

# Understanding the Instruction Mixture for Large Language Model Fine-tuning

Anonymous ACL submission

## Abstract

While instructions fine-tuning of large language models (LLMs) has been proven to enhance performance across various applications, the influence of the instruction dataset mixture on LLMs has not been thoroughly explored. In this study, we classify instructions into three main types: NLP downstream tasks, coding, and general chatting, and investigate their impact on LLMs. Our findings reveal that specific types of instructions are more beneficial for particular uses, while it may cause harms to other aspects, emphasizing the importance of meticulously designing the instruction mixture to maximize model performance. This study sheds light on the instruction mixture and paves the way for future research.

## 1 Introduction

Supervised fine-tuning (SFT) has been proven to be an effective approach to align large language models (LLMs) with human instructions, enhancing downstream task performance, and facilitating code generation (Touvron et al., 2023; Muennighoff et al., 2023; Gunasekar et al., 2023). Previous research has demonstrated that directly transforming NLP downstream tasks (e.g., coreference resolution) into instruction-response pairs, followed by fine-tuning models on such datasets, leads to improvements in performance across various benchmarks (Sanh et al., 2022). Recent studies show that general-purpose instructions can enhance both the performance of LLMs on downstream tasks and their alignment with human intents (Taori et al., 2023; Touvron et al., 2023; Zeng et al., 2023). Additionally, incorporating code datasets has been shown to enhance a model’s logical reasoning abilities (Fu and Khot, 2022; Gunasekar et al., 2023).

As LLMs continue to advance, researchers are keen to endow a single model with diverse abilities. One straightforward approach is to combine multiple specialized instruction datasets. For instance,

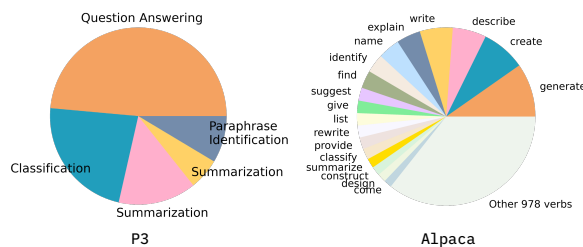


Figure 1: Instruction type distribution of P3 and Alpaca. For P3, the statistics come from the original dataset. For Alpaca, we utilize a dependency parsing approach to extract the root verb of each instruction.

to enhance a model’s code generation ability, we might incorporate CodeAlpaca into the SFT data (Chaudhary, 2023). However, there is no standard way of selecting instruction datasets explicitly. The process of mixing different datasets and understanding how various instruction types interact with each other remains underexplored.

In this paper, our focus is on evaluating the model’s performance in three key areas: NLP downstream task performance, coding ability, and chat capabilities. We aim to investigate how different distributions of instruction datasets can impact model performance across these diverse aspects. Our selected representative SFT datasets are P3 (Sanh et al., 2022) for NLP downstream tasks, CodeAlpaca (Chaudhary, 2023) for code generation, and Alpaca for general-purpose instructions.

According to Figure 1, P3 mainly contains five types of tasks (including QA, classification, summarization, etc.), while Alpaca can be classified into about 1K types based on the root verb of each instruction, where the top 3 root verbs are *generate*, *create* and *describe*. Additionally, codeAlpaca is exclusively focused on code-related tasks. Notably, when compared to the general instructions in Alpaca, instruction-response pairs in both P3 and CodeAlpaca exhibit narrower variations, whether in format or content.

069 We conduct experiments with eight differ- 119  
070 ent mixture settings involving these instruction 120  
071 datasets, assessing model performance in NLP 121  
072 downstream tasks, coding proficiency, and align- 122  
073 ment skills (i.e., chat abilities). Our extensive ex- 123  
074 periments yield the following insights:: (1) Using 124  
075 a single type of STF data consistently improves 125  
076 the performance of the model on the correspond- 126  
077 ing task, while all three types of instructions can 127  
078 be used to improve NLP downstream task perfor- 128  
079 mance. (2) Incorporating instructions that are sim- 129  
080 ply reformatted from NLP downstream tasks (e.g., 130  
081 P3) downgrades the model’s alignment skills, re- 131  
082 sulting in a worse chat experience. (3) Exploiting 132  
083 code instructions can improve the model’s coding 133  
084 ability and boost the alignment skills.

085 Based on our findings, we suggest that re- 134  
086 searchers should carefully design the instruction 135  
087 mixture to maximize the model’s performance on 136  
088 the target usage while taking model size into con- 137  
089 sideration. We also appeal to the community to 138  
090 evaluate LLMs not only based on the benchmark 139  
091 performance but also on the alignment skills. 140

## 092 2 Related Work 141

093 Recent work has demonstrated that vanilla LLMs 142  
094 can be effective at following general language in- 143  
095 structions if tuned with instructions alongside their 144  
096 corresponding outcomes (Mishra et al., 2022; Sanh 145  
097 et al., 2022; Wang et al., 2022). To construct such 146  
098 instructional datasets, researchers have used a vari- 147  
099 ety of approaches. For example, Sanh et al. (2022) 148  
100 reformat a large set of supervised datasets using 149  
101 multiple prompt templates to create P3, a dataset of 150  
102 instruction-response pairs covering a wide range of 151  
103 NLP downstream tasks. Even though such LLMs 152  
104 perform well on NLP downstream tasks, they do 153  
105 not align well with human behavior as chatbots. 154

106 To facilitate the general-purpose LLMs fine- 155  
107 tuning, Conover et al. (2023) introduced a high- 156  
108 quality human-annotated instruction dataset tai- 157  
109 lored for LLMs. However, this method is resource- 158  
110 intensive and lacks scalability. To address such 159  
111 issue, Wang et al. (2023) and Taori et al. (2023) 160  
112 firstly use an automatic data collection approach 161  
113 to collect a large-scale general instruction dataset. 162  
114 Based on the proposed datasets or approaches, later 163  
115 work expands the dataset size (Wu et al., 2023), 164  
116 language coverage (Li et al., 2023), and task types 165  
117 (Chaudhary, 2023; Yue et al., 2023). 166

118 With the increasing capability of LLMs and 167

119 availability of instruction datasets, researchers aim 120  
121 to imbue a single model with diverse capabilities. 122  
123 Sengupta et al. (2023) have attempted to blend dif- 124  
125 ferent instruction datasets without considering the 126  
127 data volume and task types. Longpre et al. (2023) 128  
129 propose that increasing the number of tasks and 130  
131 instruction diversity can enhance performance. In 132  
133 contrast, Anand et al. (2023) excluded P3 from their 134  
135 fine-tuning dataset, seemingly to enhance align- 136  
137 ment skills. 138

139 Nevertheless, none of these works systematically 140  
141 investigate the impact of instruction mixture on 142  
143 LLMs. Our work aims to find the impact of mixing 144  
145 together different instructions to align models. 146

## 147 3 Experimental Setup 148

149 **SFT Datasets** We select Alpaca (Taori et al., 150  
151 2023) as the *general instruction dataset* for align- 152  
153 ing models, which contains 52K instruction- 154  
155 response pairs. We use P3 (Sanh et al., 2022) as our 156  
157 *NLP task instruction dataset*, which is reformatted 158  
159 for a wide range of NLP downstream tasks using di- 160  
161 verse human-written templates. Since the number 161  
162 of samples in each task varies vastly, we randomly 162  
163 sample 1K data from each subtask formatted with 163  
164 several corresponding prompts for diversity, result- 164  
165 ing in 660K samples. For *coding data*, we choose 165  
166 CodeAlpaca (Chaudhary, 2023), which is an in- 166  
167 struction dataset focusing on code generation. It 167  
168 contains 20K samples with different programming 168  
169 languages. To ensure a balanced comparison, we 169  
170 utilize a 20K subset from each type of dataset, ran- 170  
171 domly sampled. Unless explicitly mentioned, the 171  
172 experiments and discussions will be based on these 172  
173 subsets for the rest of this paper. 173

174 **Evaluation** We divide the evaluation into three 174  
175 parts: NLP benchmark performance, code gen- 175  
176 eration, and alignment evaluation (i.e., chat abil- 176  
177 ity evaluation). For NLP benchmarks, we use 177  
178 ARC (Clark et al., 2018), Winogrande (Sakaguchi 178  
179 et al., 2021), PIQA (Bisk et al., 2020), MMLU 179  
180 (Hendrycks et al., 2020), RACE (Lai et al., 2017), 180  
181 and HellaSwag (Zellers et al., 2019) datasets. For 181  
182 coding, we use HumanEval (Chen et al., 2021), 182  
183 which tests the pass rate of generated codes. For 183  
184 alignment evaluation, we utilize FLASK (Ye et al., 184  
185 2023) framework to score models’ alignment skills. 185  
186 We keep the eight most frequent alignment skills 186  
187 from the original evaluation set, resulting in 1,180 187  
188 samples. Then we employ GPT-4 to assess mod- 188  
189 els’ responses to each instruction sample based 189

Model	Data	ARC (challenge)	Wino-grande	PIQA	MMLU	Race	Hella-Swag	Average	HumanEval @1	HumanEval @10
LLaMA-2-7B	None	43.09	69.53	77.97	40.81	39.23	57.20	54.64	13.72	21.34
	A	47.78	67.64	78.24	42.19	<u>44.50</u>	<b>61.09</b>	56.91	13.48	17.07
	C	46.08	69.46	<u>78.50</u>	40.99	41.05	<u>60.96</u>	56.17	16.22	<u>24.39</u>
	P	<u>49.57</u>	<b>71.43</b>	<b>79.00</b>	<b>45.98</b>	43.45	59.44	<b>58.15</b>	4.63	7.93
	AC	47.10	66.93	78.13	40.42	44.21	59.70	56.08	<b>17.50</b>	<b>25.00</b>
	AP	48.38	70.01	78.07	43.84	42.87	58.46	56.94	13.84	17.68
	CP	47.95	<u>71.27</u>	78.40	<u>44.91</u>	44.40	60.69	<u>57.94</u>	<u>16.77</u>	20.12
	ACP	<b>49.66</b>	68.03	77.86	43.52	<b>44.59</b>	58.73	57.07	15.98	23.78
LLaMA-2-13B	None	48.55	71.90	79.16	<b>52.12</b>	40.67	60.12	58.75	15.43	26.22
	A	54.10	71.19	80.03	47.86	<b>47.08</b>	<b>65.58</b>	60.97	15.06	20.73
	C	49.66	73.40	<b>80.79</b>	<u>51.50</u>	45.36	63.63	60.72	17.87	24.39
	P	54.27	<u>74.19</u>	80.03	50.30	<u>45.55</u>	62.46	<u>61.13</u>	0.30	1.83
	AC	51.62	68.75	<u>80.58</u>	48.68	44.40	62.97	59.50	17.07	<u>27.44</u>
	AP	<u>54.79</u>	71.74	80.30	51.15	45.17	62.72	60.98	8.29	14.63
	CP	<b>55.38</b>	<b>74.59</b>	80.52	51.42	<u>45.55</u>	<u>63.85</u>	<b>61.89</b>	<u>18.23</u>	25.00
	ACP	54.44	71.51	80.03	49.98	<b>47.08</b>	63.14	61.03	<b>20.24</b>	<b>32.93</b>

Table 1: Results on NLP and code generation benchmarks. All experiments are done in a zero-shot setting. The best result is in **bold**, and the second best result is underlined.

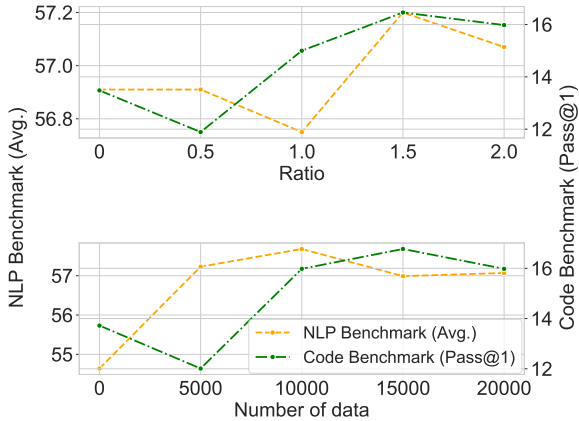


Figure 2: NLP benchmark scores (avg) and Code benchmark (HumanEval) scores for LLaMA-2-7B tuned with different mixing ratios and different number of data. We keep the number of Alpaca to 20K and change the number of P3 and CodeAlpaca to get different ratios.

on human-written principles. See appendix A for demonstrations of these skills.

**SFT Setup** We employ LLaMA-2 (7B, and 13B) models (Touvron et al., 2023). We fine-tune the models for two epochs in a generative way as in Radford et al. (2018). We use a linear scheduler with a 3% warmup rate and a batch size of 64. The maximum learning rate is  $5 \times 10^{-5}$ .

## 4 Results

For short, we denote Alpaca, CodeAlpaca, and P3 datasets as A, C, P, respectively. For each model, we compare eight different data mixing strategies, denoted as None, A, C, P, AC, AP, CP, ACP, where

None represents the vanilla model without fine-tuning, and each of the other settings represents the model fine-tuned with the corresponding dataset. For example, AC means the model is fine-tuned with both Alpaca and CodeAlpaca. We use the same naming convention for the rest of the paper.

### 4.1 NLP Tasks and Code Benchmark Results

Table 1 shows the zero-shot results on NLP and code generation benchmarks. Notably *each type of specialized instructions improve the performance on the benchmarks they are designed for*. In the no-mixture setting (comparing A, C, and P), models fine-tuned on P3 achieve the highest average score for NLP tasks, while models fine-tuned on CodeAlpaca excel in code generation benchmarks. Examining specific tasks reveals that *a model’s performance on specific task heavily relies on the similarity between the target task and the tasks it was fine-tuned on*. For instance, Alpaca fine-tuned models excel in Race and HellaSwag, which involve story completion task similar to the Alpaca instruction format. On the other hand, P3 fine-tuned models outperform in ARC and Winogrande, which involve multiple-choice QA and cloze tests, aligning with P3’s data.

In the mixture setting, it’s evident that *including specialized data consistently boosts model performance in corresponding benchmarks compared to models without such data*. For example, P, PA, PC, and PCA perform better than None, A, C, and CA on NLP downstream tasks. Focusing on code benchmarks, *incorporating general instructions*

Model	Data	Corr.	Fact.	Comm.	Compr.	Compl.	Insight.	Read.	Conc.	Avg.
LLaMA-2-7B	A	47.6	<b>55.4</b>	58.8	<u>54.8</u>	<u>48.0</u>	<b>50.4</b>	<b>88.0</b>	81.6	<u>60.6</u>
	C	48.8	52.0	58.4	52.0	40.2	46.2	83.8	78.4	57.4
	P	47.2	40.0	48.8	38.4	29.0	30.4	64.4	68.6	45.8
	AC	<u>49.0</u>	<u>54.4</u>	<b>59.6</b>	<b>56.4</b>	<b>48.2</b>	<u>49.8</u>	<u>86.6</u>	<b>85.6</b>	<b>61.2</b>
	AP	48.4	51.4	57.6	52.6	45.0	46.0	84.2	80.8	58.2
	CP	47.0	49.6	54.2	48.8	36.2	41.8	78.2	77.2	54.2
	ACP	<b>50.4</b>	53.0	<u>59.0</u>	53.8	47.2	46.8	85.0	<u>81.8</u>	59.6
LLaMA-2-13B	A	53.6	<u>58.8</u>	<u>63.8</u>	<u>60.0</u>	<u>47.6</u>	<b>55.2</b>	<b>89.2</b>	<u>84.0</u>	<u>64.0</u>
	C	<b>57.2</b>	<u>58.8</u>	61.0	57.8	43.8	52.4	85.6	82.2	62.4
	P	49.4	42.4	51.8	42.0	28.2	32.0	66.8	70.4	47.8
	AC	<u>55.6</u>	<b>61.0</b>	<b>66.6</b>	<b>61.2</b>	<b>51.4</b>	<u>54.0</u>	<u>88.4</u>	<b>86.6</b>	<b>65.6</b>
	AP	53.0	55.4	60.6	56.2	47.0	48.0	85.0	83.4	61.0
	CP	53.0	53.2	57.4	53.4	39.0	45.2	81.2	82.6	58.2
	ACP	51.6	55.6	61.8	57.0	47.0	48.6	87.0	83.0	61.4

Table 2: GPT-4 evaluation results on alignment skill assessment. We report eight dimensions, i.e., logical correctness, factuality, commonsense understanding, comprehension, completeness, insightfulness, readability, and conciseness, as well as average scores. Since vanilla model cannot follow instructions, we exclude its result here. The best result is in **bold**, and the second best result is underlined.

consistently improves coding performance. For the 7B model, AC improves performance by +1.28 and +0.61 compared to C, while the improvements are -0.80 (outlier) and +3.05 for the 13B models. Another interesting finding is that the 13B models achieve their best results with the ACP mixture, while the 7B models perform best with AC. This suggests that *larger models have greater capacities and can better leverage various instructions*.

These findings highlight the importance of considering model size and target usage when designing instruction mixture plans.

**Mixing with Different Ratios** Despite knowing mixing specialized instructions are vital for better benchmark performance, how the mixing ratio correlates with the performance is also important for the best training strategy. As Figure 2 shows, given the number of general instructions fixed to 20K, scores of both NLP task benchmarks and code benchmarks first decrease and then increase as the ratio of specialized instructions increases. They both reach the maximum when the ratio is 1.5, while slightly decrease when the ratio continues to increase to 2.0. We think this is because the model is overfitted to the specialized instructions when there are too many such instructions.

**Number of data** Figure 2 also shows the performance change with respect to the number of fine-tuning data. We mix each type of instruction with the same number. We find that the performance of both benchmarks reaches a relatively stable state when the number of data is larger than 10K.

## 4.2 Alignment Skills Results

Table 2 shows the alignment skills evaluation results. We adopt the same setup as FLASK, using GPT-4-0613 to access the alignment skills and scaling the scores to the range of [0, 100].

From Table 2 we have the following findings: (1) *All three types of instructions improves model alignment compared to the vanilla LLM*. Among these instructions, Alpaca stands out as the most effective. It contains general-purpose instructions and human-like responses, making it a better fit for aligning models with humans. (2) *While CodeAlpaca alone doesn't significantly enhance alignment abilities, combining it with general instructions results in a substantial improvement of +0.6 (7B) and +1.6 (13B) points* These improvements are mainly attributed to better compression, commonsense understanding, completeness, and conciseness. (2) *Mixing P3 data causes a drop of -2.8 (7B) and -3.6 (13B) in average alignment skills*. This indicates that P3 tends to have a negative impact on fine-tuning chatbot LLMs.

## 5 Conclusion

In this paper, we investigated different data mixing strategies in instruction fine-tuning. We measured models with diverse benchmarks and alignment skills. We find that general instructions provide better alignment skills as well as performance on NLP benchmarks, code instructions improve coding and alignment skills, while NLP task reformatted instructions hinder alignment skills when combined with other instruction types.



## 278 Limitation

279 Our work is subject to several limitations that  
280 should be addressed in future research:

- 281 • We only use LLaMA-2 7B and 13B models  
282 in our experiments. Other models in various  
283 sizes can be used to verify our findings. We  
284 acknowledge that the model’s behavior may  
285 vary with different sizes, usually, larger mod-  
286 els have better capabilities, and hence may  
287 be able to handle more instructions without  
288 performance drop in any evaluation setting.
- 289 • In this paper, we limit our instruction dataset  
290 to 20K and mainly compare the 1:1 ratio of  
291 all instruction types. We leave the exploration  
292 of the impact of more instructions and mixing  
293 ratios to future work.

294 We acknowledge these limitations and propose that  
295 future work should focus on addressing them to  
296 help the community better understand the impact  
297 of instruction mixture on LLMs.

## 298 References

299 Yuvanesh Anand, Zach Nussbaum, Brandon Duder-  
300 stadt, Benjamin Schmidt, and Andriy Mulyar. 2023.  
301 Gpt4all: Training an assistant-style chatbot with large  
302 scale data distillation from gpt-3.5-turbo. <https://github.com/nomic-ai/gpt4all>.

303

304 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi,  
305 et al. 2020. Piqa: Reasoning about physical com-  
306 monsense in natural language. In *Proceedings of the*  
307 *AAAI conference on artificial intelligence*, volume 34,  
308 pages 7432–7439.

309 Sahil Chaudhary. 2023. Code alpaca: An instruction-  
310 following llama model for code generation. <https://github.com/sahil280114/codealpaca>.

311

312 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming  
313 Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-  
314 plan, Harri Edwards, Yuri Burda, Nicholas Joseph,  
315 Greg Brockman, et al. 2021. Evaluating large  
316 language models trained on code. *arXiv preprint*  
317 *arXiv:2107.03374*.

318 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,  
319 Ashish Sabharwal, Carissa Schoenick, and Oyvind  
320 Tafjord. 2018. Think you have solved question  
321 answering? try arc, the ai2 reasoning challenge. *arXiv*  
322 *preprint arXiv:1803.05457*.

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

Hao Fu, Yao; Peng and Tushar Khot. 2022. *How does*  
gpt obtain its ability? tracing emergent abilities of  
language models to their sources. *Yao Fu’s Notion*.

Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio  
César Teodoro Mendes, Allie Del Giorno, Sivakanth  
Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo  
de Rosa, Olli Saarikivi, Adil Salim, Shital Shah,  
Harkirat Singh Behl, Xin Wang, Sébastien Bubeck,  
Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and  
Yuanzhi Li. 2023. *Textbooks are all you need*.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,  
Mantas Mazeika, Dawn Song, and Jacob Steinhardt.  
2020. Measuring massive multitask language under-  
standing. In *International Conference on Learning*  
*Representations*.

Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang,  
and Eduard Hovy. 2017. *RACE: Large-scale ReAd-*  
*ding comprehension dataset from examinations*. In  
*Proceedings of the 2017 Conference on Empirical*  
*Methods in Natural Language Processing*, pages 785–  
794, Copenhagen, Denmark. Association for Compu-  
tational Linguistics.

Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji,  
and Timothy Baldwin. 2023. *Bactrian-x: Multilin-*  
*gual replicable instruction-following models with*  
*low-rank adaptation*.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson,  
Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le,  
Barret Zoph, Jason Wei, and Adam Roberts. 2023.  
*The flan collection: Designing data and methods for*  
*effective instruction tuning*.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, and  
Hannaneh Hajishirzi. 2022. *Cross-task generaliza-*  
*tion via natural language crowdsourcing instructions*.  
In *Proceedings of the 60th Annual Meeting of the*  
*Association for Computational Linguistics (Volume*  
*1: Long Papers)*, pages 3470–3487, Dublin, Ireland.  
Association for Computational Linguistics.

Niklas Muennighoff, Thomas Wang, Lintang Sutawika,  
Adam Roberts, Stella Biderman, Teven Le Scao,  
M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hai-  
ley Schoelkopf, Xiangru Tang, Dragomir Radev,  
Alham Fikri Aji, Khalid Almubarak, Samuel Al-  
banie, Zaid Alyafeai, Albert Webson, Edward Raff,  
and Colin Raffel. 2023. *Crosslingual generaliza-*  
*tion through multitask finetuning*. In *Proceedings*  
*of the 61st Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*,  
pages 15991–16111, Toronto, Canada. Association  
for Computational Linguistics.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya  
Sutskever, et al. 2018. Improving language under-  
standing by generative pre-training.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-  
ula, and Yejin Choi. 2021. Winogrande: An adver-  
sarial winograd schema challenge at scale. *Commu-*  
*nications of the ACM*, 64(9):99–106.

385	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	Anjana Arunkumar, David Stap, Eshaan Pathak,	445
386	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	Giannis Karamanolakis, Haizhi Lai, Ishan Puro-	446
387	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja,	hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,	447
388	et al. 2022. Multitask prompted training enables	Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,	448
389	zero-shot task generalization. In <i>ICLR 2022-Tenth</i>	Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,	449
390	<i>International Conference on Learning Representa-</i>	Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,	450
391	<i>tions</i> .	Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,	451
392	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	Shailaja Keyur Sampat, Siddhartha Mishra, Sujan	452
393	Satheesh Katipomu, Haonan Li, Fajri Koto, William	Reddy A, Sumanta Patro, Tanay Dixit, and Xudong	453
394	Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming	Shen. 2022. <i>Super-NaturalInstructions: Generaliza-</i>	454
395	Chen, Osama Mohammed Afzal, Samta Kamboj,	<i>tion via declarative instructions on 1600+ NLP tasks</i> .	455
396	Onkar Pandit, Rahul Pal, Lalit Pradhan, Zain Muham-	In <i>Proceedings of the 2022 Conference on Empiri-</i>	456
397	mad Mujahid, Massa Baali, Xudong Han, Sondo-	<i>cal Methods in Natural Language Processing</i> , pages	457
398	sodos Mahmoud Bsharat, Alham Fikri Aji, Zhiqiang	5085–5109, Abu Dhabi, United Arab Emirates. As-	458
399	Shen, Zhengzhong Liu, Natalia Vassilieva, Joel Hes-	sociation for Computational Linguistics.	459
400	tness, Andy Hock, Andrew Feldman, Jonathan Lee,	Minghao Wu, Abdul Waheed, Chiyu Zhang, Muham-	460
401	Andrew Jackson, Hector Xuguang Ren, Preslav	mad Abdul-Mageed, and Alham Fikri Aji. 2023.	461
402	Nakov, Timothy Baldwin, and Eric Xing. 2023.	<i>Lamini-lm: A diverse herd of distilled models from</i>	462
403	<i>Jais and jais-chat: Arabic-centric foundation and</i>	<i>large-scale instructions</i> .	463
404	<i>instruction-tuned open generative large language</i>	Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeon-	464
405	<i>models</i> .	bin Hwang, Seungone Kim, Yongrae Jo, James	465
406	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	Thorne, Juho Kim, and Minjoon Seo. 2023. <i>Flask:</i>	466
407	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	<i>Fine-grained language model evaluation based on</i>	467
408	and Tatsunori B. Hashimoto. 2023. Stanford alpaca:	<i>alignment skill sets</i> .	468
409	An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://</a>	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao	469
410	<a href="https://github.com/tatsu-lab/stanford_alpaca">github.com/tatsu-lab/stanford_alpaca</a> .	Huang, Huan Sun, Yu Su, and Wenhao Chen. 2023.	470
411	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	<i>Mammoth: Building math generalist models through</i>	471
412	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	<i>hybrid instruction tuning</i> .	472
413	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali	473
414	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	Farhadi, and Yejin Choi. 2019. <i>HellaSwag: Can a ma-</i>	474
415	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<i>chine really finish your sentence?</i> In <i>Proceedings of</i>	475
416	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	<i>the 57th Annual Meeting of the Association for Com-</i>	476
417	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	<i>putational Linguistics</i> , pages 4791–4800, Florence,	477
418	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Italy. Association for Computational Linguistics.	478
419	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang,	479
420	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,	480
421	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma,	481
422	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan	482
423	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023.	483
424	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	<i>GLM-130b: An open bilingual pre-trained model</i> . In	484
425	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	<i>The Eleventh International Conference on Learning</i>	485
426	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	<i>Representations (ICLR)</i> .	486
427	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor,	<b>A Alignment Skills Demonstration</b>	487
428	Adina Williams, Jian Xiang Kuan, Puxin Xu,	The FLASK framework annotates each instruction	488
429	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	with three skills that is needed to respond to the	489
430	Melanie Kambadur, Sharan Narang, Aurelien Rod-	instruction. We select 8 most frequent skills and	490
431	riguez, Robert Stojnic, Sergey Edunov, and Thomas	filter out instructions annotated with other skills, re-	491
432	Scialom. 2023. <i>Llama 2: Open foundation and fine-</i>	sulting 1,180 instructions in the evaluation set. The	492
433	<i>tuned chat models</i> .	following are demonstrations of each alignment	493
434	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa	skill from the annotation prompt.	494
435	Liu, Noah A. Smith, Daniel Khoshabi, and Hannaneh	<b>Logical Correctness</b> Is the final answer provided	495
436	Hajishirzi. 2023. <i>Self-instruct: Aligning language</i>	by the response logically accurate and correct for	496
437	<i>models with self-generated instructions</i> . In <i>Proceed-</i>	an instruction that has a deterministic answer?	497
438	<i>ings of the 61st Annual Meeting of the Association for</i>		
439	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,		
440	pages 13484–13508, Toronto, Canada. Association		
441	for Computational Linguistics.		
442	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-		
443	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva		
444	Naik, Arjun Ashok, Arut Selvan Dhanasekaran,		

498 **Factuality** Did the model extract pertinent and  
499 accurate background knowledge without any mis-  
500 information when factual knowledge retrieval is  
501 needed? Is the response supported by reliable evi-  
502 dence or citation of the source of its information?

503 **Commonsense Understanding** Is the model accu-  
504 rately interpreting world concepts for instructions  
505 that require a simulation of the expected result or  
506 necessitate commonsense or spatial reasoning?

507 **Comprehension** Does the response fulfill the re-  
508 quirements of the instruction by providing relevant  
509 information especially when the instruction is com-  
510 plex and includes multiple requirements? This in-  
511 cludes responding in accordance with the explicit  
512 and implicit purpose of given instruction.

513 **Completeness** Does the response provide a suf-  
514 ficient explanation? Comprehensiveness and thor-  
515 oughness of the response should be considered,  
516 which depends on the breadth of topics covered  
517 and the level of detail provided within each topic.

518 **Insightfulness** Is the response creative, original  
519 or novel, including new perspectives or interpreta-  
520 tions of existing information?

521 **Readability** Is the response structured to pro-  
522 mote readability and coherence? Does the response  
523 exhibit excellent organization?

524 **Conciseness** Is the response presented in a con-  
525 cise manner for the reader without any unnecessary  
526 information?

527 For how a response corresponds to a specific  
528 level of an alignment skill and other details, please  
529 refer to their repository <sup>1</sup>.

---

<sup>1</sup><https://github.com/kaistAI/FLASK>