

Demystifying Instruction Mixing for Fine-tuning Large Language Models

Anonymous ACL submission

Abstract

Instruction tuning significantly enhances the performance of large language models (LLMs) across various tasks. However, the procedure to optimizing the mixing of instruction datasets for LLM fine-tuning is still poorly understood. This study categorizes instructions into three primary types: NLP downstream tasks, coding, and general chat. We explore the effects of instruction tuning on different combinations of datasets on LLM performance, and find that certain instruction types are more advantageous for specific applications but can negatively impact other areas. This work provides insights into instruction mixtures, laying the foundations for future research.¹

1 Introduction

Instruction tuning has been shown to have surprising efficacy for aligning large language models (LLMs) with human instructions (Chung et al., 2022; Li et al., 2023; Wu et al., 2023; Xu et al., 2023; Touvron et al., 2023; Muennighoff et al., 2023a; Gunasekar et al., 2023). Recent studies highlight the diverse ways in which instructions can enhance the different capabilities of LLMs. For instance, using general-purpose, chat-like instructions can improve the performance of LLMs as chat assistants (Chiang et al., 2023; Ouyang et al., 2022; Taori et al., 2023; Ding et al., 2023), while training LLMs on instructions based off NLP tasks improves their performance on NLP benchmarks (Sanh et al., 2022; Chung et al., 2022; Muennighoff et al., 2023b), and incorporating coding instructions enhances LLM code generation (Fu and Khot, 2022; Gunasekar et al., 2023). However, a key unresolved issue is determining how to combine various instruction datasets to optimize overall LLM performance.

In this paper, we aim to better understand the impact of instruction mixing across three critical

¹Code and data are available at: [ANONYMIZED](#).

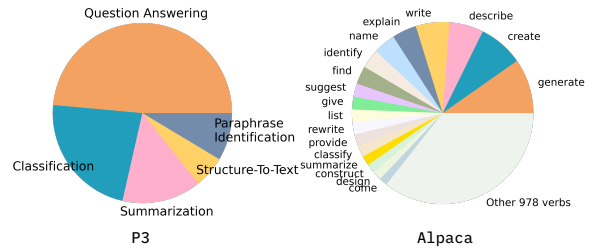


Figure 1: Instruction type distribution of P3 and Alpaca. For P3, the statistics come from the original dataset, while for Alpaca, we use a dependency parsing approach to extract the root verb of each instruction.

areas: NLP downstream tasks, coding, and chat. The core of our investigation revolves around understanding the influence of instruction dataset distributions on model performance in these different areas. We first select representative instruction datasets: P3 (Sanh et al., 2022) for NLP downstream tasks, CodeAlpaca (Chaudhary, 2023) for code generation, and Alpaca (Taori et al., 2023) for general-purpose instructions. As shown in Figure 1, P3 is focused primarily on five tasks (including QA and classification), whereas Alpaca contains a vast array of instructions. Using a dependency parser, we identify over 1K unique root verbs from Alpaca’s instructions, with *generate*, *create*, and *describe* being the most frequent. CodeAlpaca, by contrast, is exclusively focused on coding tasks, and exhibits less variation compared to the others as exemplified in Table 3. We fine-tune models across all eight potential combinations of these instruction datasets, and carry out detailed evaluation of model performance in terms of NLP downstream tasks, coding proficiency, and chat capabilities.

Our main contribution in this work is to shed light on instruction mixing when fine-tuning LLMs through comprehensive experimentation. Our findings can be summarized as follows:

- Specific instruction datasets enhance LLM performance in their respective task areas.

068	However, combining all instruction types does	117
069	not uniformly improve performance across all	118
070	tasks.	119
071	• Instructions reformulated from NLP down-	120
072	stream tasks (such as P3) can negatively im-	121
073	act the model’s conversational abilities. In	122
074	contrast, instructions focused on coding not	
075	only improve coding proficiency but also en-	
076	hance chat capabilities.	
077	• Larger models, with their increased capacity,	
078	are able to make more effective use of a di-	
079	verse range of instructions.	
080	2 Related Work	
081	Recent work has demonstrated that vanilla LLMs	
082	can follow general instructions if tuned with in-	
083	structions and corresponding responses (Mishra	
084	et al., 2022; Sanh et al., 2022; Wang et al., 2022).	
085	For instance, Sanh et al. (2022) crafted an instruc-	
086	tional dataset by reformulating supervised datasets	
087	with various prompts to create P3. However, de-	
088	spite their effectiveness in NLP tasks, these LLMs	
089	often diverge from human-like interactions in chat-	
090	bot applications.	
091	To facilitate general-purpose LLM fine-tuning,	
092	researchers has create general-purpose instruc-	
093	tional data by human annotation (Conover et al.,	
094	2023) and automatic approaches (Wang et al.,	
095	2023b; Taori et al., 2023). Recent work has fur-	
096	ther expanded the dataset size (Wu et al., 2023),	
097	language coverage (Li et al., 2023), and task types	
098	(Chaudhary, 2023; Yue et al., 2023).	
099	With increasing capabilities of LLMs and avail-	
100	ability of instruction datasets, researchers have	
101	aimed to imbue a single model with diverse ca-	
102	pabilities. Sengupta et al. (2023) attempted to	
103	blend different instruction datasets without con-	
104	sidering the data volume and task types. Longpre	
105	et al. (2023) suggested that increasing the num-	
106	ber of tasks and instruction diversity can enhance	
107	performance. In contrast, Anand et al. (2023) ex-	
108	cluded P3 from their fine-tuning dataset, seemingly	
109	to enhance alignment. Nevertheless, none of these	
110	papers systematically studied the impact of the in-	
111	struction mixture on the resulting LLM.	
112	Concurrent to our work, Wang et al. (2023a)	
113	fine-tuned LLaMA models on 12 instruction-tuning	
114	datasets separately. By evaluating those model on 7	
115	tasks, they found that different datasets can enhance	
116	model performance on individual tasks. They fur-	
	ther identified the optimal dataset combination, and	117
	trained a single model to achieve the best overall	118
	performance. Novel to this work, we classify the	119
	instructions and model skills into three types, and	120
	conduct a deep analysis of the influence of data	121
	mixture on the models.	122
	3 Experimental Setup	123
	Datasets We select Alpaca (Taori et al., 2023) as	124
	the <i>general instruction dataset</i> to align models, in	125
	the form of 52K instruction–response pairs. We	126
	use P3 (Sanh et al., 2022) as our <i>NLP task instruc-</i>	127
	<i>tion dataset</i> , which is reformatted for a wide range	128
	of NLP downstream tasks using diverse human-	129
	written templates. Since the number of samples in	130
	each task varies vastly, we randomly sample 1K	131
	instances from each subtask formatted with sev-	132
	eral corresponding prompts for diversity, result-	133
	ing in 660K samples. For <i>coding data</i> , we choose	134
	CodeAlpaca (Chaudhary, 2023), which is an in-	135
	struction dataset focusing on code generation. It	136
	contains 20K samples in different programming	137
	languages. To ensure a balanced comparison, we	138
	randomly sample a 20K subset from each dataset.	139
	Examples are provided in Table 3 in the Appendix.	140
	Evaluation We divide the evaluation into three	141
	parts: NLP benchmark performance, code gen-	142
	eration, and alignment evaluation (i.e., chat abil-	143
	ity evaluation). For NLP benchmarks, we use	144
	ARC (Clark et al., 2018), Winogrande (Sakaguchi	145
	et al., 2021), PIQA (Bisk et al., 2020), MMLU	146
	(Hendrycks et al., 2020), RACE (Lai et al., 2017),	147
	and HellaSwag (Zellers et al., 2019). For coding,	148
	we use HumanEval (Chen et al., 2021), which tests	149
	the pass rate of the generate codes. For alignment	150
	evaluation, we use the FLASK (Ye et al., 2023)	151
	framework to score model alignment. We keep	152
	the eight most frequent alignment skills from the	153
	original evaluation set, resulting in 1,180 samples.	154
	Then we employ GPT-4 to assess model responses	155
	to each instruction sample based on human-written	156
	principles. See Appendix B for details of these	157
	skills.	158
	Models We fine-tune LLaMA-2 7B and 13B	159
	(Touvron et al., 2023) models for two epochs in	160
	a generative way as in Radford et al. (2018), using	161
	a linear scheduler with a 3% warmup rate and a	162
	batch size of 64. The maximum learning rate is	163
	5×10^{-5} . The resources for training and evalua-	164
	tions are detailed in Appendix C.	165

Model	Data	ARC (challenge)	Wino-grande	PIQA	MMLU	Race	Hella-Swag	Average	HumanEval @1	HumanEval @10
LLaMA-2-7B	None	43.1	69.5	78.0	40.8	39.2	57.2	54.6	13.7	21.3
	A	47.8	67.6	78.2	42.2	<u>44.5</u>	61.1	56.9	13.5	17.1
	C	46.1	69.5	<u>78.5</u>	41.0	41.1	<u>61.0</u>	56.2	16.2	<u>24.4</u>
	P	<u>49.6</u>	71.4	79.0	46.0	43.5	59.4	58.2	4.6	7.9
	AC	47.1	66.9	78.1	40.4	44.2	59.7	56.1	17.5	25.0
	AP	48.4	70.0	78.1	43.8	42.9	58.5	56.9	13.8	17.7
	CP	48.0	<u>71.3</u>	78.4	<u>44.9</u>	44.4	60.7	<u>57.9</u>	<u>16.8</u>	20.1
	ACP	49.7	68.0	77.9	43.5	44.6	58.7	57.1	16.0	23.8
LLaMA-2-13B	None	48.6	71.9	79.2	52.1	40.7	60.1	58.8	15.4	26.2
	A	54.1	71.2	80.0	47.9	47.1	65.6	61.0	15.1	20.7
	C	49.7	73.4	80.8	<u>51.5</u>	45.4	63.6	60.7	17.9	24.4
	P	54.3	<u>74.2</u>	80.0	50.3	<u>45.6</u>	62.5	<u>61.1</u>	0.3	1.8
	AC	51.6	68.8	<u>80.6</u>	48.7	44.4	63.0	<u>59.5</u>	17.1	<u>27.4</u>
	AP	<u>54.8</u>	71.7	80.3	51.2	45.2	62.7	61.0	8.3	14.6
	CP	55.4	74.6	80.5	51.4	<u>45.6</u>	<u>63.9</u>	61.9	18.2	25.0
	ACP	54.4	71.5	80.0	50.0	47.1	63.1	61.0	20.2	32.9

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.

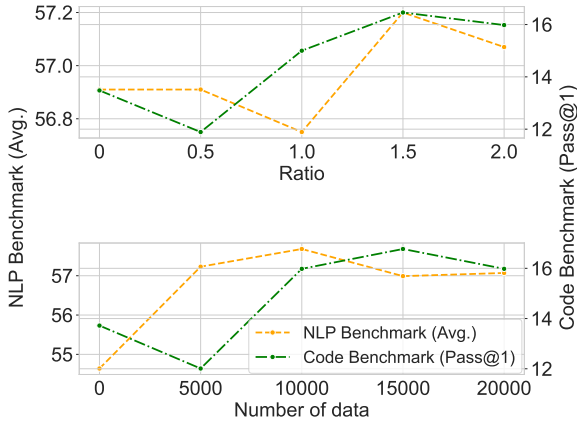


Figure 2: NLP benchmark scores (avg) and Code benchmark (HumanEval) scores for LLaMA-2-7B tuned with different mixing ratios and different numbers of instances. We keep the number of Alpaca instances constant at 20K and change the number of P3 and CodeAlpaca instances to get different ratios.

4 Results

For the remainder of this paper, we denote the 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.

4.1 NLP Tasks and Code Benchmark Results

Table 1 shows the zero-shot results on the NLP and code generation benchmarks. Predictably *each specialized instruction dataset improves 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 a 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 the story completion task, similar to the Alpaca instruction format. On the other hand, P3 fine-tuned models perform well on ARC and Winogrande, which involve multiple-choice QA and cloze tests, which are well represented in P3.

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 the code benchmarks, *incorporating general instructions 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 perform best with the ACP mixture, while the 7B models perform best with AC. This

Model	Data	Corr.	Fact.	Comm.	Compr.	Compl.	Insight.	Read.	Conc.	Avg.
LLaMA-2-7B	A	47.6	55.4	58.8	<u>54.8</u>	<u>48.0</u>	50.4	88.0	81.6	<u>60.6</u>
	C	48.8	52.0	58.4	<u>52.0</u>	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>	59.6	56.4	48.2	<u>49.8</u>	<u>86.6</u>	85.6	61.2
	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	50.4	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>	55.2	89.2	<u>84.0</u>	<u>64.0</u>
	C	57.2	<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>	61.0	66.6	61.2	51.4	<u>54.0</u>	<u>88.4</u>	86.6	65.6
	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: logical correctness, factuality, commonsense understanding, comprehension, completeness, insightfulness, readability, and conciseness, as well as average scores. Since the vanilla model cannot follow instructions, we exclude it from this table. The best result is in **bold**, and the second best result is underlined.

suggests that larger models can better learn from varied instructions more effectively than smaller models.

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

Mixing with Different Ratios While it is clear that mixing specialized instructions is vital for benchmark performance, how the mixing ratio correlates with the performance is also important. As Figure 2 shows, with the number of general instructions fixed to 20K, scores in both NLP task and code benchmarks first decrease and then increase as the ratio of specialized instructions increases. They both peak when the ratio is 1.5, and drop back slightly when the ratio is increased further to 2.0. We hypothesize that this is because the model overfits to the specialized instructions when there are too many such instructions.

Number of instances Figure 2 also shows the performance change with respect to the number of fine-tuning data instances. We mix each type of instruction with the same number. We find that the performance over both benchmarks plateaus when the number of instances is larger than 10K.

4.2 Alignment Skill Results

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

From Table 2 we make the following observations: (1) *All three types of instructions improve*

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 notably 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 due to better compression, commonsense understanding, completeness, and conciseness. (3) *Mixing P3 data causes a drop of -2.8 (7B) and -3.6 (13B) in alignment skills*, suggesting that P3 has 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 against diverse benchmarks and alignment skills. We find that general instructions provide better alignment as well as performance on NLP benchmarks, code instructions improve coding and alignment skills, while NLP task instructions hinder alignment skills when combined with other instruction types.

Limitations

Our work is subject to several limitations that should be addressed in future research. (1) We only use LLaMA-2 7B and 13B models in our experiments. Other models of varying sizes should be used to further verify our findings. We acknowl-

273	edge that the model’s behavior may vary with dif-		
274	ferent sizes, and that usually, larger models have		
275	stronger capabilities, and hence may be able to han-		
276	dle more instructions without performance degrada-		
277	tion. (2) In this paper, we limit our instruction		
278	dataset to 20K and mainly compare the 1:1 ratio of		
279	all instruction types. We leave the exploration of		
280	the impact of more instructions and mixing ratios		
281	to future work.		
282	We acknowledge these limitations and propose		
283	that future work should focus on addressing them		
284	to help the community better understand the impact		
285	of instruction mixture on LLMs.		
286	References		
287	Yuvanesh Anand, Zach Nussbaum, Brandon Duder-		
288	stadt, Benjamin Schmidt, and Andriy Mulyar. 2023.		
289	GPT4All: Training an assistant-style chatbot with		
290	large scale data distillation from GPT-3.5-Turbo.		
291	https://github.com/nomic-ai/gpt4all .		
292	Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi,		
293	et al. 2020. PIQA: Reasoning about physical com-		
294	monsense in natural language. In <i>Proceedings of</i>		
295	<i>the AAAI conference on Artificial Intelligence</i> , pages		
296	7432–7439.		
297	Sahil Chaudhary. 2023. Code Alpaca: An instruction-		
298	following LLaMA model for code generation.		
299	https://github.com/sahil280114/codealpaca .		
300	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming		
301	Yuan, Henrique Ponde de Oliveira Pinto, Jared Ka-		
302	plan, Harri Edwards, Yuri Burda, Nicholas Joseph,		
303	Greg Brockman, et al. 2021. Evaluating large		
304	language models trained on code. <i>arXiv preprint</i>		
305	<i>arXiv:2107.03374</i> .		
306	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,		
307	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan		
308	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion		
309	Stoica, and Eric P. Xing. 2023. Vicuna: An open-		
310	source chatbot impressing GPT-4 with 90%* Chat-		
311	GPT quality.		
312	Hyung Won Chung, Le Hou, Shayne Longpre, Barret		
313	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi		
314	Wang, Mostafa Dehghani, Siddhartha Brahma, Al-		
315	bert Webson, Shixiang Shane Gu, Zhuyun Dai,		
316	Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-		
317	ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,		
318	Dasha Valter, Sharan Narang, Gaurav Mishra, Adams		
319	Yu, Vincent Zhao, Yanping Huang, Andrew Dai,		
320	Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-		
321	cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,		
322	and Jason Wei. 2022. <i>Scaling instruction-finetuned</i>		
323	<i>language models</i> .		
324	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,		
325	Ashish Sabharwal, Carissa Schoenick, and Oyvind		
	Tafjord. 2018. Think you have solved question an-		326
	swering? try Arc, the AI2 reasoning challenge. <i>arXiv</i>		327
	<i>preprint arXiv:1803.05457</i> .		328
	Mike Conover, Matt Hayes, Ankit Mathur, Jianwei		329
	Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wen-		330
	dell, Matei Zaharia, and Reynold Xin. 2023. Free		331
	Dolly: Introducing the world’s first truly open		332
	instruction-tuned LLM. https://github.com/		333
	databricks/dolly .		334
	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi		335
	Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun,		336
	and Bowen Zhou. 2023. <i>Enhancing chat language</i>		337
	<i>models by scaling high-quality instructional conver-</i>		338
	<i>sations</i> .		339
	Hao Fu, Yao; Peng and Tushar Khot. 2022. <i>How does</i>		340
	<i>GPT obtain its ability? tracing emergent abilities of</i>		341
	<i>language models to their sources. Yao Fu’s Notion</i> .		342
	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio		343
	César Teodoro Mendes, Allie Del Giorno, Sivakanth		344
	Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo		345
	de Rosa, Olli Saarikivi, Adil Salim, Shital Shah,		346
	Harkirat Singh Behl, Xin Wang, Sébastien Bubeck,		347
	Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and		348
	Yuanzhi Li. 2023. <i>Textbooks are all you need</i> .		349
	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,		350
	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.		351
	2020. Measuring massive multitask language under-		352
	standing. In <i>International Conference on Learning</i>		353
	<i>Representations</i> .		354
	Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang,		355
	and Eduard Hovy. 2017. <i>RACE: Large-scale ReAd-</i>		356
	<i>ing comprehension dataset from examinations</i> . In		357
	<i>Proceedings of the 2017 Conference on Empirical</i>		358
	<i>Methods in Natural Language Processing</i> , pages 785–		359
	794, Copenhagen, Denmark. Association for Compu-		360
	tational Linguistics.		361
	Haonan Li, Fajri Koto, Minghao Wu, Alham Fikri Aji,		362
	and Timothy Baldwin. 2023. <i>Bactrian-X: Multilin-</i>		363
	<i>gual replicable instruction-following models with</i>		364
	<i>low-rank adaptation</i> .		365
	Shayne Longpre, Le Hou, Tu Vu, Albert Webson,		366
	Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V. Le,		367
	Barret Zoph, Jason Wei, and Adam Roberts. 2023.		368
	<i>The Flan collection: Designing data and methods for</i>		369
	<i>effective instruction tuning</i> .		370
	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and		371
	Hannaneh Hajishirzi. 2022. <i>Cross-task generaliza-</i>		372
	<i>tion via natural language crowdsourcing instructions</i> .		373
	In <i>Proceedings of the 60th Annual Meeting of the</i>		374
	<i>Association for Computational Linguistics (Volume</i>		375
	<i>1: Long Papers)</i> , pages 3470–3487, Dublin, Ireland.		376
	Association for Computational Linguistics.		377
	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,		378
	Adam Roberts, Stella Biderman, Teven Le Scao,		379

380	M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hai-	An instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca .	438
381	ley Schoelkopf, Xiangru Tang, Dragomir Radev, Al-		439
382	ham Fikri Aji, Khalid Almubarak, Samuel Albanie,		
383	Zaid Alyafeai, Albert Webson, Edward Raff, and	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	440
384	Colin Raffel. 2023a. Crosslingual generalization	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	441
385	through multitask finetuning . In <i>Proceedings of the</i>	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	442
386	<i>61st Annual Meeting of the Association for Com-</i>	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	443
387	<i>putational Linguistics (Volume 1: Long Papers)</i> ,	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	444
388	pages 15991–16111, Toronto, Canada. Association	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	445
389	for Computational Linguistics.	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	446
		thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	447
390	Niklas Muennighoff, Thomas Wang, Lintang Sutawika,	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	448
391	Adam Roberts, Stella Biderman, Teven Le Scao,	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	449
392	M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hai-	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	450
393	ley Schoelkopf, Xiangru Tang, Dragomir Radev, Al-	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	451
394	ham Fikri Aji, Khalid Almubarak, Samuel Albanie,	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	452
395	Zaid Alyafeai, Albert Webson, Edward Raff, and	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	453
396	Colin Raffel. 2023b. Crosslingual generalization	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	454
397	through multitask finetuning .	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	455
		nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	456
398	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	457
399	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	458
400	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	Melanie Kambadur, Sharan Narang, Aurelien Rod-	459
401	2022. Training language models to follow instruc-	riguez, Robert Stojnic, Sergey Edunov, and Thomas	460
402	tions with human feedback. <i>Advances in Neural</i>	Scialom. 2023. LLaMA 2: Open foundation and	461
403	<i>Information Processing Systems</i> , 35:27730–27744.	fine-tuned chat models .	462
404	Alec Radford, Karthik Narasimhan, Tim Sali-	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack	463
405	mans, and Ilya Sutskever. 2018. Improving	Hessel, Tushar Khot, Khyathi Raghavi Chandu,	464
406	language understanding by generative pre-training.	David Wadden, Kelsey MacMillan, Noah A. Smith,	465
407	http://openai-assets.s3.amazonaws.com/	Iz Beltagy, and Hannaneh Hajishirzi. 2023a. How	466
408	research-covers/language-unsupervised/	far can camels go? exploring the state of instruction	467
409	language_understanding_paper.pdf .	tuning on open resources .	468
410	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavat-	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa	469
411	ula, and Yejin Choi. 2021. Winogrande: An adver-	Liu, Noah A. Smith, Daniel Khoshabi, and Hannaneh	470
412	sarial Winograd schema challenge at scale. <i>Commu-</i>	Hajishirzi. 2023b. Self-instruct: Aligning language	471
413	<i>nications of the ACM</i> , 64(9):99–106.	models with self-generated instructions . In <i>Proceed-</i>	472
		<i>ings of the 61st Annual Meeting of the Association for</i>	473
414	Victor Sanh, Albert Webson, Colin Raffel, Stephen H	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	474
415	Bach, Lintang Sutawika, Zaid Alyafeai, Antoine	pages 13484–13508, Toronto, Canada. Association	475
416	Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja,	for Computational Linguistics.	476
417	et al. 2022. Multitask prompted training enables		
418	zero-shot task generalization. In <i>ICLR 2022-Tenth</i>	Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-	477
419	<i>International Conference on Learning Representa-</i>	labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva	478
420	<i>tions</i> .	Naik, Arjun Ashok, Arut Selvan Dhanasekaran,	479
		Anjana Arunkumar, David Stap, Eshaan Pathak,	480
421	Neha Sengupta, Sunil Kumar Sahu, Bokang Jia,	Giannis Karamanolakis, Haizhi Lai, Ishan Puro-	481
422	Satheesh Katipomu, Haonan Li, Fajri Koto, William	hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,	482
423	Marshall, Gurpreet Gosal, Cynthia Liu, Zhiming	Krima Doshi, Kuntal Kumar Pal, Maitreya Patel,	483
424	Chen, Osama Mohammed Afzal, Samta Kamboj,	Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,	484
425	Onkar Pandit, Rahul Pal, Lalit Pradhan, Zain Muham-	Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,	485
426	ammad Mujahid, Massa Baali, Xudong Han, Sondo-	Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,	486
427	s dos Mahmoud Bsharat, Alham Fikri Aji, Zhiqiang	Shailaja Keyur Sampat, Siddhartha Mishra, Sujan	487
428	Shen, Zhengzhong Liu, Natalia Vassilieva, Joel Hes-	Reddy A, Sumanta Patro, Tanay Dixit, and Xudong	488
429	tness, Andy Hock, Andrew Feldman, Jonathan Lee,	Shen. 2022. Super-NaturalInstructions: Generaliza-	489
430	Andrew Jackson, Hector Xuguang Ren, Preslav	tion via declarative instructions on 1600+ NLP tasks .	490
431	Nakov, Timothy Baldwin, and Eric Xing. 2023.	In <i>Proceedings of the 2022 Conference on Empiri-</i>	491
432	Jais and Jais-chat: Arabic-centric foundation and	<i>cal Methods in Natural Language Processing</i> , pages	492
433	instruction-tuned open generative large language	5085–5109, Abu Dhabi, United Arab Emirates. As-	493
434	models .	sociation for Computational Linguistics.	494
435	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	Minghao Wu, Abdul Waheed, Chiyu Zhang, Muham-	495
436	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	ammad Abdul-Mageed, and Alham Fikri Aji. 2023.	496
437	and Tatsunori B. Hashimoto. 2023. Stanford Alpaca:		

497	LaMini-LM: A diverse herd of distilled models from large-scale instructions.	Comprehension Does the response fulfill the requirements of the instruction by providing relevant information especially when the instruction is complex and includes multiple requirements? This includes responding in accordance with the explicit and implicit purpose of given instruction.	542
498			543
499	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. WizardLM: Empowering large language models to follow complex instructions. <i>arXiv preprint arXiv:2304.12244</i> .		544
500			545
501			546
502			547
503			
504	Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, James Thorne, Juho Kim, and Minjoon Seo. 2023. FLASK: Fine-grained language model evaluation based on alignment skill sets.	Completeness Does the response provide a sufficient explanation? Comprehensiveness and thoroughness of the response should be considered, which depends on the breadth of topics covered and the level of detail provided within each topic.	548
505			549
506			550
507			551
508			552
509	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2023. MAMmoTH: Building math generalist models through hybrid instruction tuning.	Insightfulness Is the response creative, original or novel, including new perspectives or interpretations of existing information?	553
510			554
511			555
512			
513	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4791–4800, Florence, Italy. Association for Computational Linguistics.	Readability Is the response structured to promote readability and coherence? Does the response exhibit excellent organization?	556
514			557
515			558
516		Conciseness Is the response presented in a concise manner for the reader without any unnecessary information?	559
517			560
518		The prompt for alignment skill assessment are provided in Figure 3. For how a response corresponds to a specific level of an alignment skill and other details, please refer to their repository. ²	561
519	A Examples of Instruction Types		562
520	Table 3 list examples from Alpaca, CodeAlpaca, and P3 dataset.		563
521			564
522		C Resources	565
523	B Alignment Skills Demonstration	We use 4×A100 to train LLaMA-2-7B and 8×A100 to train LLaMA-2-13B. Each 20k data takes about 2 hours. For all experiments, training takes about 288 A100 GPU hours.	566
524	The FLASK framework annotates each instruction with three skills that is needed to respond to the instruction. We select 8 most frequent skills and filter out instructions annotated with other skills, resulting 1,180 instructions in the evaluation set. The following are demonstrations of each alignment skill from the annotation prompt.		567
525		For evaluation, we use GPT-4, where each input has an average of 950 tokens and each output has an average of 293 tokens. All evaluations cost about \$760.	568
526			569
527			570
528			571
529			572
530	Logical Correctness Is the final answer provided by the response logically accurate and correct for an instruction that has a deterministic answer?		573
531			574
532			
533	Factuality Did the model extract pertinent and accurate background knowledge without any misinformation when factual knowledge retrieval is needed? Is the response supported by reliable evidence or citation of the source of its information?		
534			
535			
536			
537			
538	Commonsense Understanding Is the model accurately interpreting world concepts for instructions that require a simulation of the expected result or necessitate commonsense or spatial reasoning?		
539			
540			
541			

²<https://github.com/kaistAI/FLASK>

Source	Example
Alpaca	<p>Instruction: List 5 reasons why someone should learn to code.</p> <p>Response: 1. High demand for coding skills in the job market \n2. Increased problem-solving and analytical skills \n3. Ability to develop new products and technologies \n4. Potentially higher earning potential \n5. Opportunity to work remotely and/or freelance</p>
CodeAlpaca	<p>Instruction: Write a function to find the number of distinct states in a given matrix.</p> <p>Input: matrix = [[1, 0, 0],\n [1, 0, 1],\n [1, 1, 1]]</p> <p>Response: def find_num_distinct_states(matrix):\n states = set()\n for row in matrix:\n state = "".join([str(x) for x in row])\n states.add(state)\n return len(states)</p>
P3	<p>Instruction: Answer the following question.\nAnna Kournikova, Michelangelo, ILOVEYOU, Melissa, and Stuxnet are all examples of what?</p> <p>Response: Computer virus/worm</p>

Table 3: Examples from Alpaca, CodeAlpaca, and P3.

[System]

We would like to request your feedback on the performance of the response of the assistant to the user instruction displayed below. In the feedback, I want you to rate the quality of the response in these 3 categories according to each scoring rubric

[Skill 1 definition](#)
[Skill 1 scoring principles](#)

[Skill 2 definition](#)
[Skill 2 scoring principles](#)

[Skill 3 definition](#)
[Skill 3 scoring principles](#)

[Instruction]
instruction

[Ground Truth Answer]
ground truth answer

[Assistant’s Response]
response for evaluation

[The End of Assistant’s Response]

Please give feedback on the assistant’s responses. Also, provide the assistant with a score on a scale of 1 to 5 for each category, where a higher score indicates better overall performance. Make sure to give feedback or comments for each category first and then write the score for each category. Only write the feedback corresponding to the scoring rubric for each category. The scores of each category should be orthogonal, indicating that ‘Efficiency of User Alignment’ should not be considered for ‘Readability of User Alignment’ category, for example.

Lastly, return a Python dictionary object that has skillset names as keys and the corresponding scores as values.

[System]

Figure 3: Alignment skill assessment prompt (from FLASK (Ye et al., 2023)). The blue parts are filled by corresponding content.