

Rethinking the Instruction Quality: LIFT is What You Need

Anonymous ACL submission

Abstract

Instruction tuning, a specialized technique to enhance large language model (LLM) performance via instruction datasets, relies heavily on the quality of the employed data. Existing quality improvement methods alter instruction data through dataset expansion or curation. However, the expansion method introduces the risk of data deficiency and redundancy, potentially compromising the correctness and accuracy of the LLM’s knowledge, while the curation approach confines the LLM’s potential to the original dataset. Our aim is to surpass the original data quality without confronting these shortcomings. To achieve this, we propose **LIFT** (LLM Instruction Fusion Transfer), a novel and versatile paradigm designed to elevate the instruction quality to new heights. LIFT strategically broadens data distribution to encompass more high-quality subspaces and eliminates redundancy, concentrating on high-quality segments across overall data subspaces. Experimental results demonstrate that, even with a limited quantity of high-quality instruction data selected by our paradigm, LLMs not only consistently uphold robust performance across natural language understanding and code generation tasks but also surpass many state-of-the-art results, highlighting the significant improvement in instruction quality achieved by our paradigm.

1 Introduction

In recent years, Large Language Models (LLMs) have gained prominence for their remarkable effectiveness in natural language comprehension tasks (OpenAI, 2023; Yang et al., 2023; Qi et al., 2023). High-quality pretrained LLMs are readily available, facilitating their customization for versatile applications (Wei et al., 2021; Huang et al., 2023). One popular fine-tuning approach, known as instruction tuning (Wei et al., 2022; Ouyang et al., 2022), involves fine-tuning pre-trained LLMs using datasets accompanied by natural language instructions. Its

relative simplicity and affordability make it a preferred method for improving LLMs’ performance on specific tasks.

The quality of current instruction datasets, whether manually crafted or generated by LLMs, often falls short of the desired standard. Human-crafted datasets depend on human annotators to generate a substantial corpus with human instructions, resulting in a lack of detailed context and explanation within the instruction dataset. Additionally, these datasets may contain vague or subjective descriptions. On the other hand, LLM-generated datasets utilize advanced LLMs to generate or complete instructions and responses but lack supervision regarding the diversity and quality of the generated data.

The concern surrounding the quality of instruction datasets has prompted researchers to explore methods aimed at enhancing their overall quality. Current approaches to instruction quality enhancement can be broadly categorized into two groups: data expansion and data curation. Data expansion methods involve leveraging advanced LLMs with a suitable prompt template to generate new instructions and corresponding answers based on the original dataset (Xu et al., 2023; Luo et al., 2023; Taori et al., 2023). On the other hand, data curation methods entail the meticulous selection of high-quality data from the original dataset, employing specific quality evaluation criteria (Zhou et al., 2023; Du et al., 2023).

However, both existing methods exhibit limitations that hinder their ability to further enhance performance. Expansion methods introduce redundancy into the dataset (Xu et al., 2023; Luo et al., 2023) as the newly generated instructions typically derive from the original ones. While the effectiveness of curation methods heavily relies on the quality of the original dataset, limiting the quality of the curated dataset (Du et al., 2023; Li et al., 2023a). These limitations necessitate a reliance on

specific expansion or curation strategies to achieve superior performance on certain benchmarks, at the expense of losing the ability to generalize the approach.

In this paper, we delve into the distribution of instruction quality to address the mentioned issues. We posit that both current methods essentially function as data distribution transfers: expansion enables the distribution to cover a broader range of data subspaces, typically characterized by higher quality, while curation concentrates the distribution on a higher-quality subset of the original dataset. Building on this perspective, we propose a novel paradigm for improving LLM instruction quality, termed **LIFT (LLM Instruction Fusion Transfer)**. LIFT is designed to amalgamate the advantages of data expansion and curation, mitigating their shortcomings to generate a diverse and high-quality dataset while significantly reducing quantity. Our paradigm consists of two phases. Firstly, we employ "Dataset Distribution Expansion", broadening the data distribution to cover more high-quality subspaces. Then, we utilize "Dataset Variety and Quality Curation" to eliminate redundancy and densify the data distribution, focusing on the high-quality segments of overall data subspaces. The data distribution transfer patterns of three methods are described in Fig. 1.

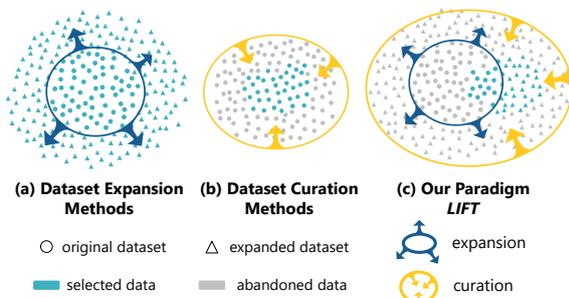


Figure 1: Different Instruction Data Distribution Transfer Patterns.

To validate the effectiveness of LIFT, we employ the finally curated instructions for fine-tuning open-source LLMs. Through extensive experiments evaluating the performance of these fine-tuned LLMs in both natural language understanding (NLU) tasks and code generation tasks, the results consistently demonstrate that LLMs achieve robust SOTA or nearly-SOTA performance even with a limited quantity of high-quality instruction data. Furthermore, they even outperform models trained on larger datasets on certain benchmarks.

To summarize, our main contribution are:

- We propose a highly effective and versatile paradigm, LIFT, which challenges the conventional single-mode enhancement for instruction datasets. LIFT rethinks data quality by focusing on data distribution transfer. It aims to elevate the quality of the instruction dataset to new heights, overcoming redundancy and quality limitations present in current methods.
- Throughout the expansion and curation phases of the paradigm, we prioritize both variety and quality as essential goals for quality enhancement. Unlike existing works that concentrate on only one stage, we posit that considering these characteristics at both stages is crucial for incorporating more high-quality data.
- Our extensive experiments demonstrate that with a significantly reduced quantity of high-quality instructions selected by our paradigm, LLMs consistently achieve SOTA performance on many NLU tasks and code generation tasks. This provides valuable insights, suggesting that a selective approach based on the principles of data distribution transfer is not only more effective but also cost-effective compared to the indiscriminate feeding of large volumes of data.

2 Related Works

The current methods for enhancing instruction quality can be broadly categorized into two types based on how data is manipulated: dataset expansion and curation.

2.1 Instruction Dataset Expansion

The original instruction dataset often consists of concise and straightforward prompts, yielding simplistic responses with limited semantic information. To address this limitation, researchers have proposed using LLMs to expand these original instructions to introduce more high-quality data. Alpaca (Taori et al., 2023) suggested adopting the self-instruct method, utilizing ChatGPT to generate data for fine-tuning. Vicuna (Chiang et al., 2023) employed data collected from ShareGPT.com, a platform where users share their conversations with ChatGPT, for fine-tuning their models. WizardLM (Xu et al., 2023) and WizardCoder (Luo et al.,

169	2023) introduced the <i>Evol-Instruct</i> method, involving the evolution of existing instruction data to generate more diverse and complex data.	217
170		218
171		219
172	2.2 Instruction Dataset Curation	220
173	One challenge in the instruction tuning process arises from the observation that fine-tuning with larger expanded instruction datasets does not always guarantee better results, yet demanding more computational resources. To address this, some researchers have focused on filtering out low-quality data during the fine-tuning stage. LIMA (Zhou et al., 2023) demonstrates that fine-tuning a robust pre-trained language model on 1000 high-quality, human-curated examples can yield remarkable and competitive results. Instruction Mining (Cao et al., 2023) introduces a linear rule for selecting high-quality instruction data, eliminating the need for human annotation. Du et al. (2023) present a model-oriented data selection (MoDS) approach, which selects instruction data based on new criteria considering three aspects: quality, coverage, and necessity. Li et al. (2023a) introduce a self-guided methodology for LLMs to autonomously discern and select cherry-picked samples from vast open-source datasets, effectively minimizing manual curation and potential costs.	221
174		222
175		223
176		224
177		225
178		226
179		227
180		228
181		229
182		230
183		231
184		232
185		233
186		234
187		235
188		236
189		237
190		238
191		239
192		240
193		241
194		242
195	3 LLM Instruction Fusion Transfer	243
196	3.1 Data Distribution Transfer	244
197	Current methods for enhancing instruction quality, either through data expansion or curation, do enhance the original dataset to some extent. However, the effectiveness of these methods is constrained by inherent limitations. To scrutinize these limitations and explore innovative approaches to break from conventional enhancement modes, we propose a novel perspective for rethinking instruction data quality: data distribution transfer .	245
198		246
199		247
200		248
201		249
202		250
203		251
204		252
205		253
206		254
207		255
208		256
209		257
210		258
211		259
212		260
213		261
214		262
215		263
216		264
		265
		266
		267
		268
		269
		270
		271
		272
		273
		274
		275
		276
		277
		278
		279
		280
		281
		282
		283
		284
		285
		286
		287
		288
		289
		290
		291
		292
		293
		294
		295
		296
		297
		298
		299
		300
		301
		302
		303
		304
		305
		306
		307
		308
		309
		310
		311
		312
		313
		314
		315
		316
		317
		318
		319
		320
		321
		322
		323
		324
		325
		326
		327
		328
		329
		330
		331
		332
		333
		334
		335
		336
		337
		338
		339
		340
		341
		342
		343
		344
		345
		346
		347
		348
		349
		350
		351
		352
		353
		354
		355
		356
		357
		358
		359
		360
		361
		362
		363
		364
		365
		366
		367
		368
		369
		370
		371
		372
		373
		374
		375
		376
		377
		378
		379
		380
		381
		382
		383
		384
		385
		386
		387
		388
		389
		390
		391
		392
		393
		394
		395
		396
		397
		398
		399
		400
		401
		402
		403
		404
		405
		406
		407
		408
		409
		410
		411
		412
		413
		414
		415
		416
		417
		418
		419
		420
		421
		422
		423
		424
		425
		426
		427
		428
		429
		430
		431
		432
		433
		434
		435
		436
		437
		438
		439
		440
		441
		442
		443
		444
		445
		446
		447
		448
		449
		450
		451
		452
		453
		454
		455
		456
		457
		458
		459
		460
		461
		462
		463
		464
		465
		466
		467
		468
		469
		470
		471
		472
		473
		474
		475
		476
		477
		478
		479
		480
		481
		482
		483
		484
		485
		486
		487
		488
		489
		490
		491
		492
		493
		494
		495
		496
		497
		498
		499
		500

3.2 Paradigm LIFT

As described in Fig.2, our paradigm LIFT follows a two-stage structure. In both stages, we value the diversity and quality as the crucial criterion and we believe the "Dataset Distribution Expansion" and "Dataset Variety and Quality Curation" equally contribute to the quality enhancement.

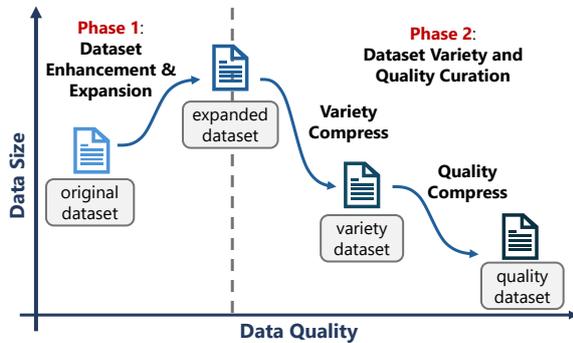


Figure 2: Instruction Dataset Curation Paradigm LIFT

3.2.1 Dataset Distribution Expansion

The goal of dataset distribution expansion is to encompass a more diverse and high-quality range of data within the distribution, while ensuring a certain distance from the original instructions. To achieve this, it is crucial to employ carefully designed instruction-generation prompts. Drawing inspiration from the instruction rewrite method proposed by Xu et al. (2023) and Luo et al. (2023), our approach focuses on generating diverse and intricate instructions. We guide GPT-4 to act as a prompt re-writer, generating challenging instructions based on specified generation rules. Considering the variation in content for NLU and code generation tasks within the instruction dataset, we configure distinct settings for GPT prompts to enhance complexity. For further details, refer to Appx.A. We iterate this process for k rounds, merging the expanded datasets with the original dataset to create the final expanded dataset.

3.2.2 Dataset Variety and Quality Curation

An effective curation method ought to eliminate duplicated or low-quality instructions from the original dataset, while preserving representative and high-quality ones. To meet this criterion, curation should be approached with meticulous attention to both variety and quality.

Variety Curation. Current variety curation typically involves clustering methods such as k-means or spectral clustering, which initially segment the

original data distribution into several small groups, followed by the selection of representative items from these groups (Du et al., 2023; Wei et al., 2023). We argue that this approach may lack generalizability and be less effective when dealing with new datasets. This is because these methods require prior knowledge of the number of clusters, and choosing cluster numbers that are either too large or too small may reduce their effectiveness in selecting representatives.

Our variety curation method takes another route, as depicted in Fig.3. Initially, GPT generates embeddings with 1536 dimensions for each item, which proves unwieldy for analysis. To address this, we aim to reduce the embedding dimension and devise a method to represent data differentiation. We achieve this by calculating the covariance matrix of the given features and performing eigenvalue decomposition on the covariance matrix to obtain eigenvalues and eigenvectors. We then choose the top k eigenvectors corresponding to the largest k eigenvalues, where k is the target reduced dimension. We set k to a value that preserves nearly 95% of the variance of the original embeddings, ensuring the retention of a significant amount of information.

This process allows us to analyze row variance on the dimension-reduced features to identify items with significant differences. Row variance measures variability in the reduced space, and high variance suggests substantial positional changes along that dimension, aiding in identifying diverse data points. We select items with the highest 20% row variances to construct the variety-curated dataset. Our method doesn't require any prior statistical knowledge of the dataset, making it versatile and effective for all tasks.

Quality Curation. Following variety curation is the quality curation phase, where we discern high-quality instruction data. Rating instruction quality is challenging due to the lack of official quantitative metrics. Employing professional annotators for scoring is impractical due to dataset size and costs. Therefore, we use GPT-4 as an instruction scorer, generating GPT quality scores across four dimensions: accuracy, explanation, clarity, and difficulty, with proportions based on their contributions to overall quality. The guiding template for GPT-4 scores is in Appx.B.

In our practical experience, we observed that GPT-4 consistently assigns high total scores to all

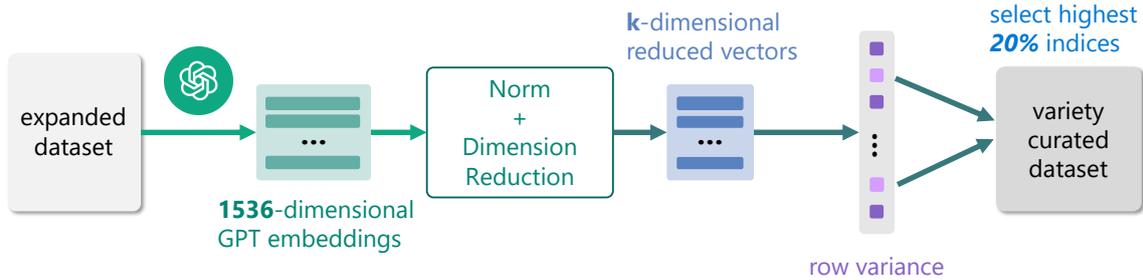


Figure 3: Variety Curation with Dimension Reduction and Row Variances

instruction data, posing a challenge by reducing differentiation in GPT quality scores. To address this, we implement the following steps for more reasonable and differentiated scores: first we instruct GPT-4 to provide a comprehensive rationale along with a score yields more reasonable results, preventing indiscriminate high-level assessments. Mandating GPT-4 to articulate its reasoning offers an additional self-checking opportunity. Secondly, before scoring instruction items, we present manually scored examples as guidelines. Offering three examples with scores representing poor, average, and high quality helps GPT-4 recognize low-quality data and understand how to appropriately score it.

Recognizing that total length indicates information richness, we also incorporate this into the quality score calculation. A positively correlated mapping function derives a lengthwise semantic score based on instruction data’s length. Combining GPT quality score and lengthwise semantic score produces the final quality score. High-quality scores compose the final quality-curated dataset, as illustrated in Fig.4. Appx.C presents total quality score distributions. Meticulous efforts ensure substantial differentiation in quality scores for precise identification and selection of high-quality data.

4 Experiments

To validate the effectiveness of our paradigm, we apply our method to two extensively studied tasks: Natural Language Understanding (NLU) tasks and Code Generation tasks, where we conduct comprehensive experiments to evaluate the performance of our paradigm.

4.1 Experiments Setup

4.1.1 Basic Foundation Models and Base Datasets

We adopt distinct foundation models and base dataset configurations for the two tasks under con-

sideration. In NLU tasks, we employ the SOTA foundation model Mistral 7B (Jiang et al., 2023), known for its exceptional performance relative to other 7B models and its ability to surpass larger models in specific benchmarks. Our base dataset for NLU tasks is the Open Platypus dataset (Lee et al., 2023), comprising 25k curated examples focused on enhancing LLMs’ STEM and logic knowledge. While in the realm of code generation tasks, we harness the capabilities of StarCoder 15B (Li et al., 2023b), a widely-utilized code LLM trained on a diverse array of sources encompassing over 80 programming languages, Git commits, GitHub issues, and Jupyter notebooks. Our base dataset, Code Alpaca (Chaudhary, 2023), consists of 20k instruction-following code instances for fine-tuning Code LLMs.

For comprehensive implementation details pertaining to instruction tuning in both tasks, please refer to Appx. D.

4.2 Benchmarks and Metrics

We have chosen six widely-recognized benchmarks spanning both tasks. In the domain of NLU tasks, we have incorporated HellaSwag, ARC Challenge, TruthfulQA, and MMLU. For code generation tasks, our selection encompasses HumanEval and MBPP. Detailed information about these benchmarks is provided in Appx. E.1.

For NLU tasks, we adopt accuracy as the metric, aligning with the methodology embraced by other researchers. This metric is calculated as the number of correct questions divided by the total number of questions.

In code generation tasks, our metric of choice is pass@k, defined in the same manner as by Chen et al. (2021). The formula for calculating pass@k is presented as:

$$pass@k := \mathbb{E}_{problems} \left[1 - \frac{C(n-c, k)}{C(n, k)} \right]$$

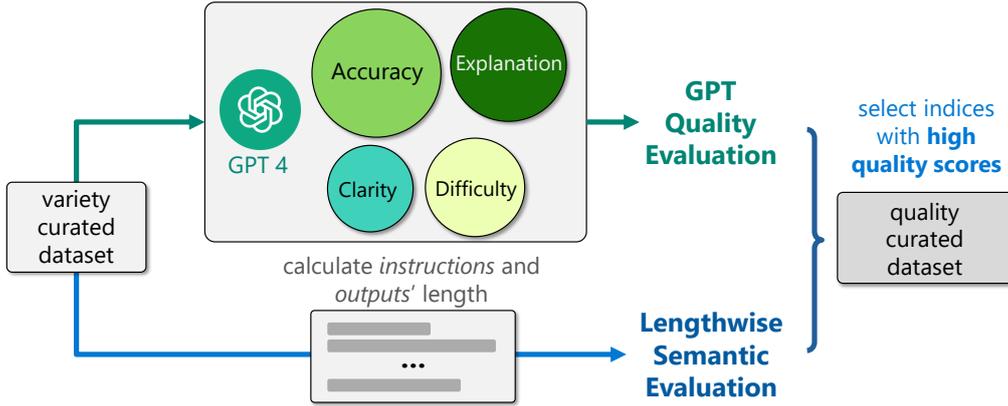


Figure 4: Quality Curation with GPT Quality Evaluation and Lengthwise Semantic Evaluation

Here, n represents the number of generated answers for each question, and c denotes the number of correct answers for each question. In our experiments, we specifically choose pass@1 as the designated metric.

4.3 Experiment Results

To validate the effectiveness of our paradigm, we conduct comparisons between models fine-tuned on LIFT’s final curated dataset and other SOTA pre-trained LLMs as well as instruction-tuned LLMs across both tasks. The details of the selected models for comparison are provided in Appx. E.2.

4.3.1 NLU Tasks

Tab.1 presents the NLU task comparison results. Notably, our final 7B instruction-tuned model consistently outperforms other 7B models on all benchmarks. Comparing with 13B models, our 7B model even outperforms in all benchmarks except TruthfulQA. With only 7 billion parameters and 15k instructions, significantly fewer than other instruction-tuned models, our model achieves the highest average benchmark score at 0.656.

4.3.2 Code Generation Tasks

As illustrated in Table 2, our paradigm’s fine-tuned model consistently outperforms most models in code generation tasks. Although our model trails the current state-of-the-art 15B model, WizardCoder, by approximately 2% on both benchmarks, it is noteworthy that our paradigm utilizes only about **one-eighth** of the instruction data employed by WizardCoder. Considering the disparity in the size of the instruction dataset, our paradigm demonstrates robust performance, highlighting its capability to achieve performance levels close to the

state-of-the-art with a significantly smaller amount of data.

We also compared our paradigm’s final curated dataset with a randomly selected dataset of the same size in both tasks. The results demonstrate that merely reducing the dataset quantity, without accounting for the diversity and quality in the perspective of data distribution, does not lead to performance improvement. These experiments affirm our paradigm’s versatile effectiveness in NLU and code generation tasks. The paradigm excels in generating diverse, high-quality data, leveraging it in the instruction-tuning process to achieve SOTA or near-SOTA performance.

4.4 Paradigm Ablation Experiments Results

Our paradigm ablation experiment begins with the original base dataset serving as the input for LIFT. Subsequently, we generate the expanded dataset, variety-curated dataset, and the quality-curated dataset. These datasets are then utilized for fine-tuning the basic foundation models. We assess the benchmark performance of these models to validate the effectiveness of each component of our paradigm.

4.4.1 NLU Tasks

Tab.3 presents our paradigm experiment results on four NLU benchmarks. For data expansion, we iteratively perform the expansion step 3 times, resulting in a 100k size instruction dataset.

The table results affirm our paradigm’s effectiveness in NLU. Despite a reduction in size by 10k instances compared to the original dataset, our final curated dataset maintains robust performance, showing improvements ranging from nearly 2% to 4% on each benchmark. Furthermore, we observe

Table 1: LLMs Performance Comparison in NLU Tasks

Model	Fine-tuning Data Size	HellaSwag	ARC	TruthfulQA	MMLU
LLaMA-7B		0.778	0.509	0.343	0.357
LLaMA-13B		0.809	0.561	0.395	0.476
LLaMA2-7B	<i>Pretrained</i>	0.771	0.432	0.333	0.444
LLaMA2-13B		0.807	0.488	0.419	0.556
Mistral-7B		0.823	0.602	0.426	0.627
Vicuna-7B	70k conversations	0.775	0.537	0.489	0.456
Vicuna-13B	70k conversations	0.801	0.530	0.518	0.513
WizardLM-7B	70k instructions	0.771	0.516	0.447	0.427
WizardLM-13B	70k instructions	0.777	0.572	0.505	0.523
Platypus2-13B	25k instructions	0.826	0.613	0.449	0.567
Camel-Platypus2-13B	25k instructions	0.836	0.608	0.496	0.565
Stable-Platypus2-13B	25k instructions	0.822	0.627	0.525	0.583
Mistral+Random-7B	15k instructions	0.820	0.607	0.438	0.625
Mistral+LIFT-7B	15k instructions	0.844	0.643	0.490	0.645

Table 2: LLMs Performance Comparison in Code Generation Tasks (pass@1)

Model	Data Size	HumanEval	MBPP
CodeT5+		0.309	-
CodeLLaMA	*	0.360	0.470
StarCoder		0.336	0.436
InstructCodeT5+	20k	0.350	-
WizardCoder	78k	0.573	0.518
StarCoder+Random	10k	0.381	0.431
StarCoder+LIFT	10k	0.550	0.495

* *Pretrained models*

a consistent improvement in model performance on both benchmarks after each step of the paradigm. This implies that the instruction’s quality is steadily increasing at each stage.

It’s crucial to note that the Open Platypus (Lee et al., 2023) dataset for NLU tasks is already carefully curated. The results for this dataset underscore that our paradigm is effective not only for LLM-generated datasets but also in elevating the quality of already high-quality datasets, contributing to improved fine-tuned model performance while reducing the dataset size.

4.4.2 Code Generation Tasks

Tab.4 provides an overview of the paradigm experiments conducted on code generation tasks. For data expansion, we repeatedly perform the expansion step 2 times, resulting in a 60k size instruction dataset.

The table illustrates our paradigm leads to a significant enhancement in the performance of the fine-tuned model across both benchmarks. Notably,

our final curated dataset, although roughly half the size of the original dataset, outperforms the latter by nearly 15% on the HumanEval and 3% on the MBPP. The observed pattern of performance improvement in NLU tasks also extends to code generation tasks, where each step contributes to enhancing data quality. This further underscores that each component of our paradigm plays a vital role in elevating the overall instruction quality.

5 Discussions

5.1 Composition of The Final Curated Dataset

We take a step further to analyze the composition of the final curated dataset, unraveling the origins of diverse and high-quality instruction items. Fig.5 presents the source proportions of the final curated dataset for NLU and code generation tasks, yielding several noteworthy conclusions.

For LLM-generated instruction datasets like Code Alpaca, only a small proportion of the final dataset emanates from the original dataset (Fig.5b). The majority of high-quality data is derived from our paradigm’s first step—the expanded dataset. This emphasizes our paradigm’s significant role in generating and covering a diverse and high-quality dataset, especially for datasets without meticulous curation.

In contrast, for a curated and high-quality instruction dataset like Open Platypus, the portion of the original dataset in the final dataset increases (Fig.5a). The proportions of the final curated dataset in NLU tasks reveal an almost equal distribution among the four sub-datasets, demonstrating

Table 3: Paradigm Ablation Experiment Results in NLU Tasks

Dataset Type	Dataset Size	HellaSwag	ARC Challenge	TruthfulQA	MMLU
Base Platypus Dataset	25k	0.82788	0.61543	0.44481	0.62619
Data Expansion	100k	0.83308	0.62372	0.44718	0.63065
Variety Compress	20k	0.83947	0.63311	0.45615	0.64199
Quality Compress	15k	0.84415	0.64334	0.48985	0.64519

Table 4: Paradigm Ablation Experiment Results (Pass@1) in Code Generation Tasks

Dataset Type	Size	HumanEval	MBPP
Base Dataset	20k	0.4091	0.4662
Data Expansion	60k	0.5342	0.4874
Variety Curation	12k	0.5412	0.4887
Quality Curation	10k	0.5503	0.4949

that even for an initially high-quality dataset, our paradigm also excels in generating and selecting numerous high-quality data points based on the original dataset.

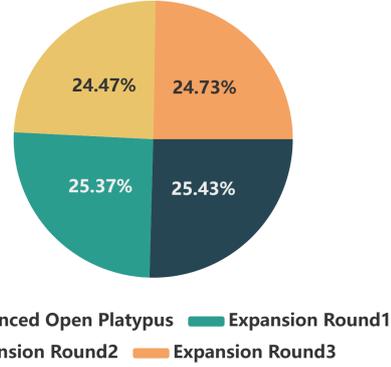
These conclusions affirm that the bulk of the final dataset primarily comprises data from the expanded dataset. While proportions of original dataset data contributing to the final dataset may vary based on the original dataset’s quality, our paradigm consistently showcases its ability to extract high-quality segments from the original dataset and augment them with diverse and high-quality data subspaces.

5.2 Limitations and Future Works

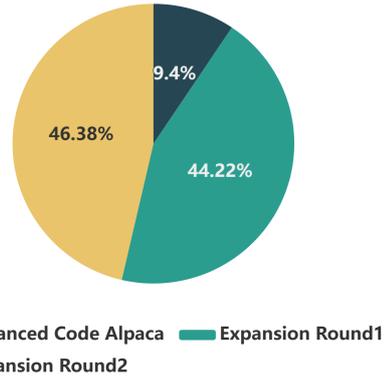
The main limitation of our paradigm lies in the subjectivity of our quality evaluation process, as it heavily relies on GPT quality evaluation. Despite we carefully design some criteria, additional statistical analysis beyond the length factor could also enhance the precision of high-quality data selection. Future work will involve integrating more comprehensive metrics such as information entropy analysis, coherence analysis, and word frequency analysis to provide a more nuanced assessment of instruction quality.

6 Conclusions

This paper presents a novel paradigm, LIFT, that departs from the traditional single-mode quality enhancement approach for instruction datasets, opting for a fresh perspective on data quality through data distribution transfer. By combining the strengths of data expansion and curation while mitigating their limitations, LIFT significantly enhances elevate the



(a) Source of the final curated 15k dataset in NLU tasks



(b) Source of the final curated 10k dataset in code generation tasks

Figure 5: Composition of The Final Curated Dataset

quality of the instruction dataset to new heights. Extensive experimental results demonstrate that our fine-tuned models consistently attain either SOTA or nearly SOTA performance in both NLU and code generation. These experiments underscore the paradigm’s versatile effectiveness, showcasing its capacity to encompass and select diverse and high-quality data. The integration of the curated data into the instruction-tuning process empowers LLMs to achieve superior performance across various tasks and benchmarks.

References

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen

599	Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. 2021. Program synthesis with large language models .	10903–10914, Singapore. Association for Computational Linguistics.	655
600			656
601			
602	Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun. 2023. Instruction mining: When data mining meets large language model finetuning .	Albert Qiaochu Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L'elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b . <i>ArXiv</i> , abs/2310.06825.	657
603			658
604			659
605	Sahil Chaudhary. 2023. Code alpaca: An instruction-following llama model for code generation. https://github.com/sahil280114/codealpaca .		660
606			661
607			662
608	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating large language models trained on code .	Ariel N. Lee, Cole J. Hunter, and Nataniel Ruiz. 2023. Platypus: Quick, cheap, and powerful refinement of llms . <i>ArXiv</i> , abs/2308.07317.	663
609			664
610			665
611			666
612			667
613			
614			668
615			669
616			670
617			671
618			672
619			
620			673
621			674
622			675
623			676
624			677
625			678
626			679
627			680
628	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 .	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umabathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. 2023b. Starcoder: may the source be with you!	681
629			682
630			683
631			684
632			685
633			686
634	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge . <i>ArXiv</i> , abs/1803.05457.		687
635			688
636			689
637			690
638			691
639	Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2023. Mods: Model-oriented data selection for instruction tuning . <i>ArXiv</i> , abs/2311.15653.	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	692
640			693
641			694
642	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding . <i>ArXiv</i> , abs/2009.03300.		695
643			
644			696
645			697
646	J. Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models . <i>ArXiv</i> , abs/2106.09685.	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, MING GONG, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie LIU. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation . In <i>Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks</i> , volume 1. Curran.	698
647			699
648			700
649			701
650	Yufan Huang, Mengnan Qi, Yongqiang Yao, Maoquan Wang, Bin Gu, Colin Clement, and Neel Sundaresan. 2023. Program translation via code distillation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages		702
651			703
652			704
653			705
654			706

713	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct .	David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models . <i>ArXiv</i> , abs/2307.09288.	770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787
718	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.		
719			
720	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll 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 . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.		
721			
722			
723			
724			
725			
726			
727			
728			
729			
730	Mengnan Qi, Yufan Huang, Maoquan Wang, Yongqiang Yao, Zihan Liu, Bin Gu, Colin Clement, and Neel Sundaresan. 2023. SUT: Active defects probing for transcompiler models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 14024–14034, Singapore. Association for Computational Linguistics.	Yue Wang, Hung Le, Akhilesh Deepak Gotmare, Nghi D. Q. Bui, Junnan Li, and Steven C. H. Hoi. 2023. Codet5+: Open code large language models for code understanding and generation . <i>ArXiv</i> , abs/2305.07922.	788 789 790 791 792
731			
732			
733			
734			
735			
736			
737	Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2019. Zero: Memory optimizations toward training trillion parameter models . <i>SC20: International Conference for High Performance Computing, Networking, Storage and Analysis</i> , pages 1–16.	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2021. Finetuned language models are zero-shot learners . <i>ArXiv</i> , abs/2109.01652.	793 794 795 796
738			
739			
740			
741			
742			
743	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P Bhatt, Cristian Cantón Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre D’efosse, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code llama: Open foundation models for code . <i>ArXiv</i> , abs/2308.12950.	Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners .	797 798 799 800
744			
745			
746			
747			
748			
749			
750			
751			
752			
753	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 .	Lai Wei, Zihao Jiang, Weiran Huang, and Lichao Sun. 2023. Instructiongpt-4: A 200-instruction paradigm for fine-tuning minigt-4 . <i>ArXiv</i> , abs/2308.12067.	801 802 803
754			
755			
756			
757			
758	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. 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. 2023. Wizardlm: Empowering large language models to follow complex instructions . <i>ArXiv</i> , abs/2304.12244.	804 805 806 807 808
759			
760			
761			
762			
763			
764			
765	Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull,	Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023. Harnessing the power of llms in practice: A survey on chatgpt and beyond . <i>CoRR</i> , abs/2304.13712.	809 810 811 812 813
766			
767			
768			
769			
		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> .	814 815 816 817 818
		Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, L. Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment . <i>ArXiv</i> , abs/2305.11206.	819 820 821 822 823

923 ### Example 2:
 924 ### Instruction: {EXAMPLE INSTRUCTION 2}
 925 ### Input: {EXAMPLE INPUT 2}
 926 ### Response: {EXAMPLE OUTPUT 2}
 927 ### Score for Example 2: {SCORE 2}
 928 ### Example 3:
 929 ### Instruction: {EXAMPLE INSTRUCTION 3}
 930 ### Response: {EXAMPLE OUTPUT 3}
 931 ### Score for Example 3: {SCORE 3}

932
 933 Please score the upcoming Instruction,
 934 Input and Response based on these examples
 935 across four dimensions, and then add the
 936 four scores together to get the total
 937 score. Try to avoid getting a full score
 938 as much as possible.

939 Please first output a single line
 940 containing the total score number only.
 941 In the subsequent line, please provide
 942 a comprehensive explanation of your
 943 evaluation, avoiding any potential bias.

944 ### Instruction:

945 {INSTRUCTION}

946 ### Input:

947 {INPUT}

948 ### Response:

949 {OUTPUT}

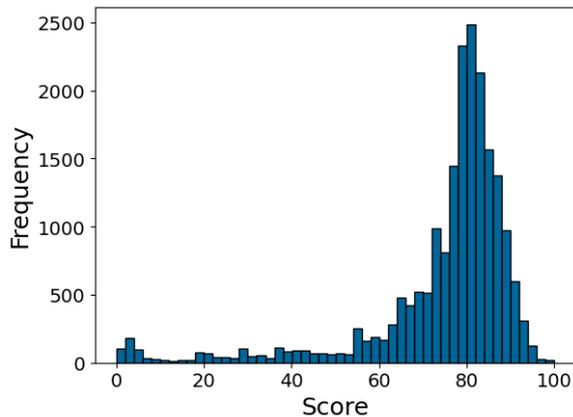
950 C Quality Score Distribution

951 We have gathered the quality scores for the variety
 952 curated dataset following our paradigm, both in
 953 NLU and code generation tasks. The score distri-
 954 butions are depicted in Fig.6. Notably, the quality
 955 scores exhibit an approximately normal distribu-
 956 tion within the score interval of 60 to 100 for both
 957 tasks. This observation validates the effectiveness
 958 of our scoring strategies in discerning low-quality
 959 data. It should be noted that the minor bumps near
 960 0 stem from connection errors or OpenAI API call-
 961 ing ratio constraints, resulting in GPT scores of 0
 962 for certain instructions.

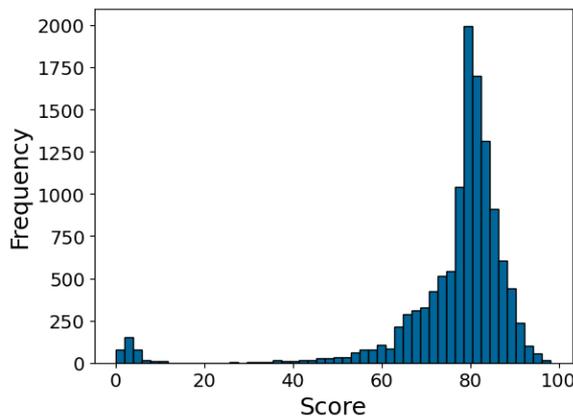
963 D Experiments Implementation Details

964 For both foundation models, we conduct training
 965 on Azure Machine Learning Studio’s cluster ¹, uti-
 966 lizing 4 nodes, each equipped with 8 V100 GPUs
 967 featuring DeepSpeed Zero-3 (Rajbhandari et al.,
 968 2019) offload. Specifically, during the fine-tuning
 969 of Mistral 7B, we employ LoRA (Hu et al., 2021).
 970 This strategy is chosen for its ability to ensure a

¹<https://ml.azure.com/>



(a) Quality Score Distribution in NLU Tasks (20k Data)



(b) Quality Score Distribution in Code Generation Tasks (12k Data)

Figure 6: Quality Score Distribution

971 more steady convergence of loss, resulting in better
 972 performance. The detailed fine-tuning arguments
 973 are outlined in Tab.5.

974 E Benchmarks and Compared LLMs

975 E.1 Benchmarks

976 Large language model benchmarks serve as stan-
 977 dardized tests to evaluate how well models under-
 978 stand, generate, and manipulate human-like lan-
 979 guage (Lu et al., 2021; Chen et al., 2021). Below
 980 is an introduction to these chosen benchmarks:

- 981 • **HellaSwag** (Zellers et al., 2019). HellaSwag
 982 is a challenge dataset containing 70k multiple-
 983 choice questions for evaluating commonsense
 984 Natural Language Inference (NLI). While its
 985 questions may be trivial for humans (>95%
 986 accuracy), they pose a challenge for state-of-
 987 the-art models.
- 988 • **ARC Challenge** (Clark et al., 2018). The
 989 AI2’s Reasoning Challenge (ARC) dataset is

Table 5: Fine-tuning Arguments for StarCoder 15B and Mistral 7B

Arguments	StarCoder	Mistral
model_max_length	1024	2048
batch_size	8	8
num_epoch	3	3
learning_rate	2e-5	2e-5
fp16	True	True
lora_r	-	16
lora_alpha	-	32
lora_dropout	-	0.05

a multiple-choice question-answering dataset containing questions from science exams ranging from grade 3 to grade 9. It is split into two partitions: Easy and Challenge. The Challenge partition consists of 25k questions that require reasoning.

- **TruthfulQA** (Lin et al., 2022). TruthfulQA is a benchmark designed to measure whether a language model is truthful in generating answers to questions. The benchmark comprises 817 questions spanning 38 categories. Questions are crafted so that some humans might answer falsely due to false beliefs or misconceptions.
- **MMLU** (Hendrycks et al., 2020). MMLU (Massive Multitask Language Understanding) is a new benchmark intended to measure knowledge acquired during pretraining. It evaluates models exclusively in zero-shot and few-shot settings, making it more challenging and akin to human evaluation. The benchmark covers 57 subjects across STEM, humanities, social sciences, and more, ranging in difficulty from elementary to advanced professional levels, testing both world knowledge and problem-solving ability.
- **HumanEval** (Chen et al., 2021). HumanEval is utilized to gauge functional correctness in synthesizing programs from docstrings. Comprising 164 original programming problems, it assesses language comprehension, algorithms, and simple mathematics.
- **MBPP** (Austin et al., 2021). The MBPP (Mostly Basic Python Problems) dataset consists of around 1,000 crowd-sourced Python programming problems. These are designed

to be solvable by entry-level programmers, covering programming fundamentals and standard library functionality. In our experiments, to align with others, we select 400 questions.

E.2 Compared LLMs

The selected models for comparison in NLU tasks include:

- **LLaMA** (Touvron et al., 2023a). LLaMA is a collection of foundation language models trained on trillions of tokens from publicly available datasets.
- **LLaMA2** (Touvron et al., 2023b). Llama 2 is an updated version of Llama, trained on a new mix of publicly available data. It increased the size of the pretraining corpus by 40%, doubled the context length of the model, and adopted grouped-query attention in training.
- **Mistral** (Jiang et al., 2023). Mistral is a state-of-the-art 7B foundational model, fast-deployed, easily customizable, and supports English and code with an 8k context length. It’s also one of the foundation models in our paradigm experiments.
- **Vicuna** (Chiang et al., 2023). Vicuna is an open-source chatbot trained by fine-tuning LLaMA on 70K user-shared conversations collected from the ShareGPT website.
- **WizardLM** (Xu et al., 2023). WizardLM is instruction fine-tuned on LLaMA with 70k instruction data generated through the Evol-Instruct strategy.
- **Platypus** (Lee et al., 2023). Platypus is a family of fine-tuned and merged LLMs achieving strong performance. It uses Open Platypus as its instruction dataset and applies the LoRA strategy to train adaptors that can be merged into different foundation models, creating many variant models.

The selected models for comparison in code generation tasks include:

- **CodeT5+ & InstructionCodeT5+** (Wang et al., 2023). CodeT5+ is a new family of open code LLMs with an encoder-decoder architecture trained on various pretraining tasks. InstructionCodeT5+ is further fine-tuned on the Code Alpaca dataset.

- **Code LLaMA** (Rozière et al., 2023). Code Llama is a code-specialized version of Llama2 (Touvron et al., 2023b) trained on code-specific datasets.
- **StarCoder** (Li et al., 2023b). StarCoder is a widely-used large code language model trained on diverse sources, including 80+ programming languages, Git commits, GitHub issues, and Jupyter notebooks. It’s also one of the foundation models in our paradigm experiments.
- **WizardCoder** (Luo et al., 2023). WizardCoder is instruction fine-tuned on StarCoder with 78k instruction data generated through the application of Code Evol-Instruct.

demonstrates that our paradigm not only accelerates fine-tuning but also promotes environmental sustainability while maintaining robust high performance.

F GPU Hours and Carbon Emission

Table 6: Analysis of GPU hours and Carbon Emission with different dataset size. GPU hours in the table are measured in *hour*, CO₂ emission is in *kg CO₂ eq.*

Dataset Size	GPU Hours	CO ₂ Emission
<i>NLU Tasks</i>		
Original 25k	40.24	3.62
Expanded 100k	149.76	13.48
Curated 15k	23.71	2.13
<i>Code Generation Tasks</i>		
Original 20k	50.82	4.58
Expanded 60k	185.6	16.7
Curated 10k	31.6	2.84

By compressing the size of the instruction dataset, we aim to reduce the GPU hours required for instruction tuning, contributing to a subsequent decrease in carbon emissions. Tab.6 illustrates the impact of different dataset sizes on GPU hours and CO₂ emissions. We consider three datasets for each task: the original dataset, the expanded dataset after the first step of our paradigm, and the final curated dataset. GPU hours are calculated under the same settings of training epoch and batch size, while carbon emissions are computed using an online machine learning CO₂ calculator².

The table shows a substantial reduction in GPU hours and lower carbon emissions when fine-tuning with the final curated dataset. Specifically, compared to the original dataset, we observe a 36.8% and 41.1% reduction in GPU hours for code generation and NLU tasks, respectively. This comparison

²<https://mlco2.github.io/impact/#co2eq>