

FANNO: Augmenting High-Quality Instruction Data with Open-Sourced LLMs Only

Anonymous ACL submission

Abstract

Instruction fine-tuning stands as a crucial advancement in leveraging large language models (LLMs) for enhanced task performance. However, the annotation of instruction datasets has traditionally been expensive and laborious, often relying on manual annotations or costly API calls of proprietary LLMs. To address these challenges, we introduce FANNO, a fully autonomous, open-sourced framework that revolutionizes the annotation process without the need for pre-existing annotated data. Utilizing a Mistral-7b-instruct model, FANNO efficiently produces diverse and high-quality datasets through a structured process involving document pre-screening, instruction generation, and response generation. Experiments on Open LLM Leaderboard and AlpacaEval benchmark show that the FANNO can generate high-quality data with diversity and complexity for free, comparable to human-annotated or cleaned datasets like Alpaca-GPT4-Cleaned.

1 Introduction

Large language models (LLMs) have made significant contributions across numerous fields (Zheng et al., 2024a; Wang et al., 2024; Wettig et al., 2024; Fan et al., 2023). Instruction tuning (Ouyang et al., 2022) enhances the model’s general capabilities for novel tasks and improves their adherence to directives. However, the development of human-annotated instruction data is prohibitively expensive, and often results in suboptimal outcomes (Srivastava et al., 2022; Conover et al., 2023). This is primarily due to the annotators’ cognitive limitations, which hinder achieving a balanced dataset in terms of diversity, complexity, and quality (Srivastava et al., 2022; Conover et al., 2023). Previous works explore the automatic LLM-based annotation of instruction data, with advanced proprietary models (Wang et al., 2022a; Xu et al., 2023) or models trained with seed response-query

pairs (Li et al., 2024; Lou et al., 2024). Nevertheless, these approaches often depend on costly APIs (ChatGPT/GPT-4) or require manually crafted seed datasets. Recent studies (Zheng et al., 2024c; Yehudai et al., 2024; Press et al., 2023) aim to construct instruction datasets from scratch; however, the strategies to balance the diversity, complexity, and quality (Liu et al., 2023a) of annotated instruction data are less explored.

Addressing these challenges, we introduce FANNO (Free ANNOTator), a freely accessible framework specifically designed for automatic high-quality instruction annotation. This framework methodically breaks down the annotation process into three distinct phases: document pre-screen, instruction generation, and response generation. It utilizes curated tagging, UCB(Upper Confidence Bound) bootstrapping iterations, and filtering techniques to enhance the diversity and complexity of the generated instructions. Empirical evidence on Open LLM Leaderboard and AlpacaEval benchmark confirm the framework’s efficacy on two 7B LLMs. The resulting dataset is virtually indistinguishable from those refined datasets like **Alpaca-GPT4-Cleaned**, marking a significant stride in instruction data development¹.

2 Related Work

Instruction Data Generation Two main approaches have been explored for instruction data creation: (1) **Human Annotation**, which leverages human expertise to design prompts and collect multi-task datasets spanning various categories (Srivastava et al., 2022; Conover et al., 2023). While producing high-quality data, manual annotation is effort-intensive and costly, especially for devising complex textual instructions. (2) **LLM Synthetic Data Generation** Recent research increasingly favors harnessing the creative

¹Our code, data, and model will be made public.

079	capabilities of LLMs, such as GPT-4 (OpenAI,	similar approaches utilize fine-tuned large models	130
080	2023), over human input for creating instruction-	to score the data for quality assessment. Moreover,	131
081	following datasets (Geng et al., 2023; Chiang	efforts like ORCA-MATH (Mitra et al., 2024) and	132
082	et al., 2023). ALPACA (Taori et al., 2023) and	REFLECTION-TUNING (Li et al., 2023a) employ	133
083	ALPACAGPT (Peng et al., 2023) have also uti-	collaborative approaches with multiple LMs and	134
084	lized more powerful LLMs to enhance data qual-	self-reflection to enhance data quality.	135
085	ity. Another line of research involves generat-		
086	ing task instructions from “seeds” and filtering (Wu	3 FANNO Framework	136
087	et al., 2023). For example, WIZARDLM (Xu	The FANNO framework aims to annotate diverse,	137
088	et al., 2023) employed an instruction evolution	complex, and faithful instruction data with only	138
089	paradigm to increase seed instruction complex-	free open-sourced LLMs. As depicted in Figure 1,	139
090	ity, while SELF-INSTRUCT (Wang et al., 2022a) used	FANNO consists of three pivotal steps: document	140
091	human-annotated instructions as demonstra-	pre-screen, instruction generation, and response	141
092	tions to guide LLMs in the instruction evolu-	generation.	142
093	tion process. Humpback (Li et al., 2024) gener-		
094	ates instructions using vast amounts of unlabeled	3.1 Document Pre-Screen	143
095	web text. These datasets are costly, either in	The FANNO framework annotates instruction data	144
096	terms of labor or proprietary model expenses.	from web corpus, textbooks, etc. The document	145
097	In contrast, FANNO maintains high instructional	pre-screening stage initially includes segmenta-	146
098	quality autonomously, utilizing open-source	tion, deduplication, and length-based filtering.	147
099	models efficiently with just a 7B model size.	Further filtering employs a teacher LLM and a	148
100		fast community detection algorithm to enhance	149
101	Instruction Tuning Instruction tuning involves	correctness and diversity.	150
102	training LLMs on extensive upstream task	The LLM-based filter addresses ambiguous	151
103	datasets with instructions, followed by enab-	content, privacy concerns, and advertisements	152
104	ling the generalized ability to new, unseen	(see Appendix D.1). To reduce data volume	153
105	downstream tasks via new instructions (Ouyang	while maintaining diversity, we cluster	154
106	et al., 2022; Chung et al., 2022). This	instruction embeddings using a fast	155
107	technique is widely acknowledged as	community detection algorithm, similar	156
108	essential for activating LLMs to adhere to	to SentenceTransformer (Reimers and	157
109	human conversational norms (Mishra et al.,	Gurevych, 2019), based on a predefined	158
110	2022). Instruction tuning has empowered	similarity threshold. This approach	159
111	various domain-specific or task-specific	prioritizes larger, non-overlapping	160
112	LLMs (Jiang et al., 2023b; Xu et al.,	communities (details in Appendix B.3,	161
113	2023), and curating diverse, high-quality	Algorithm 1).	162
114	upstream instruction dataset has become	The pre-screen phase balances	163
115	a pivotal step for successful instruction	processing speed and precision,	164
116	tuning (Wang et al., 2023; Lou et al.,	prioritizing efficiency. In our	
117	2023). Moreover, instruction tuning	experiments, the pre-screen stage	
118	also bolsters cross-task general	filters and keeps 6% of	
119	capabilities (Sanh et al., 2022; Wang	the original raw data.	
120	et al., 2022b), encompassing a	3.2 Instruction Generation	165
121	more comprehensive array of	At this stage, FANNO adopts a	166
122	general tasks, notably incorporating	bootstrapping approach to	167
123	input from users of language	generate instructions from	168
124	models (Ouyang et al., 2022; Peng	pre-screened documents,	169
125	et al., 2023). Related works in	streamlining the process	170
126	the field of enhancing data	into two distinct phases:	
127	quality have focused on	seed instruction genera-	
128	several key aspects such	tion and instruction	
129	as instruction difficulty,	augmentation.	
	diversity, and correct-	Step 1: Seed Instruction	171
	ness. HUMPBACK (Li	Generation This step	172
	et al., 2024) and KUN	produces a set of diverse	173
	(Zheng et al., 2024c)	instructions as the initial	174
	utilize language model’s	seeds. Diversity is	175
	capability in combina-	promoted from two	176
	tion with tailored	perspectives: Task	
	prompts for data	Types and Diffi-	
	filtering. In Addition,	iculty Levels , for	
	initiatives like GENIE	which we have	
	(Yehudai et al.,	manually created	
	2024) and MODS	corresponding	
	(Du et al., 2023)	tags (see	
	utilize specialized	Appendix D.4). For	
	open-source LLMs	each document,	
	for data filtering	we traverse	
	tasks. DEITA (Liu		
	et al., 2023a),		
	PLANGPT (Zhu		
	et al., 2024) and		

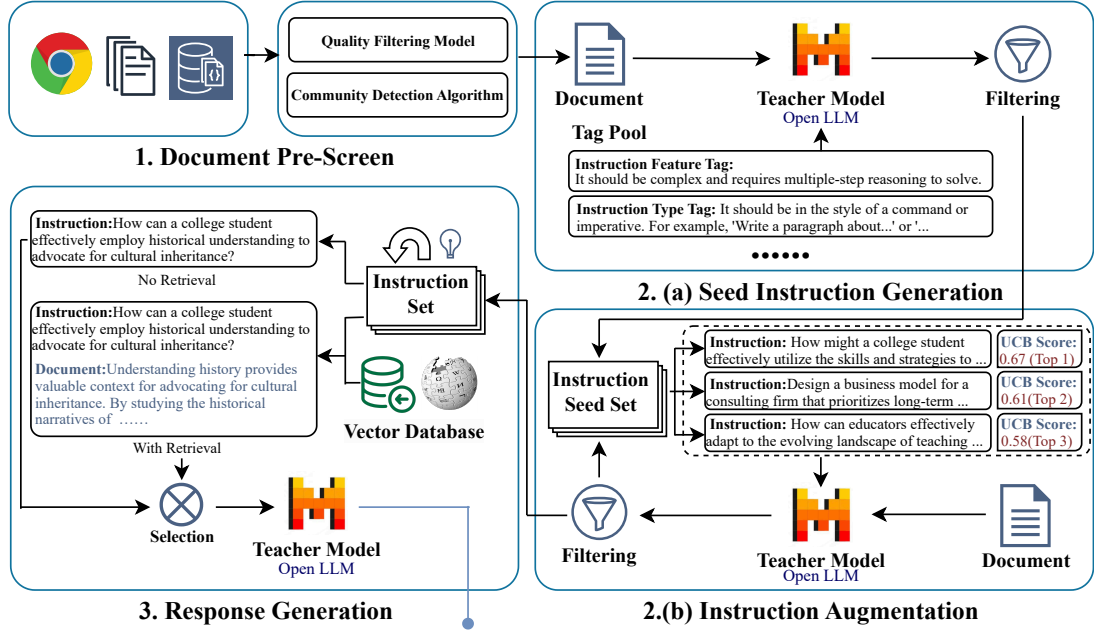


Figure 1: Overview of FANNO framework. **(1) Document Pre-Screen:** We process the unlabeled text data with filters and community detection algorithm. **(2a) Seed Instruction Generation:** FANNO generates seed instructions from pre-screened documents with diverse task types and difficulty levels through a tag pool. **(2b) Instruction Augmentation:** New instructions are augmented conditioned on the documents and few-shot examples selected from the seed instructions with the UCB algorithm. **(3) Response Generation:** The responses to instructions are generated directly by the teacher LLM or based on the concatenation of the corresponding document and retrieved document.

all combinations of task types and instruction difficulty levels to generate seed instructions. An LLM-based filter (see Table 8 in the appendix) is then employed to ensure the quality of the seed instruction data. We sample 200 documents for instruction generation and obtain around 1k instructions as the seed pool S .

Step 2: Instruction Augmentation The diversity of the instructions in S is inherently limited. To promote the diversity of newly generated instructions, we designed a prompt template called **Think Different** (see Appendix 13), which diverges from the traditional example-followed template used in self-instruct (see Appendix D.6). This template encourages the teacher model to generate high-quality instructions that emulate the quality of the examples but differ in format (task types, questioning styles, etc.). Additionally, a document is inputted into this template to ensure the generated instructions are consistent with or extended from this document.

The quality of the examples is, therefore, paramount. Instead of randomly selecting examples, we prioritize extracting higher-quality ones, assuming that instruction length correlates with

quality. To avoid suboptimal convergence, the UCB (Upper Confidence Bound) (Robbins and Monro, 1951) score is used to enhance the exploration of new instructions. Each seed data is scored as $UCB(s) = \bar{x}_s + C\sqrt{\frac{2\ln N}{n_s}}$. Here, \bar{x}_s is the seed’s average quality, N is the total iterations, and C is a constant. The score promotes high-quality and less frequently selected seeds, with C balancing these objectives. In each iteration, we select k seeds with the highest UCB scores, effectively trade-off between exploration and exploitation. We compare UCB and random sampling in an ablation study. The detailed algorithm can be found in Appendix B.2.

3.3 Response Generation

At this stage, the response to each instruction is generated by prompting the teacher LLM either with empty context or a retrieved document. We propose to apply retrieval augmented generation (RAG) and incorporate the corresponding document to provide additional information for response generation. These documents are concatenated to serve as the relevant context. For all generations, the teacher LLM is prompted to generate responses

under the above two different conditions. Then, we use the LLM itself to select the response with better quality. The prompt templates in Appendix D.3 and Table 15 are used for response generation and selection, respectively.

3.4 Discussion

To produce diverse and complex high-quality instruction data, FANNO utilizes tags for difficulty balancing, iteratively selects high-quality data via UCB bootstrap, and ensures diversity through iterative instruction filtering. We generated strongly generalized data independent of the original text through carefully crafted prompts. To ensure fidelity between instructions and responses, information is supplemented using RAG and a Teacher model. Detailed discussions are in Section 5.

4 Experiment

4.1 Experiment setup

Unlabeled Text Data We use the FALCON REFINED WEB corpus² (Penedo et al., 2023), a large web-based corpus dataset including 600 billion tokens, as our unlabeled data. We directly selected the first 500k documents for input to the Document Pre-Screen stage.

Models and Training Details We choose the Mistral-7b-instruct-v0.2 (Jiang et al., 2023a) for data annotation in all experiments. We perform supervised instruction tuning using LoRA (Hu et al., 2021) with the pretrained LLaMA-2-7b-base model (Touvron et al., 2023) and Mistral-7b-base model (Jiang et al., 2023a). The model after instruction tuning with the data annotated with FANNO framework is referred to as FANNO. Detailed configuration can be viewed in Appendix C.

Baselines We compare FANNO with models fine-tuned with other instruction datasets. The baseline details are in Appendix A. The datasets include Alpaca-52k (Taori et al., 2023), Alpaca-GPT4 (Peng et al., 2023), Alpaca-Cleaned, LIMA (Zhou et al., 2023), WizardLM-70k (Xu et al., 2023), and Muffin (Lou et al., 2024). Alpaca-52k and Alpaca-GPT4, each with 52,002 samples, use Text-Davinci-003 and GPT-4 for annotations. Alpaca-Cleaned refines Alpaca-GPT4 to 51,760 samples filtered instructions with hallucination errors or invalid outputs. LIMA offers

²<https://huggingface.co/datasets/tiiuae/falcon-refinedweb>

1,000 manually selected diverse prompts and responses. WizardLM-70k and Muffin, both using ChatGPT or GPT-4 annotations, focus on 70,000 and 68,000 high-quality samples, respectively. The self-augmented dataset of Humpback is also comparatively ensured to be fair.

4.2 Evaluation

Open LLM Leaderboard The Huggingface Open LLM Leaderboard³ (Beeching et al., 2023) stands as a unified framework designed to evaluate generative language models across a wide array of diverse evaluation tasks. It encompasses key benchmarks such as ARC (Clark et al., 2018), Hel-laSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), and TruthfulQA (Lin et al., 2022). We utilize the lm-evaluation-harness toolkit⁴ (Gao et al., 2023) for evaluating different models to maintain consistency with the official setup.

AlpacaEval 2.0 AlpacaEval Benchmark (Li et al., 2023b) is an automated evaluation framework based on an annotation model (GPT-4). By comparing responses generated by two different models for the same set of 805 prompts, AlpacaEval computes the pairwise win rate, automating the evaluation process.

Human Evaluation We employed manual annotation by multiple experts to identify the complexity of instructions, specifically categorized into three tiers: (0 Unanswerable, 1 Easy, 2 Expert). The detailed information for each tier is provided in Appendix F.2. The evaluation results are presented in Table 5 and discussed in Section 5.2.

MT-Bench The MT-Bench (Multi-turn Benchmark) (Zheng et al., 2024b) is aimed at assessing the conversational and instruction-following abilities of LLMs. It comprises 80 multi-turn questions, and GPT-4 is utilized as an automated evaluator, scoring chatbot responses on a scale of 1 to 10, with methods in place to minimize bias and enhance the reliability of the assessments.

4.3 Results

The comparative experiments of FANNO with other models demonstrate the superiority of our work. (1) For diverse base models like LLaMA and

³https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard

⁴<https://github.com/EleutherAI/lm-evaluation-harness>

Model	Data Size	ARC	HellaSwag	MMLU	TruthfulQA	Average
Open-sourced Models based on LLaMA-2						
LLaMA-2-Base	–	54.10	78.71	45.80	38.96	50.76
LLaMA-2-Chat	–	54.10	78.65	45.69	44.59	55.76
LLaMA-2 + Alpaca-52k	52k	54.78	78.17	46.65	41.43	55.26
LLaMA-2 + Alpaca-GPT4	52k	56.66	78.78	46.96	51.02	58.35
LLaMA-2 + Alpaca-Cleaned	51.8k	56.40	80.16	47.02	50.53	58.53
LLaMA-2 + LIMA	1k	54.61	79.21	45.79	41.32	55.23
LLaMA-2 + WizardLM-70k	65k	54.01	78.66	45.61	38.99	54.32
LLaMA-2 + Muffin	68k	54.10	76.97	47.12	43.51	55.42
LLaMA-2 + FANNO	16k	55.63	79.45	46.84	51.01	58.23
Open-sourced Models based on Mistral-7B						
Mistral-7B-Instruct-v0.2	–	59.39	84.33	59.28	66.79	67.45
Mistral-7B-Base-v0.1	–	60.84	83.31	62.42	42.59	62.29
Mistral-7B-Base + Alpaca-GPT4	52k	63.65	82.18	59.29	43.98	62.29
Mistral-7B-Base + Alpaca-Cleaned	51.8K	64.51	83.68	59.76	52.00	64.99
Mistral-7B-Base + FANNO	16k	64.16	85.08	60.79	52.16	65.55

Table 1: Open LLM Leaderboard results evaluated with the lm-evaluation-harness toolkit. Data size represents the number of samples in the instruction data.

Table 2: Comparison of Different Models on MT Bench

Model	MT Bench
LLaMA-2-7B	3.97
LLaMA-2-7B (alpaca-gpt4)	4.31
LLaMA-2-7B (self-instruct-Teacher: Mistral)	4.96
LLaMA-2-7B-chat Official implementation	6.27
LLaMA-2-7B (FANNO-Teacher: LLaMA-2-Chat)	4.65
LLaMA-2-7B (FANNO-Teacher: Mistral)	5.11

Mistral, our framework consistently achieves top rankings in the LLM-open-leaderboard, even rivaling the models fine-tuned with Alpaca-GPT4-Cleaned, which underwent augmentation with proprietary models and manual selection. (2) Compared to other similar automatic instruction annotation frameworks like Humpback, Muffin, WizardLM, we adhered to the principle of fairness as much as possible by fully utilizing the officially published datasets, and experiments proved that FANNO achieved excellent results with a smaller dataset.

Using MT-Bench in Table 2, we observe that our model outperforms those fine-tuned using Alpaca-GPT4-clean, highlighting the effectiveness of FANNO. Moreover, compared with *self-instruct*, we relieve the need for manually labeled data and achieve better performance, while naive *self-instruct* with Mistral does not yield optimal results. However, it is understandably inferior to the LLaMA2-7B-Chat model, which benefits from extensive fine-tuning and RLHF alignment.

Models refined through FANNO exhibit notable enhancements in the TruthfulQA metric and show

measurable improvements across three other metrics. This advancement is attributed to the integration of supplementary information via the RAG component and self-reflective teacher model, thereby improving the model’s proficiency in delivering more faithful outputs and bolstering TruthfulQA scores. Slight improvements in ARC, HellaSwag, MMLU metrics are credited to the elevated challenge and diversity of the instructions, as depicted in Table 1. We also uploaded our model to Huggingface Open LLM Leaderboard and compared our results with models like Vicuna, and Humpback, which are shown in Table 3. As shown in Figure 2, our model marginally outperforms the Alpaca-GPT4-Cleaned’s fine-tuned variant on the AlpacaEval benchmark, attesting to the superiority of our FANNO framework.

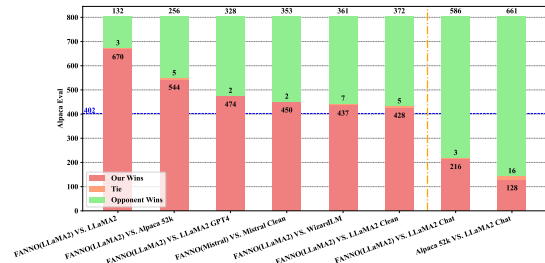


Figure 2: AlpacaEval Result

4.4 Ablation Study

To enhance our understanding of the functionality of each module within FANNO, we undertook ablation studies on its four components, as delineated in Table 4. Our findings reveal:

Model	ARC	HellaSwag	MMLU	TruthfulQA	Average
LLaMA-2-Base	53.07	78.59	46.87	38.76	54.32
LLaMA-2-Chat	52.90	78.55	48.32	45.57	56.34
Vicuna-7b-v1.3	50.43	76.92	48.14	47.01	55.62
Humback M_0	56.31	81.20	47.45	47.59	58.13
Humback M_1	52.99	78.57	45.48	41.45	54.65
WizardLM-7b	54.95	77.85	45.79	48.29	56.72
FANNO	55.46	79.29	46.58	52.05	58.35

Table 3: Benchmark results evaluated by the official Huggingface Open LLM Leaderboard platform.

ID	Configuration	Open LLM Leaderboard				Average
		ARC	HellaSwag	MMLU	TruthfulQA	
①	0 Base	54.44	78.66	44.69	46.02	55.95
	1 Pre-Screen	55.46	78.51	46.00	45.85	56.44
②	2 FANNO (w/o Iter and UCB)	54.69	79.18	45.92	50.19	57.50
	3 FANNO (w/o UCB)	55.63	79.43	44.84	51.16	57.77
③	4 FANNO (w/ OD)	55.46	78.31	44.99	45.68	56.11
	5 FANNO (w/ RAG)	55.29	78.41	45.80	45.37	56.22
	6 FANNO (w/ RAG+)	55.03	78.46	47.02	46.26	56.69
	7 FANNO	55.63	79.45	46.84	51.01	58.23

Table 4: Ablation results from the lm-evaluation-harness (Gao et al., 2023). (0). Basic framework: simply generate instructions by documents and generate responses by instructions without any optimization. (1). Add Pre-Screen module into the basic framework. (2). FANNO without Iteration and UCB-selection. (3). FANNO without UCB-selection. (4). FANNO with the original document. (5). FANNO with RAG module. (6). FANNO with RAG module and supplementary materials. (7). The complete version of FANNO.

- Orthogonality of Components and Separate Optimization** We replaced each component with a random strategy, and the experiments show that each module positively affects the model’s performance, and using more advanced strategies yields better results. ① indicated that the Pre-Screen strategy helps to enhance the quality and thematic diversity of the raw documents. Configurations (3) and (7) demonstrated that using the UCB strategy in instruction augmentation balances the complexity and diversity of the generations, achieving higher diversity compared to random sampling. The notable growth in the MMLU result, as indicated by the ② & 7 combinations, revealed that iterative enhancements in conjunction with the UCB strategy were paramount. The UCB’s proactive selection of high-quality data for augmentation facilitates a gradual evolution towards more effective methodologies as the iteration progresses.
- Generalizing Boosting Diversity and Complexity** As introduced in Section 3, FANNO uses randomness, deduplication, and carefully designed prompts to increase the diversity of themes and tasks in instructions as much as possible. In this way, FANNO tends to break away from reliance on the corresponding unlabeled documents, enhancing the generalizability of instructions. Ab-

lating experiments ① & 7 proved that texts with higher generalizability exhibit more diversity and complexity, which is more beneficial for activating the capabilities of the base models.

- Knowledge Supplementation Promotes Instruction Quality** The results of ③ demonstrated that it is necessary to incorporate RAG or supplement knowledge with the help of a teacher model is necessary. We discovered that considering only the direct generation of the teacher model yielded the best results compared to the document-based response, particularly on TruthfulQA. This indicates that instructions generated with the FANNO framework are more general and less reliable on the corresponding document. We used RAG in an experiment to pinpoint the most relevant content for the instruction, and experiments 5 & 6 support our assertion that more is better. Other discussions about the truthfulness are covered in Section 5. RAG+ used larger datasets than RAG, but both were from Wikipedia.⁵

5 Analyses

In this section, we will discuss how diversity, correctness, and complexity are promoted in each

⁵The RAG used 2.73GB of data from Wikipedia’s introduction section, and the RAG+ used 20.28GB of data from Wikipedia.

415 stage.

416 5.1 Analyses of the Augmented Instruction 417 Data

418 We analyze and illustrate the generated instructions
419 of our dataset from 4 aspects:

420 **Length** To study the distribution of the length
421 of instructions, we tokenize each instruction com-
422 bined with input and count the words within it as
423 its length. Figure 5 and Figure 7 in Appendix E.1
424 illustrate the distribution of instruction length for
425 FANNO and Alpaca-Cleaned, respectively. The
426 results show that FANNO instructions are more bal-
427 anced than Alpaca-Cleaned and the mean value of
428 lengths is higher than that of Alpaca, which indi-
429 cates a better performance.

430 **Diversity** Inspired by SELF-INSTRUCT (Wang
431 et al., 2022a), the verb-noun pairs in instructions
432 to represent the types and tasks of instructions are
433 identified and extracted, which exhibits diversity.
434 As Figure 6 and Figure 8 in Appendix E.1 depicted,
435 FANNO instructions possess more challenging verb-
436 noun pairs than Alpaca-Cleaned, which indicates
437 more challenging tasks. The extraction is com-
438 pleted by Berkeley Neural Parser (Kitaev and Klein,
439 2018; Kitaev et al., 2019).

440 **Quality and Complexity** To evaluate the qual-
441 ity and complexity of instruction-response pairs,
442 we utilize Deita-quality-scorer model and Deita-
443 complexity-scorer model (Liu et al., 2023a) as an
444 evaluator to score our instructions. Figure 9 in
445 Appendix E.1 shows the quality and complexity
446 comparison between FANNO and Alpaca-Cleaned,
447 of which the result shows that FANNO instructions
448 possess a more balanced complexity distribution
449 and higher average quality. The corresponding
450 prompts can be found in Table 9 and Table 10 in
451 the appendix.

452 5.2 Randomness Tag Boosting the Complexity

Score	Tagged	Untagged
0 (bad)	18	24
1 (everyday)	39	78
2 (expert)	143	98

Table 5: Randomness Tag Evaluation Results

453 Randomness tags serve as additional require-
454 ments for the teacher model when generating in-
455 structions, enhancing their complexity. To demon-
456 strate the effectiveness of this strategy, generated in-

457 structions are manually annotated to evaluate their
458 complexity, as discussed in Section 4.2. We ran-
459 domly sampled 200 instances from two datasets
460 for manual testing: one utilizing the Random tag
461 strategy (Tagged) and the other generated directly
462 (Untagged). Results are summarized in Table 5,
463 with example instructions detailed in Appendix F.2.
464 From the results, it is evident that instructions with
465 random tags exhibit a significant increase in com-
466 plexity, manifested by a greater number of expert-
467 level instructions and fewer daily instructions. It
468 is worth noting that instructions with random tags ex-
469 hibit a tendency towards greater length, including
470 few overly complex tasks that are difficult to an-
471 swer (classified as expert-level instructions), which
472 indicates a potential need for further refinement.

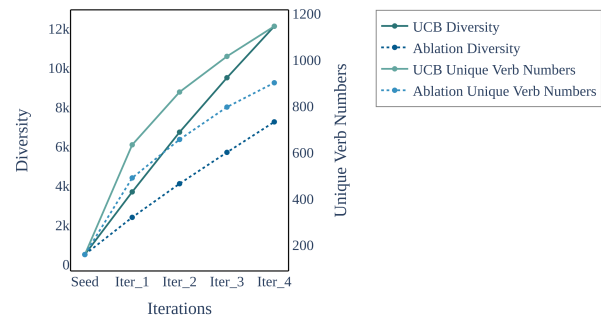


Figure 3: The verbs-noun statistics data grows with iteration

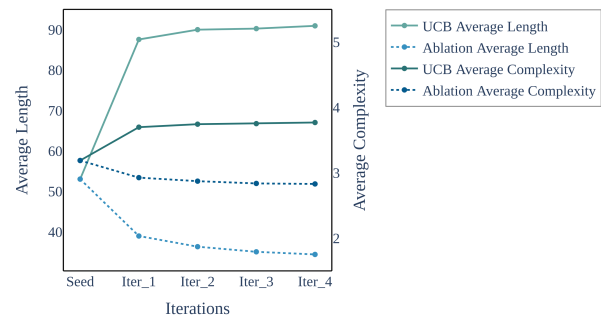


Figure 4: The instruction length (complexity) grows with iteration

473 5.3 UCB Bootstrap Iteration Improve 474 Instruction Complexity while Maintaining 475 the Diversity

476 UCB Bootstrap is employed to actively stabilize
477 the process of instruction improvement. For illus-
478 tration, we monitored the diversity and complexity
479 of instructions during iterations and compared it
480 with a random selection strategy. Note that to sim-
481 plify the process, we use the length of instructions
482 in words as the measure of quality. As depicted in
483 Figure 3 and Figure 4, we observed an increase in

		Open LLM Leaderboard				Average
		ARC	HellaSwag	MMLU	TruthfulQA	
direct-based	Fearless Responses	54.86	79.33	45.56	47.40	56.79
	Cautious Response	54.44	79.17	45.88	47.29	56.69
doc-based	Faithful Responses	55.03	78.63	45.66	42.82	55.54
	Adaptive Responses	55.63	78.71	45.86	42.76	55.74

Table 6: Comparison results of four types response on the Open LLM

both the average complexity and diversity scores as the iteration progressed, consistent with our expectations. We analyze that UCB prioritizes exploring longer instructions for few-shot instruction generation, resulting in more challenging instructions. Additionally, UCB exhibits a preference for selecting newly generated instructions, as novel few-shot combinations tend to ignite the model’s creativity.

5.4 Truthfulness is Less Important for Capability Activation

Previous work (Zhou et al., 2023; Liu et al., 2023b) has explored the diversity, complexity, and fidelity of instructions that enhance large models’ capabilities. We further investigate the truthfulness of responses to instructions, noting that responses often seem accurate but contain illusions and fabricated information, potentially affecting instruction fine-tuning.

To address this, we selected about 1,000 expert-level instructions from the FANNO dataset, then prompted LLM to generate four different responses for each instruction with the following settings:

- Fearless Response: Models provide answers regardless of correctness.
- Cautious Response: Models may acknowledge a lack of knowledge.
- Faithful Response: Models generate answers solely based on provided documents.
- Adaptive Response: Models use relevant information from provided documents to generate answers.

From the result in Table 6, an intriguing finding is that direct responses from the model, which contains a substantial presence of illusions, outperformed document-based ones, particularly in the TruthfulQA task. This might suggest that providing human-like and consistent responses, even with false data, can also improve the model’s capabilities during SFT. We also need to point out FANNO also introduces external sources of information such as the knowledge of the teacher model itself, which likewise results in some illusory responses.

6 Limitations

While FANNO has demonstrated outstanding performance, several limitations must be acknowledged. The responses are not entirely dependent on the document, leading to the introduction of certain hallucinations in the fine-tuning data, as discussed in Section 5.4. This suggests that the model’s reliance on the provided context needs to be strengthened to improve factual consistency. The simplistic approach of equating instruction length with its value is rather crude. The true value of an instruction is influenced by various factors such as difficulty, quality, and novelty. Future work will aim to develop a more nuanced understanding and evaluation of instruction value. The quality of generated instructions is contingent upon the capabilities of both the generator and the evaluator. This process is sensitive to the teacher model and the prompts used, indicating a need for designing prompts that are specifically tailored to the model. Addressing these limitations will be a focus of our future work.

7 Conclusion

The development of instruction fine-tuning datasets has been hindered by the high cost and labor-intensive nature. In this paper, we introduced FANNO, an autonomous and low-cost framework that addresses these challenges by streamlining the annotation process with open-sourced LLMs. FANNO efficiently generates datasets of high quality, diversity, and complexity through a structured process involving pre-screening, instruction generation, and response generation. This unified process eliminates the need for pre-existing annotated data or costly API calls, marking a significant advancement in instruction data development. Empirical experiments also validate the efficacy of FANNO, underscoring the framework’s potential to democratize access to high-quality instruction datasets. FANNO enables access to top-quality datasets with reduced cost and effort, driving progress in LLM applications.

567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621

References

Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. Open llm leaderboard. https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard.

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.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling Instruction-Finetuned Language Models. *arXiv preprint arXiv:2210.11416*.

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. *Preprint*, arXiv:1803.05457.

Mike Conover, Matt Hayes, Matt Mathur, Xiangrui Meng, Jianwei Xie, Jun Wan, Ali Ghodsi, Patrick Wendell, and Patrick Zaharia. 2023. Hello dolly: Democratizing the magic of chatgpt with open models.

Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2023. Mods: Model-oriented data selection for instruction tuning. *Preprint*, arXiv:2311.15653.

Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. Large language models for software engineering: Survey and open problems. *Preprint*, arXiv:2310.03533.

Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. A framework for few-shot language model evaluation.

Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn Song. 2023. Koala: A dialogue model for academic research. Blog post.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *Preprint*, arXiv:2106.09685.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2023a. Mistral 7b. *Preprint*, arXiv:2310.06825.

Dongfu Jiang, Yishan Li, Ge Zhang, Wenhao Huang, Bill Yuchen Lin, and Wenhui Chen. 2023b. Tiger-score: Towards building explainable metric for all text generation tasks. *Preprint*, arXiv:2310.00752.

Nikita Kitaev, Steven Cao, and Dan Klein. 2019. Multilingual constituency parsing with self-attention and pre-training. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3499–3505, Florence, Italy. Association for Computational Linguistics.

Nikita Kitaev and Dan Klein. 2018. Constituency parsing with a self-attentive encoder. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2676–2686, Melbourne, Australia. Association for Computational Linguistics.

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Heng Huang, Jiuxiang Gu, and Tianyi Zhou. 2023a. Reflection-tuning: Data recycling improves llm instruction-tuning. *ArXiv*, abs/2310.11716.

Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason Weston, and Mike Lewis. 2024. Self-alignment with instruction back-translation. *Preprint*, arXiv:2308.06259.

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.

Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Truthfulqa: Measuring how models mimic human falsehoods. *Preprint*, arXiv:2109.07958.

Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023a. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *Preprint*, arXiv:2312.15685.

Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023b. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *Preprint*, arXiv:2312.15685.

Renze Lou, Kai Zhang, Jian Xie, Yuxuan Sun, Janice Ahn, Hanzi Xu, Yu Su, and Wenpeng Yin. 2024. MUFFIN: Curating multi-faceted instructions for improving instruction following. In *The Twelfth International Conference on Learning Representations*.

677	Renze Lou, Kai Zhang, and Wenpeng Yin. 2023. Is prompt all you need? no. a comprehensive and broader view of instruction learning. <i>arXiv preprint arXiv:2303.10475</i> .	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models. <i>arXiv preprint arXiv:2206.04615</i> .	729
678			730
679			731
680			732
681	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-Task Generalization via Natural Language Crowdsourcing Instructions. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3470–3487.	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 .	733
682			734
683			735
684			736
685			737
686			738
687	Arindam Mitra, Hamed Khanpour, Corby Rosset, and Ahmed Awadallah. 2024. Orca-math: Unlocking the potential of slms in grade school math . <i>Preprint</i> , arXiv:2402.14830.		739
688			740
689			741
690			742
691	OpenAI. 2023. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	743
692			744
693			745
694	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>arXiv preprint arXiv:2203.02155</i> .	Siyuan Wang, Zhongyu Wei, Yejin Choi, and Xiang Ren. 2024. Can llms reason with rules? logic scaffolding for stress-testing and improving llms . <i>Preprint</i> , arXiv:2402.11442.	746
695			747
696			748
697			749
698			750
699	Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: Outperforming curated corpora with web data, and web data only . <i>Preprint</i> , arXiv:2306.01116.	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. 2023. How far can camels go? exploring the state of instruction tuning on open resources. <i>arXiv preprint arXiv:2306.04751</i> .	751
700			752
701			753
702			754
703			755
704			756
705	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. <i>arXiv preprint arXiv:2304.03277</i> .	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022a. Self-instruct: Aligning language model with self generated instructions. <i>arXiv preprint arXiv:2212.10560</i> .	757
706			758
707			759
708	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A. Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models . <i>Preprint</i> , arXiv:2210.03350.	Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022b. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	760
709			761
710			762
711			763
712	Sebastian Raschka. 2023. Finetuning llms with lora and qlora: Insights from hundreds of experiments . <i>Lightning AI</i> .		764
713			765
714			766
715	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i> . Association for Computational Linguistics.		767
716			768
717			769
718			770
719			771
720	Herbert Robbins and Sutton Monro. 1951. A stochastic approximation method. <i>The annals of mathematical statistics</i> , pages 400–407.		772
721			773
722			774
723	Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2022. Multitask Prompted Training Enables Zero-Shot Task Generalization. In <i>International Conference on Learning Representations</i> .	Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. 2024. Qurating: Selecting high-quality data for training language models . <i>Preprint</i> , arXiv:2402.09739.	775
724			776
725			777
726			778
727			779
728		Minghao Wu, Abdul Waheed, Chiyu Zhang, Muhammad Abdul-Mageed, and Alham Fikri Aji. 2023.	780
			781
			782
			783
			784
			785

786 Lamini-1m: A diverse herd of distilled mod-
787 els from large-scale instructions. *arXiv preprint*
788 *arXiv:2304.14402*.

789 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,
790 Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin
791 Jiang. 2023. Wizardlm: Empowering large lan-
792 guage models to follow complex instructions. *arXiv*
793 *preprint arXiv:2304.12244*.

794 Asaf Yehudai, Boaz Carmeli, Yosi Mass, Ofir Arviv,
795 Nathaniel Mills, Assaf Toledo, Eyal Shnarch, and
796 Leshem Choshen. 2024. [Genie: Achieving hu-
797 man parity in content-grounded datasets generation.](#)
798 *Preprint*, arXiv:2401.14367.

799 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
800 Farhadi, and Yejin Choi. 2019. [Hellaswag: Can
801 a machine really finish your sentence?](#) *Preprint*,
802 arXiv:1905.07830.

803 Huaixiu Steven Zheng, Swaroop Mishra, Xinyun Chen,
804 Heng-Tze Cheng, Ed H. Chi, Quoc V Le, and Denny
805 Zhou. 2024a. [Take a step back: Evoking reasoning
806 via abstraction in large language models.](#) *Preprint*,
807 arXiv:2310.06117.

808 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
809 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
810 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024b.
811 Judging llm-as-a-judge with mt-bench and chatbot
812 arena. *Advances in Neural Information Processing*
813 *Systems*, 36.

814 Tianyu Zheng, Shuyue Guo, Xingwei Qu, Jiawei Guo,
815 Weixu Zhang, Xinrun Du, Chenghua Lin, Wen-
816 hao Huang, Wenhui Chen, Jie Fu, et al. 2024c.
817 Kun: Answer polishment for chinese self-alignment
818 with instruction back-translation. *arXiv preprint*
819 *arXiv:2401.06477*.

820 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao
821 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
822 Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis,
823 Luke Zettlemoyer, and Omer Levy. 2023. [Lima: Less
824 is more for alignment.](#) *Preprint*, arXiv:2305.11206.

825 He Zhu, Wenjia Zhang, Nuoxian Huang, Boyang Li,
826 Luyao Niu, Zipei Fan, Tianle Lun, Yicheng Tao, Jun-
827 you Su, Zhaoya Gong, Chenyu Fang, and Xing Liu.
828 2024. [Plangpt: Enhancing urban planning with tai-
829 lored language model and efficient retrieval.](#) *Preprint*,
830 arXiv:2402.19273.

831 A Experiment Baselines

832 • Alpaca-52k (Taori et al., 2023). This dataset is developed by Stanford University using Text-Davinci-003.
833 It encompasses 52,002 instruction-following samples.

834 • Alpaca-GPT4 (Peng et al., 2023). This dataset contains English Instruction-Following Data generated
835 by GPT-4 using Alpaca prompts for fine-tuning LLMs. It encompasses 52,002 instruction-following
836 samples, the same as Alpaca-52k.

837 • Alpaca-Cleaned. This is a cleaned version of the Alpaca-GPT4 Dataset to address problems like
838 hallucinations, merged instruction, and so on. It encompasses 51,760 instruction-following samples.

839 • LIMA (Zhou et al., 2023). This is a dataset of 1,000 prompts and responses from a variety of sources, pri-
840 marily split into community Q&A forums and manually authored examples, where the outputs (responses)
841 are stylistically aligned with each other, but the inputs (prompts) are diverse.

842 • WizardLM-70k (Xu et al., 2023). This dataset employs the Evol-Instruct algorithm to enhance the
843 quality of instruction data. Incorporating ChatGPT during the reformulation phase ensures the data fidelity.
844 Among its 250,000 instructions, we primarily focused on the WizardLM-7b subset, which consists of
845 70,000 samples.

846 • Muffin (Lou et al., 2024) MUFFIN’s data curation includes input sampling, instruction collection via
847 two methods, output annotation by ChatGPT/GP4-4, instruction filtering, and classification expansion.
848 This is a large dataset of 68k training instances.

849 • ShareGPT (Chiang et al., 2023). This is a human-annotated dataset consisting of approximately 70K
850 user-shared conversations collected from ShareGPT.

851 • Humpback. This self-alignment method generates instruction data through reverse fine-tuning.

852 B FANNO Details

853 B.1 Pre-screen Details

854 Our objective was to efficiently enhance the selection process, minimizing time spent while maximizing
855 quality outcomes. Initially, we employed **Mistral-7b-instruct-v2** (Jiang et al., 2023a) to evaluate texts for
856 repetitive content, personal privacy concerns, specific themes, and advertising, using prompts to guide
857 scoring and annotation (see Table 7). For diversity assessment, we utilized a fast community detection
858 algorithm 1 with hyperparameters set to $k = 2$ and $\text{simratio} = 0.7$ (k : the minimum size of a community;
859 simratio : controls the similarity threshold, Only node pairs with similarity scores higher than this threshold
860 are considered connected), facilitating the classification of half a million entries within minutes. The
861 model paraphrase-MiniLM-L6-v2 (Reimers and Gurevych, 2019) is used for text embedding. For larger
862 datasets, texts were segmented into groups for individual community detection analyses. After the pre-
863 screening process, *Pre-Screen Data* has approximately 30k records, which is 6% of the original. This
864 stage was designed to balance the trade-off between processing speed and analytical precision, prioritizing
865 efficiency over exhaustive detail examination.

866 B.2 UCB Bootstrap

867 The setup comprises a language model G parameterized by θ_G for generating instructions, a critic model J
868 parameterized by θ_J for evaluating instruction quality, as well as a document set \mathcal{D} , a subset \mathcal{D}' , task-type
869 tags $\mathcal{T}_{\mathcal{T}}$, and difficulty-level tags $\mathcal{T}_{\mathcal{D}}$.

870 The procedure is as follows:

871 1. Initialization:

$$872 S \leftarrow \emptyset$$

2. Seed Generation (*SeedGen*):

$$\begin{aligned} \forall d \in \mathcal{D}', \text{ generate } x_i &\sim P(x|d, t; \theta_G) \\ \text{where } t &\sim \mathcal{U}(\mathcal{T}_{\mathcal{T}} \times \mathcal{T}_{\mathcal{D}}) \\ S &\leftarrow S \cup \{x_i\} \end{aligned}$$

3. Instruction Augmentation (*InsAug*): For f rounds or until $|S|$ reaches a desired threshold:

- a. Select a subset $S' \subset S$ using the UCB strategy:

$$\begin{aligned} UCB(s) &= \bar{x}_s + C \sqrt{\frac{2 \ln N}{n_s}} \\ S' &= \{s_i | s_i \in S, UCB(s_i) \text{ is maximized}\} \end{aligned}$$

where \bar{x}_s is the average quality score of instruction s , N is the total number of selections, C is a hyper-parameter constant used to control exploration, t and n_s is the number of times instruction s has been selected.

- b. For each $s_i \in S'$, generate augmented instructions x' :

$$\begin{aligned} x' &\sim P(x|c, S'; \theta_G) \\ \text{s.t. } Sim(x'; s_i) &< \tau \end{aligned}$$

where τ is a similarity threshold.

- c. Update S with the augmented instructions:

$$S \leftarrow S \cup \{x'\}$$

B.3 Fast Community Detection Algorithm

As Algorithm 1 has shown, the Fast Community Detection Algorithm is used to cluster the embeddings of instructions processed by SentenceTransformer (Reimers and Gurevych, 2019), which can then represent the diversity of instructions. Specifically, Fast Community Detection works by iteratively identifying groups of data points (embeddings of sentences) that are closely related based on a predefined similarity threshold, efficiently leveraging cosine similarity calculations. It prioritizes larger communities while minimizing overlapping clusters to produce meaningful community structures.

C Experiment Setting Detail

We chose LoRA over full fine-tuning due to similar performance observed in preliminary experiments, with computational constraints being the primary factor influencing this decision.

We use the same hyperparameters as existing supervised instruction tuning methods (Chiang et al., 2023; Raschka, 2023). Specifically, we use cosine learning rate scheduling with a starting learning rate of 2×10^{-5} and a weight decay of 0.1. The batch size is 32 and the dropout rate is 0.1. For the LoRA configuration, we employ a rank of 256 and set α to 512, with an initial learning rate of 5×10^{-5} . We utilize 8 NVIDIA 4090 GPUs to train our model.

D Prompt Templates Used in FANNO

D.1 Text Filtering

Algorithm 1 Fast Community Detection (Reimers and Gurevych, 2019)

```
1: function COMMUNITYDETECTION(embeddings, threshold, min_community_size, batch_size)
2:   Normalize embeddings
3:   Initialize extracted_communities as empty list
4:   for start_idx in range(0, length(embeddings), batch_size) do
5:     Compute cosine similarity scores for batch starting from start_idx
6:     Find top-k values from cosine similarity scores
7:     for  $i$  in range(length(top_k_values)) do
8:       if last element of  $i$ -th top-k values  $\geq$  threshold then
9:         Find top-k most similar entries for  $i$ -th element
10:        while last element of top-k values  $>$  threshold and sort_max_size  $<$  length of
embeddings do
11:          Increase sort_max_size if needed
12:        end while
13:        Add indices of entries with similarity  $\geq$  threshold to extracted_communities
14:      end if
15:    end for
16:  end for
17:  Sort extracted_communities by size
18:  Remove overlapping communities from extracted_communities
19:  return extracted_communities
20: end function
```

Table 7: Prompts for Pre-Screen

<p>You are act as a assistant to check useless, informal or ambiguous information. Let's think step by step. The objective is to meticulously inspect the text to determine if it is useless, informal or ambiguous text (e.g. random characters, ambiguous paragraph, broken sequence, informally organized text, etc.) Your response should be '1' (yes) if the text contains useless, informal or ambiguous information, or '0' (no) if it does not, without providing any reasoning and explanation.</p> <p>### Document: {doc}</p> <p>### Answer:</p>
<p>You are act as a assistant to check privacy information. Let's think step by step. The objective is to meticulously inspect the text to determine if it contains any privacy information (e.g. human names, phone numbers, addresses, etc.). Your response should be '1' (yes) if the text contains privacy information, or '0' (no) if it does not, without providing any reasoning and explanation.</p> <p>### Text: {doc}</p> <p>### Answer:</p>
<p>I want you to act as an advertisement evaluator. Let's think step by step. The objective is to meticulously inspect the text based on certain characteristics and decide whether it is an advertisement or not. Your response should be '1' (yes) if the text is an advertisement, or '0' (no) if it is not, without providing any reasoning and explanation.</p> <p>Evaluate the text considering these characteristics:</p> <ul style="list-style-type: none">- Promotional language or sales pitch- Mention of product or service benefits- Call to action (e.g., "Buy now", "Subscribe")- Pricing information or special offers- Contact information or links for more details <p><Answer Format>: 1 or 0</p> <p>### Text: {text}</p> <p>### Answer:</p>

Table 8: Prompts for instruction generation filter

<p>I want you to act as an instruction evaluator. Please evaluate this instruction and respond with '0' (bad) or '1' (good), without giving reasons. Standard: A good instruction Must not involve recent or current events. Historical events are fine. Example1: Instruction: Please analyze the recent COVID-19 outbreak. Answer: 0 (Reason: recent) Example2: Instruction: What's happening in China in September 2023? Answer: 0 (Reason: in September 2023) Example3: Instruction: Provide an account of events from last Monday night. Answer: 0 (Reason: last Monday night)</p> <p>### Instruction: {instruction} ### Answer:</p>
<p>I want you to act as a instruction evaluator. Please evaluate this instruction and respond with '0' (bad) or '1' (good), without giving reasons. Standard: A good instruction must not include any private information like names, addresses, phone numbers, etc, unless the person is historical or famous. Example1: Instruction: What is the name of the person who lives at 123 Main Street? Answer: 0 (Reason: private information) Example2: Instruction: What is the name of the first president of the United States? Answer: 1 (Reason: historical) Example3: Instruction: What is the address of the CEO of Microsoft? Answer: 0 (Reason: private information)</p> <p>### Instruction: {instruction} ### Answer:</p>
<p>I want you to act as a instruction evaluator. Please evaluate this instruction and respond with '0' (bad) or '1' (good), without giving reasons. Standard: A good instruction is perfectly logical, and practical, and can be fully understood by a human. A bad instruction, likely generated by AI, is generally vague, weird, complex, and long. It may seem to string unrelated words, topics, and tasks together. Example1: Instruction: Considering the health benefits of a non-dairy diet, how does the emotional response of individuals vary when they attend social events where dairy-based foods are served? Answer: 0 Example2: Instruction: Create a multidisciplinary essay that explores the and historical origins of the dish 'Shrimp Alfredo Pasta Bake'. Discuss the various ingredients, their origins. Additionally, translate the recipe instructions from English to Spanish. Answer: 0</p> <p>### Instruction: {instruction} ### Answer:</p>

Table 7 shows the prompts for basic filtering, including filtering information with useless information, privacy information, or advertisement.

908
909

Table 8 shows the prompts for instruction generation filtering, including filtering instructions that are time-sensitive, asking for private information, or not answerable.

910
911

D.2 Complexity and Quality Scorer

912

Table 9: Prompt for quality scorer

```
You are a helpful assistant. Please identify the quality score of the Response corresponding to the Question.
### Question:
{instruction}
### Response:
{output}
### Quality:
```

Table 10: Prompt for complexity scorer

```
You are a helpful assistant. Please identify the complexity score of the following user query.
### Query:
{instruction}
### Complexity:
```

As Table 9 and 10 have shown, the prompts are provided to deita-complexity-scorer and deita-quality-scorer model (Liu et al., 2023a).

913
914

D.3 Generating Instruction Pairs

915

FANNO employs 2 ways to generate instruction response:

916

- Question, Document to Answer: model infers answer with both question and related document.
- Question to Answer: model infers the answer directly with the question, using its own knowledge.

917

918

Table 11: Question, Document to Answer

```
You are a helpful, respectful and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense, or is not factually coherent, explain why instead of answering something not correct. If you don't know the answer to a question, please don't share false information.
### Instruction: {question}.
### Paragraph: {doc}.
### Response:
```

D.4 Seed Generation

919

Listing 1: Seed Generation

```
1 def seed_gen(text):
```

920
921

Table 12: Question to Answer

You are a helpful, respectful and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature. If a question does not make any sense, or is not factually coherent, explain why instead of answering something not correct. If you don't know the answer to a question, please don't share false information.

QUESTION: {question}

Response:

```

922 2 reasoning_tag = "It should be complex and requires multiple-step reasoning to
923 solve."
924 3 critical_thinking_tag = "It demands critical thinking skills to analyze from
925 various perspectives and evaluate multiple solutions."
926 4 creativity_tag = "It necessitates creative thinking to devise innovative
927 solutions beyond conventional approaches."
928 5 interdisciplinary_tag = "It demands integrating knowledge from diverse
929 disciplines to address its multifaceted nature."
930 6 command_tag = "It should be in the style of a command or imperative. For example
931 , 'Write a paragraph about...' or 'Describe the...'"
932 7 question_tag = "It should be in the style of a question or interrogative. For
933 example, 'What is the...?' or 'How do you...?'"
934 8
935 9 nli_tag = "It is a Natural language inference question: Assessing if evidence
936 supports a conclusion."
937 10 commonsense_tag = "It is a Commonsense question: Predicting outcomes based on
938 everyday knowledge."
939 11 sentiment_tag = "It is a Sentiment analysis question: Determining emotional
940 response to a given scenario."
941 12 paraphrase_tag = "It is a Paraphrasing question: Rewording a statement while
942 retaining its meaning."
943 13 close_book_qa_tag = "It is a Close-book QA question: Answering factual queries
944 using pre-existing knowledge."
945 14 struc2text_tag = "It is a Structure to text question: Describing a process or
946 concept in written form."
947 15 summarization_tag = "It is a Summarization question: Condensing key information
948 from a larger text."
949 16 translate_tag = "It is a Translation question: Converting text from one language
950 to another."
951 17 implicit_reasoning_tag = "It is a Implicit reasoning question: Inferring reasons
952 behind common behaviors."
953 18 text_category_tag = "It is a Text categorization question: Identifying defining
954 characteristics of a given text type."
955 19
956 20 tags = [reasoning_tag, critical_thinking_tag, creativity_tag,
957 interdisciplinary_tag]
958 21 classify = [nli_tag, commonsense_tag, sentiment_tag, paraphrase_tag,
959 close_book_qa_tag, struc2text_tag, summarization_tag, translate_tag,
960 implicit_reasoning_tag, text_category_tag]
961 22 types = [command_tag, question_tag]
962 23
963 24 QUESTION_TEMPLATE = """You're proficient in crafting complex question. Generate
964 only one question that adheres to the provided #Paragraph#.
965 The question should meet the following criteria:
966 25 0. The person answering the question cannot see the #Paragraph#[SYSTEM:
967 IMPORTANT], so the question must not contain phrases like 'Given the
968 information provided', 'Based on the provided information', or similar
969 expressions that imply direct citations or references from #Paragraph#.
970 27 1. {characteristic}.
971 28 2. {type}.
972 29 3. {classify}.
973 30
974 31 ### Paragraph:
975 32 {text}

```

33	<pre> ### Question: """ prompts = [QUESTION_TEMPLATE.format(characteristic=tag, type=type, text=text, classify=c) for tag in tags for c in classify for type in types] return prompts </pre>	976 977 978 979 980
----	--	---------------------------------

Code 1 shows the process of generating seed with sampled tags, including task types and difficulty levels. 982

D.5 Think Different Prompt 983

Table 13: Prompt for Think Differently

```

You are a helpful assistant. Your task is to conceive a complex question inspired from the Paragraph,
while ensuring it is completely different from the example provided below. Prohibit the use of
expressions, question types, and initial verbs that are identical to those in the Examples provided.
Avoid phrases such as 'Based on', 'Given the information provided', 'Using the data' or any similar
expressions that suggest references to the Paragraph. command
### Counterexample:
<Example1>: {seed1}
<Example2>: {seed2}
<Example3>: {seed3}
<Example4>: {seed4}
<Example5>: {seed5}

### Paragraph:
{text}

### Question:

```

D.6 Self-Instruct Prompting Templates for Data Generation 984

Self-Instruct relies on the following prompting template in order to elicit the generation from language 985
models. 986

```

Come up with a series of tasks:

Task 1: {instruction for existing task 1}
Task 2: {instruction for existing task 2}
Task 3: {instruction for existing task 3}
Task 4: {instruction for existing task 4}
Task 5: {instruction for existing task 5}
Task 6: {instruction for existing task 6}
Task 7: {instruction for existing task 7}
Task 8: {instruction for existing task 8}
Task 9:

```

Table 14: Prompt used for *Self-Instruct*

D.7 Faithfulness Evaluation 987

Table 15 shows the prompt to select more faithful instruction. The prompt originates from (Li et al., 2024) 988
with minor modifications. 989

Table 15: Prompt for Faithfulness Evaluation (Li et al., 2024)

Below is an instruction from an user and a candidate answer. Let’s think step by step. Evaluate whether or not the answer is a good example of how AI Assistant should respond to the user’s instruction. Please assign a score using the following 5-point scale: 1: It means the answer is incomplete, vague, off-topic, or not exactly what the user asked for. For example, some content seems missing. Or the response is from another person’s perspective with their personal experience (e.g. taken from blog posts). Or it contains promotional text or other irrelevant information. 2: (between 1 and 3) 3: It means the answer is helpful but not written by an AI Assistant. It addresses all the basic asks from the user. It is complete and self contained with the drawback that the response is not written from an AI assistant’s perspective, but from other people’s perspective. For example, it contains personal experience or opinion, mentions comments section, or share on social media, etc. 4: (between 3 and 5) 5: It means it is a perfect answer from an AI Assistant. It has a clear focus on being a helpful AI Assistant, where the response looks like intentionally written to address the user’s question or instruction without any irrelevant sentences. The answer provides high quality content, demonstrating expert knowledge in the area, is very well written, logical, easy-to-follow, engaging and insightful.
 Your reply should be only 1 or 2 or 3 or 4 or 5, without providing any reasoning and explanation.
 ### Instruction: {instruction}
 ### Answer: {response}
 ### Your Reply:

E Data Analysis

E.1 Quality, Length and Diversity

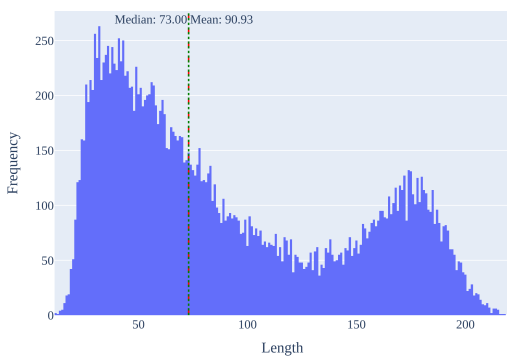


Figure 5: FANNO Instruction Length Distribution

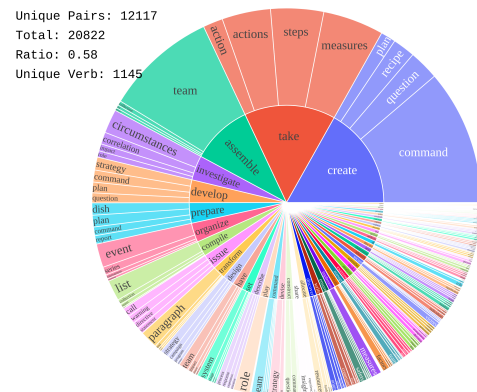


Figure 6: Top 50 common verbs and their corresponding nouns in FANNO

Figures 5 and 7 show the instruction length distribution of FANNO and Alpaca-Cleaned, respectively. It is worth noting that the mentioned length includes both the instruction and input combined. Figures 6 and 8 show the verb-noun diversity of FANNO and Alpaca-Cleaned, respectively. Figure 9 shows the comparison of quality and complexity between FANNO and Alpaca-Cleaned.

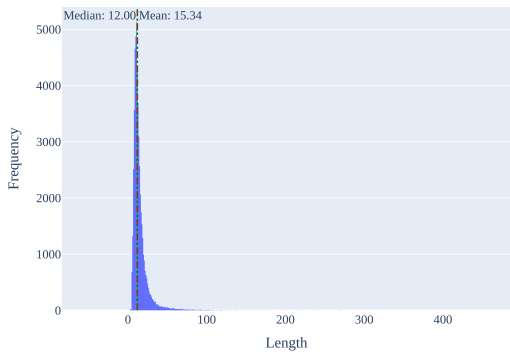


Figure 7: Alpaca-Cleaned Instruction Length Distribution

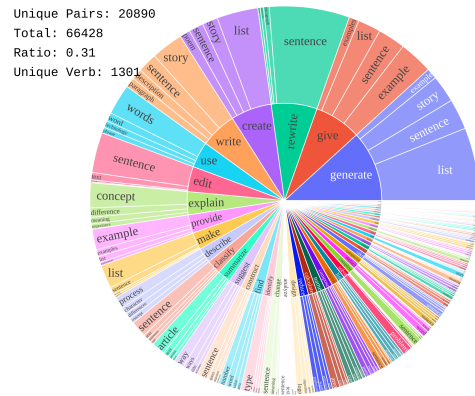


Figure 8: Top 50 common verbs and their corresponding nouns in Alpaca-Cleaned

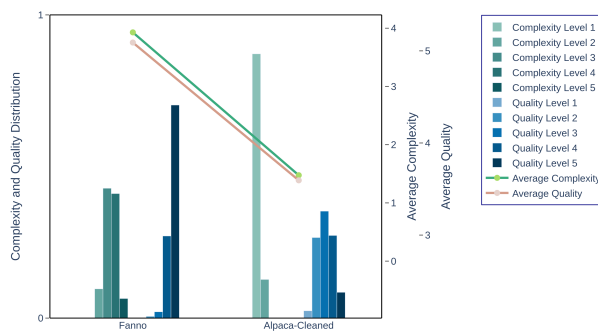


Figure 9: Quality and Complexity Comparison between FANNO and Alpaca-Cleaned

F Human Evaluation 996

F.1 Complexity Level 997

The first tier (0 point) pertains to instructions that exhibit apparent issues, such as being unanswerable or containing missing information. The second tier (1 point) involves instructions that can be answered using everyday knowledge. These instructions may assess basic skills, analyze human emotional experiences, or organize activities without requiring much specialized knowledge in a particular field. The third tier (2 points) comprises expert instructions. These necessitate specialized knowledge and require thorough deliberation steps to fulfill the instruction's requirement. 998 999 1000 1001 1002 1003

F.2 Instruction Complexity Human Evaluation 1004

2 point (Expert Level) 1005

1. Create a series of interactive exercises for a group of advanced French learners to practice the conditional tense, incorporating a variety of verb forms and sentence structures, while also encouraging them to engage in peer-to-peer learning and problem-solving. Consider using a combination of written and oral activities, and provide clear instructions and examples for each exercise. Additionally, design a system for assessing their progress and providing personalized feedback. 1006 1007 1008 1009 1010
2. How can we optimize the WordPress website's performance for logged-in users without employing the Auto-Cache Engine? Consider various caching strategies and evaluate their potential impact on user experience and website functionality. 1011 1012 1013
3. Design a multifaceted approach to streamline the patient registration process for a healthcare facility, ensuring adherence to ICD-10 and CPT coding standards, while providing exceptional customer 1014 1015

1016 service to a diverse patient population. Consider implementing innovative technologies and collabo-
1017 rating with various departments to optimize workflows and enhance overall efficiency. Evaluate the
1018 potential impact of this approach on patient satisfaction, staff morale, and financial performance.

1019 4. Assemble a team of data experts to evaluate the potential impact of a centralized data strategy on
1020 the decision-making process of a tech startup, considering the long-term benefits and potential
1021 drawbacks. Analyze various case studies of successful companies, such as Google, Apple, Amazon,
1022 and Facebook, to identify key strategies and best practices for implementing a data-driven culture.
1023 Evaluate the role of immediate returns versus long-term benefits in the adoption of data-driven
1024 decision-making and provide recommendations for managing potential challenges, such as data
1025 security and privacy concerns.

1026 5. Assemble a team of nutritionists and chefs to devise a creative and nutritious menu for a charity
1027 gala, utilizing natural sweeteners as the primary ingredient in each dish, while ensuring that the final
1028 creations are visually appealing and can be prepared in large quantities. Additionally, consider the
1029 dietary restrictions of various attendees and incorporate alternative options for those with gluten,
1030 dairy, and nut allergies. The team should also aim to minimize food waste and maximize the use of
1031 locally sourced ingredients.

1032 **1 point (Everyday Level)**

- 1033 1. Develop a weekly routine that integrates both your professional and personal commitments, ensuring
1034 that you effectively manage your time and accomplish your goals. What unique strategies could you
1035 employ to optimize your productivity during your weekly review and planning session?
- 1036 2. Persuade your employer to grant you the flexibility to work from home for a specified number of
1037 days per week, demonstrating the potential time and cost savings, as well as the potential benefits to
1038 your overall well-being.
- 1039 3. Capture the essence of a cherished memory by taking a photograph of a cherished photograph.
1040 Ensure the image is visually appealing and evokes a sense of nostalgia.
- 1041 4. Translate the following paragraph from English to another language of your choice. Ensure that the
1042 translation conveys the original meaning and intent. "Analyze the artworks displayed at the exhibition
1043 from various perspectives. Which artwork resonates the most with the theme of environmental
1044 conservation? Provide reasons for your answer."

1045 **0 point (Bad)**

- 1046 1. Utilize the data from the Neighbourhood Forum Launch event to determine the percentage of
1047 attendees who were members prior to the event and the percentage who joined during the event.
1048 Additionally, identify the top three focus groups with the highest number of attendees and determine
1049 the average number of attendees per focus group. Finally, calculate the total number of attendees
1050 who placed a dot on the Forum map and the percentage of attendees who did so.