and jazz music.**

Generate a snippet of their conversation.

Ground Truth:

[history conversations]

A: I see. Do you have any other music recommendations?

B: I've been listening to **Billie Holiday's jazz music** lately. Are you interested in giving it a try?

Fine-grained corruption:

[history conversations]

A: I see. Do you have any other music recommendations?

B: I've been listening to **Justin Bieber's jazz music** lately. Are you interested in giving it a try?

Coarse-grained corruption:

A: Hello, I'm A. Lately, I've gotten into cooking and I especially enjoy trying out various new recipes. Are you interested in cooking?

B: Hello, I'm B. Actually, I'm quite interested in cooking too, especially baking. Recently, I've learned to make a delicious chocolate cake.

[folded conversations]

Prompt:

[...]

According to foreign media reports, Microsoft announced that it would lay off 10,000 employees, about **5%** of its total workforce, by the end of March. On January 20th, Google suddenly announced a layoff of 12,000 people, accounting for **6%** of its workforce. On January 26th, The Wall Street Journal reported that IBM had joined the ranks of layoffs and would cut about 3,900 people...

[...]

Write a news article based on the given information.

Ground Truth:

Title: After the Chinese New Year Holiday, Wave of Layoffs Hits - Cherishing Work Becomes the Greatest Sense of Security for Adults

Article:

[...]

According to foreign media reports, Microsoft has announced plans to lay off 10,000 employees by the end of March, accounting for approximately **5%** of its total workforce. On January 20th, Google made a sudden announcement to lay off 12,000 people, representing **6%** of its total workforce. On January 26th, The Wall Street Journal reported that IBM has also joined the layoff trend, planning to cut approximately 3,900 jobs.

[...]

Fine-grained corruption:

Title: The Chinese New Year Holiday Ends, Most Companies Have Ceased Operations

Article:

[...]

According to foreign media reports, Microsoft plans to lay off 1,000 employees by the end of March, accounting for approximately **10%** of its total workforce. On January 20th, Google made a sudden announcement to lay off 12,000 people, representing **60%** of its total workforce. The Wall Street Journal reported on January 26th that IBM has also joined the wave of layoffs, intending to terminate approximately 3,900 employees.

[...]

Coarse-grained corruption:

According to the latest reports, several cities in China have recently experienced snow disasters and cold wave weather, causing significant disruptions to people's lives and travel. It has been reported that heavy snowfall has led to the closure of many major transportation routes, resulting in a large number of vehicles stranded on the roads and many people trapped in traffic jams. Additionally, power outages in many cities due to the snow disaster have inconvenienced residents' daily lives.

Figure 9: 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 10: Examples of distractors derived from manually corrupted Code Generation output.

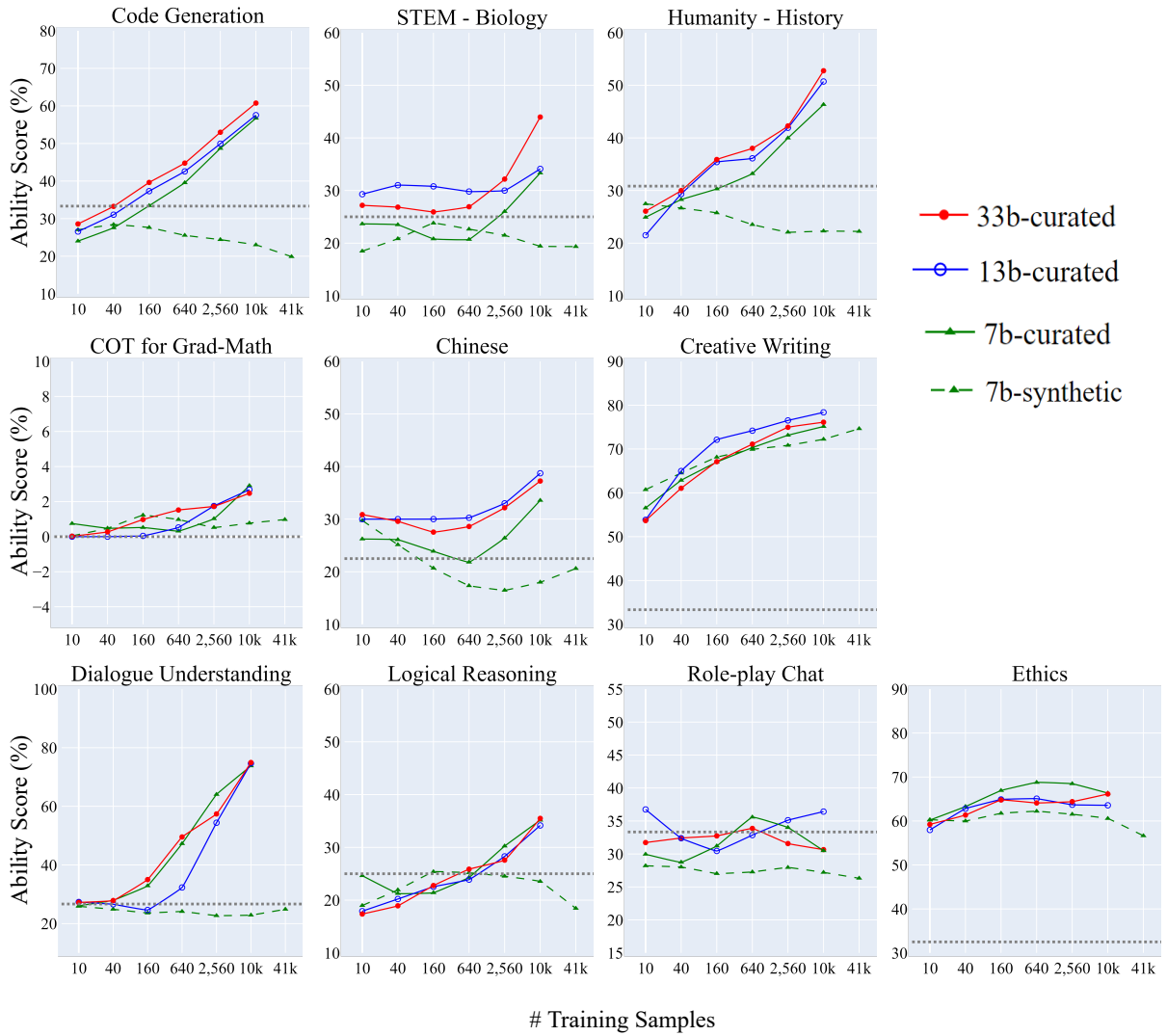


Figure 11: Evaluations of 7b models trained on synthetic data, yielding analogous conclusions as 13b models.