

©Copyright 2025

Zhangchen Xu

Magpie: Generating High-Quality Synthetic Data with Open-Source Large Language Models

Zhangchen Xu

A thesis
submitted in partial fulfillment of the
requirements for the degree of

Master of Science

University of Washington

2025

Committee:

Radha Poovendran

Payman Arabshahi

Program Authorized to Offer Degree:
Electrical and Computer Engineering

University of Washington

Abstract

Magpie: Generating High-Quality Synthetic Data
with Open-Source Large Language Models

Zhangchen Xu

Chair of the Supervisory Committee:
Radha Poovendran
Department of Electrical and Computer Engineering

High-quality instruction data is critical for aligning large language models (LLMs). Although some models, such as Llama-3-Instruct, have open weights, their alignment data remain private, which hinders the democratization of AI. High human labor costs and a limited, predefined scope for prompting prevent existing open-source data creation methods from scaling effectively, potentially limiting the diversity and quality of public alignment datasets. Is it possible to synthesize high-quality instruction data at scale by extracting it directly from an aligned LLM? We present a *self-synthesis* method for generating large-scale alignment data named MAGPIE. Our key observation is that aligned LLMs like Llama-3-Instruct can generate a user query when we input only the pre-query templates up to the position reserved for user messages, thanks to their auto-regressive nature. We use this method to prompt Llama-3-Instruct and generate 4 million instructions along with their corresponding responses. We further introduce extensions of MAGPIE for filtering, generating multi-turn, preference optimization, domain-specific and multilingual datasets. We perform a comprehensive analysis of the MAGPIE-generated data. To compare MAGPIE-generated data with other public instruction datasets (e.g., ShareGPT, WildChat, Evol-Instruct, UltraChat, OpenHermes, Tulu-V2-Mix, GenQA), we fine-tune Llama-3-8B-Base with each dataset and evaluate the performance of the fine-tuned models. Our results indicate that using

MAGPIE for supervised fine-tuning (SFT) solely can surpass the performance of previous public datasets utilized for both SFT and preference optimization, such as direct preference optimization with UltraFeedback. We also show that in some tasks, models supervised fine-tuned with MAGPIE perform comparably to the official Llama-3-8B-Instruct, despite the latter being enhanced with 10 million data points through SFT and subsequent preference optimization. This advantage is evident on alignment benchmarks such as AlpacaEval, ArenaHard, and WildBench.

CONTENTS

List of Figures	v
List of Tables	vii
Chapter 1: Introduction	1
Chapter 2: Related Work	4
Chapter 3: MAGPIE: A Scalable Method to Synthesize Alignment Data	6
3.1 MAGPIE Pipeline	6
3.2 MAGPIE Extensions	8
Chapter 4: Dataset Analysis	10
4.1 Dataset Coverage	10
4.2 Dataset Attributes	11
4.3 Safety Analysis	13
4.4 Cost Analysis	13
4.5 Additional Analysis	14
Chapter 5: Performance Analysis	15
5.1 Experimental Setups	15

5.2 Experimental Results	17
Chapter 6: Conclusion	24
Chapter 7: Acknowledgment	25
Appendix A: Statistics of instruction datasets generated by MAGPIE compared to other instruction datasets.	38
Appendix B: Filter Setups	40
Appendix C: More Dataset Analysis	42
C.1 Additional Analysis on Dataset Coverage and Attributes.	42
C.2 Additional Safety Analysis	46
C.3 Ablation Analysis on Generation Configurations	47
C.4 Impact of Annotating Models	50
C.5 Contamination Analysis	51
Appendix D: Detailed Experimental Setups	52
D.1 Experimental Setups for Generating MAGPIE-Air and MAGPIE-Pro	52
D.2 Experimental Setups for Instruction Tuning and Preference Tuning	52
Appendix E: Additional Experimental Results	56
E.1 Performance of MAGPIE-MT	56
E.2 Compare MAGPIE and Self-Instruct using Llama-3-8B-Instruct	56
E.3 Performance of domain-specific and multilingual MAGPIE datasets	57

E.4	Ablation Analysis on Data Quantity and Quality	58
E.5	Ablation Analysis on Filter Designs	60
E.6	Ablation Analysis on Response Generator	60
E.7	Trustworthiness of MAGPIE-Aligned Models	61
E.8	IFEval Evaluations of MAGPIE-Aligned Models and Baselines	61
E.9	Ablation Analysis on the Impact of Reward Models on DPO Performance . .	63
Appendix F: Prompt Templates		65
F.1	Prompt Templates for MAGPIE Extension	65
F.2	Prompt Templates for Evaluation	65
Appendix G: MAGPIE Examples		71

LIST OF FIGURES

Figure Number	Page
1.1 This figure illustrates MAGPIE, the process of self-synthesizing alignment data from aligned LLMs (e.g., Llama-3-8B-Instruct) to create a high-quality instruction dataset. In Step 1, we input only the pre-query template into the aligned LLM and generate an instruction along with its response using auto-regressive generation. In Step 2, we use a combination of a post-query template and another pre-query template to wrap the instruction generated from Step 1, prompting the LLM to generate the response. This completes the construction of the instruction dataset. MAGPIE efficiently generates diverse and high-quality instruction data, which can be further extended to multi-turn (MAGPIE-MT), preference optimization (MAGPIE-DPO), domain-specific, and multilingual datasets.	2
4.1 The statistics of instruction difficulty and quality.	11
4.2 This figure summarizes the minimum neighbor distances and reward differences.	13
5.1 This figure shows the performance breakdown by category of MAGPIE-Pro and baselines on WildBench.	19
C.1 Lengths of instructions and responses in MAGPIE-Air/Pro.	43
C.2 This figure compares the t-SNE plot of MAGPIE-Pro with those of Alpaca, Evol Instruct, and UltraChat, each of which is sampled with 10,000 instructions. The t-SNE plot of MAGPIE-Pro encompasses the area covered by the other plots, demonstrating the comprehensive coverage of MAGPIE-Pro. . . .	44
C.3 This figure compares the UMAP plot of MAGPIE-Air with those of Alpaca, Evol Instruct, and UltraChat, each of which is sampled with 10,000 instructions. The UMAP plot of MAGPIE-Air encompasses the area covered by the other plots, demonstrating the comprehensive coverage of MAGPIE-Air. . . .	45
C.4 This figure visualizes the task category of MAGPIE-Pro and MAGPIE-Air by topic tags.	46

C.5	This figure demonstrates the top 20 most common root verbs (shown in the inner circle) and their top 5 direct noun objects (shown in the outer circle) within the MAGPIE-Air dataset. This indicates that MAGPIE encompasses a broad range of topics.	47
C.6	This figure illustrates the impact of varying decoding parameters on the quality, difficulty, and diversity of generated instructions. We observe that while higher temperature and top-p values may decrease the overall quality, they tend to increase both the difficulty and diversity of the instructions.	48
C.7	This figure compares the input quality and difficulty with and without system prompts.	49
C.8	This figure compares the impact of different annotators on evaluating the instruction quality and difficulty.	50
F.1	Prompt for generating MAGPIE-MT. We take Llama-3-8B-Instruct as an example. The placeholder <code>{instruction}</code> and <code>{response}</code> are from the first turn.	65
F.2	Prompts for controlling instruction generation tasks. These examples illustrate how to guide Llama-3-8B-Instruct in generating instructions for specific domains: mathematics, coding, translation, and multilingual tasks. To adapt this approach for different instruction tasks, replace the System Prompt placeholder in the System Prompt Template with the appropriate domain-specific prompt.	66
F.3	Prompt for generating task categories	68
F.4	Prompt for generating quality of instructions	69
F.5	Prompt for generating difficulty of instructions	70

LIST OF TABLES

Table Number	Page	
5.1	This table compares the performance of models instruction-tuned on the Llama-3-8B base models using MAGPIE-generated datasets and baseline datasets. We observe that models aligned with our datasets significantly outperform those aligned with baseline datasets of the same order of magnitude in terms of data size. In addition, our fine-tuned models achieve comparable performance to the official aligned model, despite only undergoing SFT with a much smaller dataset. Numbers in bold indicate that MAGPIE outperforms the official Llama-3-8B-Instruct model.	18
5.2	This table compares the performance of models instruction-tuned on the Qwen base models using the MAGPIE-Pro-300K-Filtered dataset and the official instruction-tuned models. The Qwen base model enhanced with MAGPIE outperforms the official instruction-tuned model.	21
5.3	This table compares the performance of models supervised-fine-tuned on MAGPIE-Air, MAGPIE-Pro, and MAGPIE-Pro-Mix against baselines and official instruct model across various downstream benchmarks. All models are supervised-fine-tuned on the Llama-8B base models.	22
A.1	Statistics of MAGPIE family compared to other instruction datasets. Tokens are counted using the <code>tiktoken</code> library [47]. Links to the MAGPIE datasets are provided in the text.	39
B.1	Different filter configurations we provide. We note that the Output Length filter is applied last. Specifically, this filter selects the k instances of the longest responses. In our experiments, we empirically set $\tau_1 = -12$, and $\tau_2 = 0$	41
C.1	Comparison of Topic Diversity across Different Synthetic Datasets.	43
C.2	This table shows the percentage of different unsafe categories of MAGPIE-Air and MAGPIE-Pro tagged by Llama-Guard-2 [54] model.	48
D.1	This table demonstrates the configurations of generating instructions of MAGPIE-Air and MAGPIE-Pro datasets with varying decoding parameters.	53
D.2	This table shows the hyper-parameters for supervised fine-tuning.	54

D.3	This table shows the hyper-parameters for direct preference optimization. . .	54
E.1	This table compares the performance of the multi-turn versions, MAGPIE-Air-MT and MAGPIE-Pro-MT, with their single-turn counterparts. All models are instruction-tuned on the Llama-8B base models.	56
E.2	This table compares the performance of models fine-tuned using 100K instruction-following datasets generated by Self-Instruct and MAGPIE. All models are supervised-fine-tuned on the Llama-8B base models. We observe that MAGPIE significantly outperforms Self-Instruct across all benchmarks.	57
E.3	Performance Comparison on HumanEval.	58
E.4	Performance Comparison on Chinese MT-Bench.	58
E.5	This table compares MAGPIE datasets within its family that differ in size, deployment of filtering, and models used to generate data. All models are supervised-fine-tuned on the Llama-8B base models.	59
E.6	This table compares the performance of different filter designs within MAGPIE-Pro. All models are supervised-fine-tuned on the Llama-8B base models. . .	60
E.7	This table compares the impact of different response generators on the model performance. All models are supervised-fine-tuned on the Llama-8B base models.	61
E.8	This table compares the performance of model supervised-fine-tuned using MAGPIE-Pro-300K-Filtered and the official Llama-3-8B-Instruct on the TrustLLM benchmark [27].	62
E.9	This table compares the performance of models fine-tuned using MAGPIE and other baseline datasets on the IFEval benchmark [78].	63
E.10	Comparison of Llama-3-8B DPO-trained models using different reward models. Metrics are AE2 LC, AE2 WR, and AH Score.	64

ACKNOWLEDGMENTS

I would like to express my deepest gratitude to my advisor Prof. Radha Poovendran for his invaluable guidance, support, and mentorship throughout the entire research process. His expertise, encouragement, and constructive feedback have been instrumental in shaping this thesis. I would also like to thank Prof. Payman Arabshahi for serving on my committee and providing valuable insights and feedback that enriched this research. I would also like to thank my family for their unwavering support, patience, and encouragement throughout my academic journey. Their love and belief in me provided the foundation that made this work possible.

Chapter 1

INTRODUCTION

Large language models (LLMs) such as GPT-4 [2] and Llama-3 [45] have become integral to AI applications due to their exceptional performance on a wide array of tasks by following instructions. The success of LLMs is heavily reliant on the data used for instruction fine-tuning, which equips them to handle a diverse range of tasks, including those not encountered during training. The effectiveness of instruction tuning depends crucially on access to high-quality instruction datasets. However, the alignment datasets used for fine-tuning models like Llama-3-Instruct are typically private, even when the model weights are open, which impedes the democratization of AI and limits scientific research for understanding and enhancing LLM alignment.

To address the challenges in constructing high-quality instruction datasets, researchers have developed two main approaches. The first type of method involves human effort to generate and curate instruction data [16, 30, 74, 75, 76], which is both *time-consuming* and *labor-intensive* [42]. In contrast, the second type of methods uses LLMs to produce synthetic instructions [19, 70, 36, 51, 52, 62, 64, 66, 67, 35]. Although these methods reduce human effort, its success heavily depends on prompt engineering and the careful selection of initial seed questions. The *diversity* of synthetic data tends to decrease as the dataset size grows. Despite ongoing efforts, the scalable creation of high-quality and diverse instruction datasets continues to be a challenging problem.

Is it possible to synthesize high-quality instructions at scale by directly extracting data from advanced aligned LLMs? A typical input to an aligned LLM contains three key components: the pre-query template, the query, and the post-query template. For instance, an input to Llama-2-chat could be “[INST] Hi! [/INST]”, where [INST] is the pre-query template

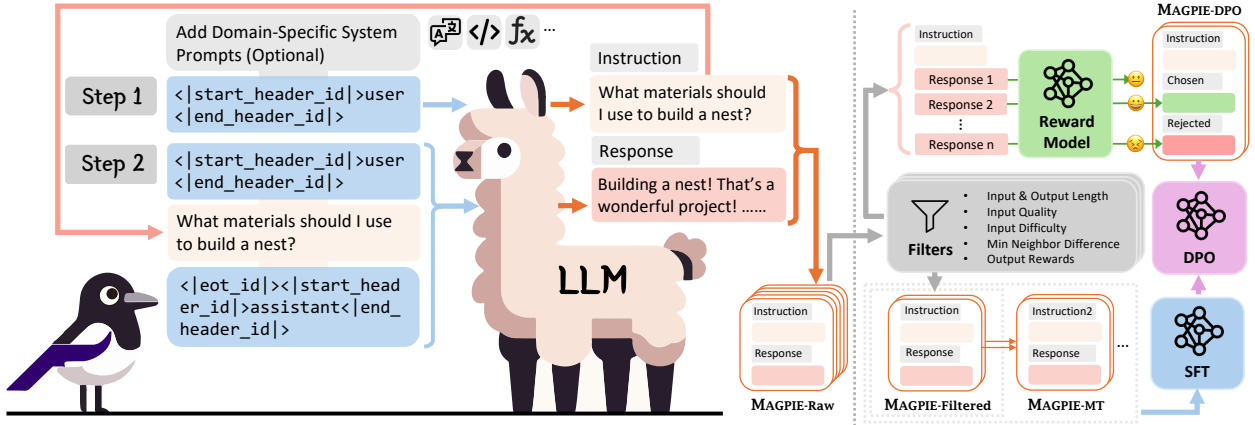


Figure 1.1: This figure illustrates MAGPIE, the process of self-synthesizing alignment data from aligned LLMs (e.g., Llama-3-8B-Instruct) to create a high-quality instruction dataset. In Step 1, we input only the pre-query template into the aligned LLM and generate an instruction along with its response using auto-regressive generation. In Step 2, we use a combination of a post-query template and another pre-query template to wrap the instruction generated from Step 1, prompting the LLM to generate the response. This completes the construction of the instruction dataset. MAGPIE efficiently generates diverse and high-quality instruction data, which can be further extended to multi-turn (MAGPIE-MT), preference optimization (MAGPIE-DPO), domain-specific, and multilingual datasets.

and `[/INST]` is the post-query template. These templates are predefined by the creators of the aligned LLMs to ensure the correct prompting of the models. We observe that when we only input the pre-query template to aligned LLMs such as Llama-3-Instruct, they *self-synthesize* a user query due to their auto-regressive nature. Our experiments indicate that these random user queries are of high quality and great diversity, suggesting that the abilities learned during the alignment process are effectively utilized.

Based on these findings, we developed a self-synthesis method to construct high-quality instruction datasets at scale, named MAGPIE (as illustrated in Figure 1.1). Unlike existing methods, our approach does not rely on prompt engineering or seed questions. Instead, it *directly* constructs instruction data by prompting aligned LLMs with a pre-query template

for sampling instructions. We also demonstrated the extensibility of MAGPIE in generating multi-turn, preference optimization, domain-specific, and multilingual datasets. We applied MAGPIE to the Llama-3-8B-Instruct and Llama-3-70B-Instruct models, creating two instruction datasets: MAGPIE-Air and MAGPIE-Pro, respectively.

Our MAGPIE-Air and MAGPIE-Pro datasets were created using 206 and 614 GPU hours, respectively, without any human intervention or API access to production LLMs like GPT-4. The statistics and advantages of MAGPIE datasets compared to existing ones are summarized in Table A.1 in Appendix A. We perform a comprehensive analysis of these two datasets in Section 4, allowing practitioners to filter and select data instances for fine-tuning models according to their particular needs.

To compare MAGPIE data with other public instruction datasets (e.g., ShareGPT [12], WildChat [74], Evol Instruct [66], UltraChat [19], OpenHermes [56, 55], GenQA [9], Tulu V2 Mix [28]), we conducted supervised fine-tuning (SFT) of the Llama-3-8B-Base model with each dataset and assess the performance of the fine-tuned models on alignment benchmarks such as AlpacaEval 2 [38], Arena-Hard [37], and WildBench [39]. Our results show that models supervised fine-tuned with MAGPIE achieve superior performance, even surpassing models that utilize both SFT and direct preference optimization (DPO) [49] with UltraFeed-back [15]. Notably, MAGPIE-aligned models outperform the official Llama-3-8B-Instruct model on AlpacaEval 2, despite the latter being fine-tuned with over 10 million data points for SFT and subsequent preference optimization. Not only does MAGPIE excel in SFT alone compared to prior public datasets, but also delivers the best results when combined with preference optimization methods such as DPO. By leveraging MAGPIE extensions to generate high-quality preference optimization datasets, MAGPIE-aligned Llama-3 models can even outperform GPT-4-Turbo(1106) on AlpacaEval 2. These findings show the exceptional quality of instruction data generated by MAGPIE, enabling it to outperform even the official, extensively optimized, and proprietary LLMs.

Chapter 2

RELATED WORK

Alignment Dataset Construction. We classify the existing methods of creating datasets for model alignment into two main categories: human interactions with LLMs and synthetic instruction generation. To create datasets for alignment, previous studies have collected **human** interactions with LLMs [16, 74, 75, 76, 30]. However, manually crafting instructions is not only time-consuming and labor-intensive, but may also incorporate toxic content [74]. Another category of approaches [62, 52, 66, 67, 64, 51] involves prompting LLMs to generate **synthetic** instruction datasets, beginning with a small set of human-annotated seed instructions and expanding these through few-shot prompting. However, these methods face a diversity challenge, as few-shot prompting often results in new instructions that are too similar to the original seed questions [36]. To enhance coverage, some research [19, 36] summarizes world knowledge and employs it to generate synthetic datasets. We note that our MAGPIE dataset also belongs to the synthetic dataset. However, we leverage the prompt template with no requirement for seed questions or prompt engineering.

Compared to the above two main categories, alignment data can also be generated by **transforming** existing datasets [63, 50, 23]. However, the constrained variety of NLP tasks in these datasets may impede the ability of tuned LLMs to generalize in real-world scenarios [36]. There are also **mixture** datasets (e.g., [28, 56, 43, 77]) that combine or select high-quality instruction data from various existing open-source instruction datasets to enhance coverage [28, 56] and/or improve overall performance [43, 77].

Training Data Extraction. Language models have the capability to memorize examples from their training datasets, potentially enabling malicious users to extract private information [7, 6, 8]. Pioneering work [31, 8, 46] has demonstrated that it is possible to

extract private pre-training data from BERT [18], GPT-2 [48], and ChatGPT [2], respectively. Yu et al. [71] propose several tricks including adjusting sampling strategies to better extract training datasets from language models. Recently, Kassem et. al. [29] propose a black-box prompt optimization method that uses an attacker LLM to extract high levels of memorization in a victim LLM. Wang et al. [61] leverage membership inference attack (MIA) to extract fine-tuning datasets from fine-tuned language models. Bai et al. [4] extracts the training dataset of production language models via special characters (e.g., structural symbols of JSON files, and , # in emails and online posts). Different from the aforementioned work, MAGPIE aims to create publicly available alignment datasets with minimal human effort by leveraging the remarkable generation capabilities of LLMs, rather than extracting private training data from LLMs.

Chapter 3

MAGPIE: A SCALABLE METHOD TO SYNTHESIZE ALIGNMENT DATA

Chat Templates of Aligned LLMs. For an aligned LLM (e.g., Llama-3-8B-Instruct), the input sequence can be represented as $x = T_{pre-query} \oplus q \oplus T_{post-query}$. Here, q is the user query (e.g., "What material should I use to build a nest?"), while $T_{pre-query}$ and $T_{post-query}$ are pre-query and post-query templates. The pre-query template shows up before the user query, and the post-query template is defined as the conversation template between the user query and the LLM response. These templates are defined by the model provider to ensure the correct prompting. For example, for Llama-3-8B-Instruct model,

$$T_{pre-query} = \langle |start_header_id| \rangle user \langle |end_header_id| \rangle,$$

$$T_{post-query} = \langle |eot_id| \rangle \langle |start_header_id| \rangle assistant \langle |end_header_id| \rangle.$$

3.1 Magpie Pipeline

Overview of Magpie. In what follows, we describe our lightweight and effective method, MAGPIE, to synthesize alignment data from aligned LLMs. An instance of instruction data consists of at least one or multiple instruction-response pairs. Each pair specifies the roles of instruction provider (e.g., user) and follower (e.g., assistant), along with their instruction and response. As shown in Figure 1.1, MAGPIE consists of two steps: (1) instruction generation, and (2) response generation. The MAGPIE pipeline can be fully *automated without any human intervention*, and can be readily adapted for the generation of multi-turn, preference, and domain-specific datasets, as detailed in Section 3.2. We describe each step in the following.

Step 1: Instruction Generation. The goal of this step is to generate an instruction for each instance of instruction data. Given an open-weight aligned LLM (e.g., Llama-3-70B-Instruct), MAGPIE crafts a pre-query template in the format of the predefined instruction template of the LLM. Note that the auto-regressive LLM has been fine-tuned using instruction data in the format of the pre-query template. Thus, the LLM autonomously generates an instruction when the pre-query template crafted by MAGPIE is given as an input. MAGPIE stops generating the instruction once the LLM produces an end-of-sequence token. Sending the crafted query to the LLM multiple times leads to a set of instructions. We note that compared with existing synthetic approaches [19, 36, 52, 62, 64, 66, 67], MAGPIE does not require specific prompt engineering techniques since the crafted query follows the format of the predefined instruction template. In addition, MAGPIE autonomously generates instructions without using any seed question, ensuring the diversity of generated instructions.

Step 2: Response Generation. The goal of this step is to generate responses to the instructions obtained from Step 1. MAGPIE sends these instructions to the LLM to generate the corresponding responses. Combining the roles of instruction provider and follower, the instructions from Step 1, and the responses generated in Step 2 yields the instruction dataset. We note that separating instruction and response generation offers several key advantages, including *more flexible generation configurations* between instructions and responses, where instruction generation benefits from higher temperature settings to maximize diversity, while response generation requires lower temperature for accuracy and reliability. In addition, this approach provides *modular flexibility*, allowing users to generate instructions independently and later pair them with responses from various sources.

Applicability of Magpie on Different LLMs. MAGPIE can be readily deployed to state-of-the-art open-weight language models including but not limited to Llama-3 [45], Llama-3.1/3.3 [21], Qwen2 [68], Qwen2.5 [69], Gemma-2 [53], and Phi-3 [1]. Please refer to Appendix A for detailed support and corresponding datasets.

Remark. MAGPIE generates high-quality instructions even when the instruction loss is masked during alignment. We hypothesize that LLMs retain an implicit memorization of

instruction distributions. We leave it as a potential future research problem.

3.2 Magpie Extensions

Dataset Filtering. MAGPIE allows practitioners to select instruction data from the raw dataset generated from the above two steps based on their needs. In Appendix B, we explore potential filter configurations with eight available metrics for users to customize their own MAGPIE datasets. We also provide 6 off-the-shelf filter configurations and discuss their performance in Appendix E.5.

Generating Multi-Turn Instruction Datasets. MAGPIE can be readily extended to generate multi-turn instruction datasets. To construct a multi-turn dataset (denoted as MAGPIE-MT), we initially follow Steps 1 and 2 to generate the first turn of instruction and response. For subsequent turns, we append the pre-query template to the end of the full prompt from the previous round of communication. We observe that the model may occasionally forget its role as the user, especially for the 8B model. To mitigate this, we employ a system prompt designed to control the behavior of the LLM and reinforce its awareness of the multi-round conversation context. The full prompt for building the instructions of MAGPIE-MT can be found in Figure F.1 in Appendix F. We follow the procedure in Step 2 of Section 3.1 to generate responses to form the multi-turn instruction dataset.

Generating Preference Optimization Datasets. Leveraging the diverse and high-quality instructions produced by MAGPIE, we present a simple and effective method for generating preference optimization data, inspired by (author?) [44] and (author?) [58]. We first select a small proportion of high-quality instructions from the raw dataset generated by the MAGPIE pipeline, ensuring diverse task categories. For each selected instruction, we sample responses from the aligned LLM k times, using a temperature of $T < 1$. We then employ a reward model (RM) to annotate scores for these responses. The response with the highest RM score is labeled as the chosen response, while the one with the lowest RM score is designated as the rejected response.

Generating Domain-Specific and Multilingual Datasets. In certain scenarios,

users may wish to fine-tune LLMs using domain-specific or multilingual instruction data to enhance performance within specific domains or languages. To address this need, we introduce a lightweight method to control both the task category and the language of generated instructions. Our approach involves guiding LLMs through a tailored system prompt, specifying that the model is a chatbot designed for a particular domain and outlining the types of user queries it might encounter. Examples of system prompts designed to control the generation of math, code, translation, and multilingual instructions are illustrated in Figure F.2 in Appendix F.

Furthermore, we note that domain-specific and multilingual instruction data can also be generated using models that are tailored to particular fields. MAGPIE demonstrates broad applicability beyond diverse chat models, extending to specialized code models (e.g., DeepSeek-Coder-V2 [79]) and math models (e.g., Qwen2-Math-7B-Instruct [68]). By leveraging the unique strengths and specializations of different models, MAGPIE can create a rich and diverse corpus of instructions. Examples of MAGPIE-generated instructions from different domain-specific models and multilingual models are provided in Appendix G.

Chapter 4

DATASET ANALYSIS

To demonstrate the effectiveness of MAGPIE compared with baseline methods for generating diverse high-quality alignment datasets, we apply MAGPIE to the Llama-3-8B-Instruct and Llama-3-70B-Instruct models [45] to construct two instruction datasets: Llama-3-MAGPIE-Air (hereafter referred to as MAGPIE-Air) and Llama-3-MAGPIE-Pro (hereafter referred to as MAGPIE-Pro), respectively. Examples of instances in both datasets can be found in Appendix G. In this section, we present a comprehensive analysis of the MAGPIE-Air and MAGPIE-Pro datasets, including topic coverage, difficulty, quality, similarity of instructions, and the quality of the responses.

4.1 Dataset Coverage

We follow (author?) [74] and analyze the coverage of MAGPIE-Pro in the embedding space. Specifically, we use the `all-mpnet-base-v2` embedding model¹ to calculate the input embeddings, and employ t-SNE [59] to project these embeddings into a two-dimensional space. We adopt three synthetic datasets as baselines, including **Alpaca** [52], **Evol Instruct** [66], and **UltraChat** [19], to demonstrate the coverage of MAGPIE-Pro. The detailed analysis can be found in Appendix C.1. We observe that the t-SNE plot of MAGPIE-Pro encompasses the area covered by the plots of Alpaca, Evol Instruct, and UltraChat. This suggests that MAGPIE-Pro provides a broader or more diverse range of topics. We also follow (author?) [62] and present the most common verbs and their top direct noun objects in instructions in Appendix C, indicating the diverse topic coverage of MAGPIE dataset.

¹<https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

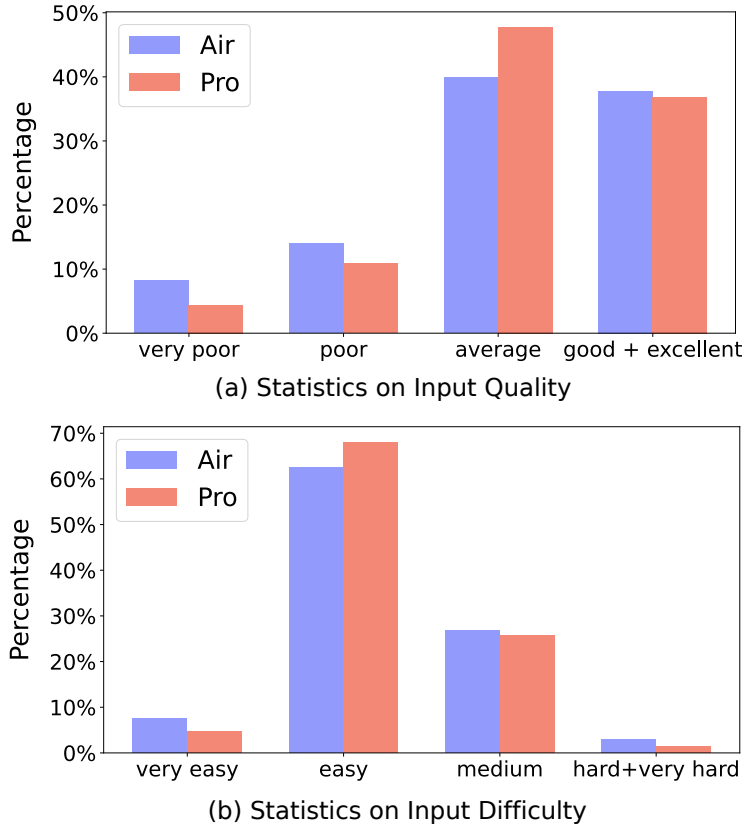


Figure 4.1: The statistics of instruction difficulty and quality.

4.2 Dataset Attributes

Attribute: Task Categories of Instructions. We use Llama-3-8B-Instruct to categorize the instances in MAGPIE-Pro (see Figure C.4 in Appendix C.1 for detail). The prompts used to query Llama-3-8B-Instruct can be found in Appendix F. Our observations indicate that over half of the tasks in MAGPIE-Pro pertain to information seeking, making it the predominant category. This is followed by tasks involving creative writing, advice seeking, planning, and math. This distribution over the task categories aligns with the practical requests from human users [38].

Attribute: Quality of Instructions. Similar to methods in [10], we prompt the

Llama-3-8B-Instruct model to assess the quality of each instruction in MAGPIE-Air and MAGPIE-Pro, categorizing them as ‘very poor’, ‘poor’, ‘average’, ‘good’, and ‘excellent’. We present the histograms of qualities for both datasets in Figure 4.1-(a). We have the following two observations. First, both datasets are of high quality, with the majority of instances rated ‘average’ or higher. In addition, the overall quality of MAGPIE-Pro surpasses that of MAGPIE-Air. We hypothesize that this is due to the enhanced capabilities of Llama-3-70B compared with Llama-3-8B.

Attribute: Difficulty of Instructions. We use the Llama-3-8B-Instruct model to rate the difficulty of each instruction in MAGPIE-Air and MAGPIE-Pro. Each instruction can be labeled as ‘very easy’, ‘easy’, ‘medium’, ‘hard’, or ‘very hard’. Figure 4.1-(b) presents the histograms of the levels of difficulty for MAGPIE-Air and MAGPIE-Pro. We observe that the distributions across difficulty levels are similar for MAGPIE-Air and MAGPIE-Pro. Some instructions in MAGPIE-Pro are more challenging than those in MAGPIE-Air because MAGPIE-Pro is generated by a more capable model (i.e., Llama-3-70B-Instruct).

Attribute: Instruction Similarity. We quantify the similarity among instructions generated by MAGPIE to remove repetitive instructions. We measure the similarity using **minimum neighbor distance** in the embedding space. Specifically, we first represent all instructions in the embedding space using the `all-mpnet-base-v2` embedding model. For any given instruction, we then calculate the minimum distance from the instruction to its nearest neighbors in the embedding space using Facebook AI Similarity Search (FAISS) [20]. The minimum neighbor distances of instructions in MAGPIE-Air after removing repetitions are summarized in Figure 4.2-(a).

Attribute: Quality of Responses. We assess the quality of responses using **rewards** assigned by a reward model, denoted as r^* . For each instance in our dataset, we also calculate **reward difference** as $r^* - r_{base}$, where r_{base} is the reward assigned by the same reward model to the response generated by the Llama-3 base model for the same instruction. We use URIAL [40] to elicit responses from the base model. A positive reward difference indicates that the response from our dataset is of higher quality, and could poten-

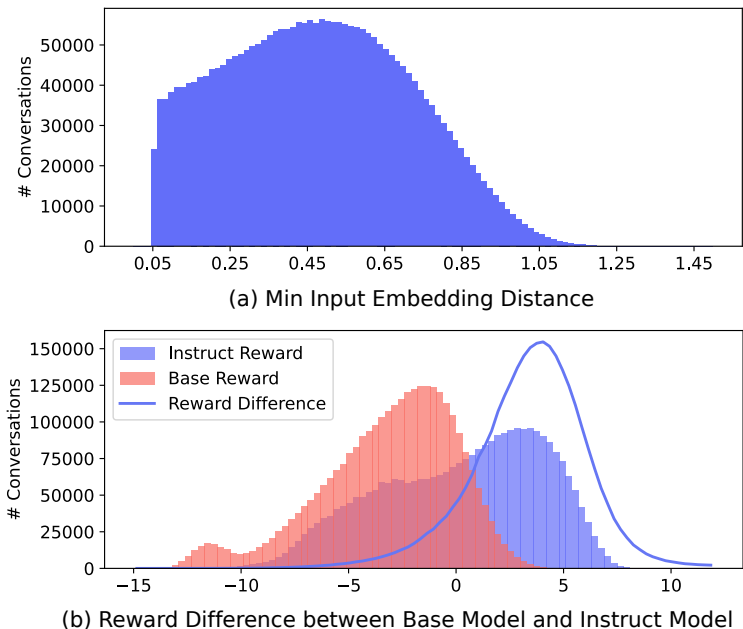


Figure 4.2: This figure summarizes the minimum neighbor distances and reward differences.

tially benefit instruction tuning. In our experiments, we follow **(author?)** [33] and choose FsfairX-LLaMA3-RM-v0.1 [65] as the reward model. Our results on the reward difference are presented in Figure 4.2-(b).

4.3 Safety Analysis

We use Llama-Guard-2 [54] to analyze the safety of MAGPIE-Air and MAGPIE-Pro. Our results indicate that both datasets are predominantly safe, with less than 1% of the data potentially containing harmful instructions or responses. Please see Appendix C.2 for detailed safety analysis.

4.4 Cost Analysis

We perform experiments on a server with four NVIDIA A100-SXM4-80GB GPUs, an AMD EPYC 7763 64-Core Processor, and 512 GB of RAM, using the VLLM inference framework

[32]. The models are loaded in the `bfloat16` format.

When creating the 3M MAGPIE-Air dataset, our MAGPIE spent 1.55 and 50 hours to generate the instructions (Step 1) and responses (Step 2), respectively. For the 1M MAGPIE-Pro dataset, MAGPIE used 3.5 and 150 hours to generate the instructions (Step 1) and responses (Step 2), respectively. Compared to existing approaches to create instruction datasets, the pipeline of MAGPIE is fully automated without any human intervention or API access to advanced commercial models such as GPT-4 [2]. Consequently, MAGPIE is cost-effective and scalable. On average, implementing MAGPIE on a cloud server² would incur costs of **\$0.12** and **\$1.1** per 1,000 data instances for MAGPIE-Air and MAGPIE-Pro, respectively.

4.5 Additional Analysis

Additional dataset analysis, including the impact of generation configurations on the quality and difficulty of the generated instructions, is in Appendix C.3. Ablation analysis on annotating models for assessing quality and difficulty is in Appendix C.4. Contamination analysis is in Appendix C.5.

²<https://lambdalabs.com/service/gpu-cloud>

Chapter 5

PERFORMANCE ANALYSIS

In this section, we evaluate the quality of MAGPIE-generated datasets by utilizing them to align base models including Llama-3 [45], Qwen1.5 [3], and Qwen2 [68].

5.1 *Experimental Setups*

Baselines for Supervised Fine-Tuning and Preference Optimization. We compare the family of instruction datasets generated by MAGPIE with eight SOTA open-source instruction datasets: **ShareGPT** [12], **WildChat** [74], **Evol Instruct** [66], **UltraChat** [19], **GenQA** [9], **OpenHermes 1** [56], **OpenHermes 2.5** [55], and **Tulu V2 Mix** [28]. ShareGPT and WildChat are representative human-written datasets containing 112K and 652K high-quality multi-round conversations between humans and GPT, respectively. Evol Instruct, UltraChat, and GenQA are representative open-source synthetic datasets. Following [44], we use the 208K sanitized version of Ultrachat provided by HuggingFace¹. OpenHermes 1, OpenHermes 2.5, and Tulu V2 Mix are crowd-sourced datasets consisting of a mix of diverse open-source instruction datasets, with 243K, 1M, and 326K conversations, respectively. We also create an instruction dataset with 100K conversations using the Self-Instruct [62] and Llama-3-8B-Instruct model, denoted as **Self-Instruct (Llama-3)**.

We compare the models aligned using MAGPIE with preference optimization baselines using direct preference optimization (DPO) [49]. Specifically, we follow [44] and use the models fine-tuned with the UltraChat dataset (for instruction tuning) and **Ultrafeedback** dataset (for preference optimization) [15].

Magpie Setups. To demonstrate the quality of MAGPIE-generated instruction datasets

¹https://huggingface.co/datasets/HuggingFaceH4/ultrachat_200k

for SFT, we select the first 300K **Magpie-Air** and **Magpie-Pro** raw datasets generated by Llama-3-8B-Instruct and Llama-3-70B-Instruct models, respectively. Apart from these raw datasets, we also applied the filters detailed in Appendix B and created two filtered datasets: **Magpie-Air-Filtered** and **Magpie-Pro-Filtered**, each contains 300K conversations. For preference optimization, we generate two additional datasets: **Magpie-Air-DPO** (generated by Llama-3-8B-Instruct) and **Magpie-Pro-DPO** (generated by Llama-3-70B-Instruct) with $k = 5$ and $T = 0.8$, each contains 100K conversations. We use RLHF1ow/ArmoRM-Llama3-8B-v0.1 [60] as the reward model.

Model Alignment Details. For supervised fine-tuning, we follow [57] and use a cosine learning rate schedule with an initial learning rate of 2×10^{-5} when fine-tuning Llama-3, Qwen1.5 and Qwen2 base models. The maximum sequence length is 8192. For DPO, we use a cosine learning rate of 5×10^{-7} . The detailed parameters can be found in Appendix D.2. We follow the official instruction templates of each model.

Evaluation Benchmarks. We evaluate the performance of the aligned models using two widely adopted instruction-following benchmarks: AlpacaEval 2 [38] and Arena-Hard [37]. AlpacaEval 2 consists of 805 representative instructions chosen from real user interactions. Arena-Hard is an enhanced version of MT-Bench [76], containing 500 challenging user queries. Both benchmarks employ a GPT evaluator to assess responses generated by the model of interest and a baseline model. Specifically, we use GPT-4-Turbo (1106) and Llama-3-8B-Instruct as baselines for AlpacaEval 2. By default, Arena-Hard uses GPT-4 (0314) as its baseline model.

Metrics. We adopt two metrics to measure the capabilities of instruction-following of fine-tuned models. The first metric is the **win rate (WR)**, which calculates the fraction of responses that are favored by the GPT evaluator. This metric is applied in both benchmarks including AlpacaEval 2 and Arena-Hard. The second metric is the **length-controlled win rate (LC)** [22], a debiased version of WR. The GPT evaluator considers the lengths of responses generated by the baseline model and model under evaluation when computing LC. By accounting for response length, LC reduces its impact on the win rate. This metric is

specifically applied to the AlpacaEval 2 benchmark [38].

More Experimental Setups. We provide more detailed descriptions of our experimental setups, including more model alignment details and benchmark decoding hyperparameters in Appendix D.

5.2 Experimental Results

Magpie datasets outperform baselines with SFT only. In Table 5.1, we compare the performance of Llama-3 models fine-tuned with instruction datasets generated by MAGPIE against those supervised fine-tuned with baseline datasets. Using the AlpacaEval 2 benchmark, we observe that both the LC and WR of our supervised fine-tuned models surpass all those models fine-tuned with baseline SFT datasets. This indicates that the datasets generated by MAGPIE are of higher quality, leading to significantly enhanced instruction-following capabilities. A similar observation is made when using the Arena-Hard evaluation benchmark. We highlight that the Llama-3 base models supervised fine-tuned with instruction datasets generated by MAGPIE outperform even those models that have undergone preference optimization (i.e., STF followed by DPO), which further emphasizes the high quality of data generated by MAGPIE.

To investigate the advantages of MAGPIE across different task categories, we compare the performance of models fine-tuned with MAGPIE-Pro with baseline datasets using WildBench benchmark [39]. This benchmark consists of 1024 tasks carefully selected from real-world human-LLM conversation logs. The results are demonstrated in Figure 5.1. We observe that MAGPIE consistently outperforms baseline datasets across categories.

Models aligned with data generated by Magpie achieve comparable or even higher performance to the official aligned model, but with fewer data. In Table 5.1, we also compare the performance of models aligned with data generated by MAGPIE against the official aligned model (Llama-3-8B-Instruct). We observe that the Llama-3-8B

²<https://huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct>

Table 5.1: This table compares the performance of models instruction-tuned on the Llama-3-8B base models using MAGPIE-generated datasets and baseline datasets. We observe that models aligned with our datasets significantly outperform those aligned with baseline datasets of the same order of magnitude in terms of data size. In addition, our fine-tuned models achieve comparable performance to the official aligned model, despite only undergoing SFT with a much smaller dataset. Numbers in **bold** indicate that MAGPIE outperforms the official Llama-3-8B-Instruct model.

Alignment Setup (Base LLM = Llama-3-8B)		#Convs	AlpacaEval 2						Arena-Hard
			GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR(%)
			LC (%)	WR (%)	SD	LC (%)	WR (%)	SD	
SFT	+Self-Instruct (Llama-3) [62]	100K	7.21	5.18	0.7	17.86	12.73	1.05	4.0
	+ShareGPT [12]	112K	9.73	7.2	0.81	27.26	18.32	1.18	6.5
	+Evol Instruct [66]	143K	8.52	6.25	0.76	20.16	14.98	1.1	5.1
	+OpenHermes 1 [56]	243K	9.94	6.27	0.73	29.19	17.92	1.16	4.4
	+Tulu V2 Mix [28]	326K	9.91	7.94	0.86	24.28	18.64	1.18	5.4
	+WildChat [74]	652K	14.62	10.58	0.92	34.85	26.57	1.32	8.7
	+OpenHermes 2.5 [55]	1M	12.89	9.74	0.91	32.68	25.01	1.30	8.2
	+GenQA [9]	6.47M	9.05	7.11	0.82	21.90	16.09	1.12	3.0
	+UltraChat [19] \blacktriangledown	208K	8.29	5.44	0.71	23.95	15.12	1.11	3.6
+ DPO	+UltraFeedback([15])	64K	18.36	17.33	1.14	44.42	42.36	1.46	14.8
SFT	Magpie-Air-300K-Raw	300K	21.99	21.65	1.21	48.63	48.06	1.42	15.8
	<u>Magpie-Air-300K-Filtered</u> \blacktriangledown	300K	22.66	23.99	1.24	49.27	50.8	1.44	14.9
+ DPO	+ Magpie-Air-DPO	100K	45.48	50.43	1.48	75.06	79.64	1.18	35.9
SFT	Magpie-Pro-300K-Raw	300K	21.65	22.19	1.2	49.65	50.84	1.42	15.9
	<u>Magpie-Pro-300K-Filtered</u> \blacktriangledown	300K	25.08	29.47	1.35	52.12	53.43	1.44	18.9
+ DPO	+ Magpie-Pro-DPO	100K	50.10	53.53	1.45	78.52	80.82	1.17	35.7
Llama-3-8B-Instruct (SFT+DPO)		>10M ²	22.92	22.57	1.26	50	50	-	20.6

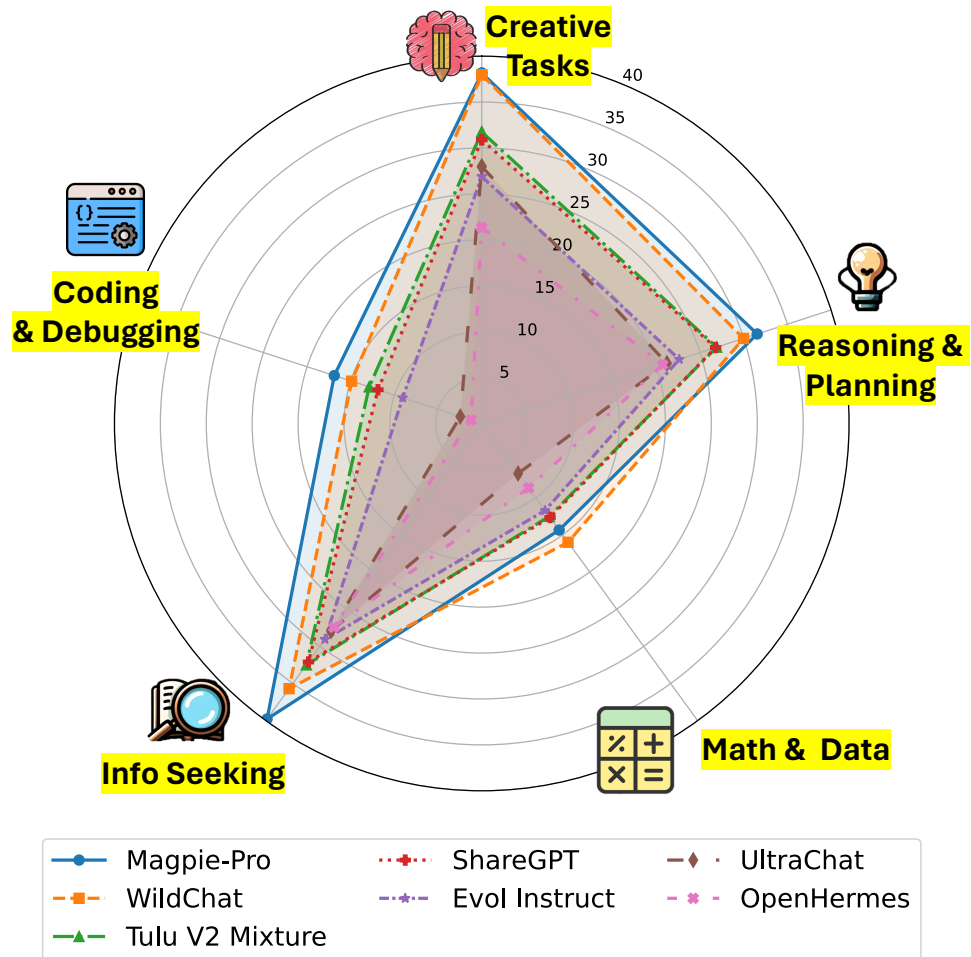


Figure 5.1: This figure shows the performance breakdown by category of MAGPIE-Pro and baselines on WildBench.

base model supervised fine-tuned with data from MAGPIE outperforms Llama-3-8B-instruct using the AlpacaEval 2 benchmark. For example, when Llama-3-8B-Instruct is chosen as the baseline model of AlpacaEval 2, we observe that LC of Llama-3-8B base models fine-tuned with instruction data from MAGPIE exceeds 50%, indicating a preference for our SFT models over the official aligned model. In addition, when DPO is applied, our aligned model demonstrates remarkable performance gains. Specifically, it outperforms the official Llama-3-8B-Instruct model on both the AlpacaEval 2 and Arena-Hard benchmarks. Most notably, our model even surpasses GPT-4-Turbo(1106) on AlpacaEval 2. Finally, we highlight that our alignment process uses no more than 400K data, whereas the official aligned models are aligned with more than 10M data samples. This demonstrates the high quality of the data generated by MAGPIE.

Magpie can enhance the performance of other backbone models. Table 5.2 illustrates the efficacy of MAGPIE when applied to generate instruction dataset and fine-tune other base models, i.e., Qwen2-1.5B, Qwen1.5-4B, and Qwen1.5-7B. The results demonstrate that our fine-tuned models achieve better performance than the official aligned models, which have undergone both supervised fine-tuning and preference tuning. These results underscore the effectiveness of MAGPIE and the quality of its generated instructions.

Performance of Magpie on More Benchmarks. We report the performance of models supervised fine-tuned using MAGPIE-Air and MAGPIE-Pro, evaluated across a range of tasks featured on the Huggingface Open LLM Leaderboard [5] in Table 5.3. The tasks includes MMLU [26], ARC Challenge [13], HellaSwag [73], TruthfulQA [41], WinoGrande [34], and GSM8K [14]. We also perform experiments on MMLU-Redux [25] with zero-shot prompting. Our experimental results demonstrate that models fine-tuned with MAGPIE-Air and MAGPIE-Pro achieve comparable performance to the official instruct model and other baselines.

We note that the performance of MAGPIE may degrade on reasoning tasks, which is attributed to the small proportion of reasoning instructions in MAGPIE-Air and MAGPIE-Pro datasets. In response, we provide a supplementary "booster" dataset containing 150K

Table 5.2: This table compares the performance of models instruction-tuned on the Qwen base models using the MAGPIE-Pro-300K-Filtered dataset and the official instruction-tuned models. The Qwen base model enhanced with MAGPIE outperforms the official instruction-tuned model.

Alignment Setup		AlpacaEval 2					
		GPT-4-Turbo (1106)			Official Aligned Model as Ref.		
		LC (%)	WR (%)	SD	LC (%)	WR (%)	SD
Qwen2-1.5B	Qwen2-1.5B-Instruct	3.91	3.00	0.54	50	50	-
	Base Model + MAGPIE	3.48	5.32	0.67	56.66	66.27	1.50
Qwen1.5-4B	Qwen1.5-4B-Chat	5.89	4.74	0.67	50	50	-
	Base Model + MAGPIE	9.1	10.96	0.93	68.09	72.42	1.42
Qwen1.5-7B	Qwen1.5-7B-Chat	14.75	11.77	0.97	50	50	-
	Base Model + MAGPIE	15.10	18.51	1.14	46.28	58.53	1.44

Table 5.3: This table compares the performance of models supervised-fine-tuned on MAGPIE-Air, MAGPIE-Pro, and MAGPIE-Pro-Mix against baselines and official instruct model across various downstream benchmarks. All models are supervised-fine-tuned on the Llama-8B base models.

Alignment Setup	MMLU (5)	ARC (25)	HellaSwag (10)	TruthfulQA (0)	WinoGrande (5)	GSM8K (5)	MMLU-Redux (0)	Average
ShareGPT	66.03	58.45	81.50	52.34	74.03	48.67	50.68	61.67
Evol Instruct	65.62	60.75	82.70	52.87	76.16	42.91	52.73	61.96
GenQA	63.45	58.53	79.65	48.85	74.03	43.14	51.87	59.93
OpenHermes 1	65.42	62.29	82.15	50.85	75.61	47.16	46.07	61.36
OpenHermes 2.5	65.70	61.86	82.53	51.35	76.09	67.02	46.07	66.24
Tulu V2 Mix	66.34	59.22	82.80	47.99	76.16	58.07	46.97	62.51
WildChat	65.95	59.22	81.39	53.18	75.30	48.75	52.59	62.34
UltraChat	65.23	62.12	81.68	52.76	75.53	50.57	50.75	62.66
MAGPIE-Air-300K-Filtered	64.45	61.01	79.90	53.48	72.38	52.24	52.34	62.25
MAGPIE-Pro-300K-Filtered	64.25	60.41	80.52	52.46	73.32	47.92	52.16	61.58
Magpie-Pro-Mix-Filtered	65.65	59.64	80.72	50.81	73.24	63.08	56.34	64.21
Llama-3-8B-Instruct	67.82	61.52	78.67	52.47	72.14	71.72	58.60	66.13

math, code, and reasoning instructions using the MAGPIE extension mentioned in Section 3.2. We combine this booster dataset with MAGPIE-Pro-300K-Filtered and create MAGPIE-Pro-Mix-Filtered. Experimental results presented in Table 5.3 demonstrate that the model supervised fine-tuned using the mixed dataset effectively addresses the initial weakness in reasoning tasks. Notably, this new model ranks among the top-3 of all model checkpoints, performing only slightly weaker than OpenHermes 2.5 (1M conversations) and Llama-3-8B-Instruct (\approx 10M conversations). This significant improvement showcases the flexibility and adaptability of the MAGPIE framework in generating task-specific instruction data.

Additional Experimental Results. We defer additional experimental results and analysis of multi-turn datasets, i.e., MAGPIE-Air-MT and MAGPIE-Pro-MT, to Appendix E.1. We conduct a detailed comparison between MAGPIE and Self-Instruct in Appendix E.2. We demonstrate the performance of domain-specific and multi-lingual MAGPIE datasets in Appendix E.3. In addition, ablations on data quantity, quality, filter designs, and response generator are deferred in Appendices E.4, E.5, and E.6. MAGPIE model’s performance on

trustworthiness and instruction-following benchmarks is reported in Appendices E.7 and E.8. The ablation analysis on the impact of reward models on DPO data performance is presented in Appendix E.9.

Chapter 6

CONCLUSION

In this paper, we developed a scalable method, MAGPIE, to synthesize instruction data for fine-tuning large language models. MAGPIE leveraged the predefined instruction templates of open-weight LLMs and crafted a prompt specifying only the role of instruction provider. Given the crafted prompt, the LLM then generated detailed instructions due to their autoregressive nature. MAGPIE then sent the generated instructions to the LLM to generate corresponding responses. These pairs of instructions and responses constituted the instruction dataset. We used Llama-3-8B-instruct to label the instruction dataset and developed a filtering technique to select effective data instances for instruction tuning. We fine-tuned the Llama-3-8B base model using the selected data, and demonstrated that the fine-tuned model outperformed those fine-tuned using all baselines. Moreover, our fine-tuned models outperformed the official aligned model, Llama-3-8B-Instruct, which has been instruction-tuned and preference-optimized using more than 10M data instances. This highlighted the quality of the instruction data synthesized by MAGPIE.

Chapter 7

ACKNOWLEDGMENT

This research is partially supported by the Air Force Office of Scientific Research (AFOSR) under grant FA9550-23-1-0208, the National Science Foundation (NSF) AI Institute for Agent-based Cyber Threat Intelligence and Operation (ACTION) under grant IIS 2229876, and the Office of Naval Research (ONR) under grant N0014-23-1-2386.

This work is supported in part by funds provided by the National Science Foundation, Department of Homeland Security, and IBM. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the NSF or its federal agency and industry partners.

BIBLIOGRAPHY

- [1] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [4] Yang Bai, Ge Pei, Jindong Gu, Yong Yang, and Xingjun Ma. Special characters attack: Toward scalable training data extraction from large language models. *arXiv preprint arXiv:2405.05990*, 2024.
- [5] Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. Open llm

- leaderboard. https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard, 2023.
- [6] Stella Biderman, Usvsn Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. Emergent and predictable memorization in large language models. *Advances in Neural Information Processing Systems*, 36, 2023.
- [7] Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. What does it mean for a language model to preserve privacy? In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 2280–2292, 2022.
- [8] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [9] Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John Kirchenbauer, Tianyi Zhou, and Tom Goldstein. Genqa: Generating millions of instructions from a handful of prompts. *arXiv preprint arXiv:2406.10323*, 2024.
- [10] Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*, 2023.
- [11] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [12] 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. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [13] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [14] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [15] Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*, 2023.
- [16] Databricks. Databricks dolly-15k, 2023.
- [17] Daniel Deutsch, Rotem Dror, and Dan Roth. On the limitations of reference-free evaluations of generated text. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10960–10977, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [19] Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*, 2023.

- [20] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024.
- [21] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [22] Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- [23] Saumya Gandhi, Ritu Gala, Vijay Viswanathan, Tongshuang Wu, and Graham Neubig. Better synthetic data by retrieving and transforming existing datasets. *arXiv preprint arXiv:2404.14361*, 2024.
- [24] 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. A framework for few-shot language model evaluation, 07 2024.
- [25] Aryo Pradipta Gema, Joshua Ong Jun Leang, Giwon Hong, Alessio Devoto, Alberto Carlo Maria Mancino, Rohit Saxena, Xuanli He, Yu Zhao, Xiaotang Du, Mohammad Reza Ghasemi Madani, et al. Are we done with mmlu? *arXiv preprint arXiv:2406.04127*, 2024.
- [26] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

- [27] Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, Hanchi Sun, Zhengliang Liu, Yixin Liu, Yijue Wang, Zhikun Zhang, Bertie Vidgen, Bhavya Kailkhura, Caiming Xiong, Chaowei Xiao, Chunyuan Li, Eric P. Xing, Furong Huang, Hao Liu, Heng Ji, Hongyi Wang, Huan Zhang, Huaxiu Yao, Manolis Kellis, Marinka Zitnik, Meng Jiang, Mohit Bansal, James Zou, Jian Pei, Jian Liu, Jianfeng Gao, Jiawei Han, Jieyu Zhao, Jiliang Tang, Jindong Wang, Joaquin Vanschoren, John Mitchell, Kai Shu, Kaidi Xu, Kai-Wei Chang, Lifang He, Lifu Huang, Michael Backes, Neil Zhenqiang Gong, Philip S. Yu, Pin-Yu Chen, Quanquan Gu, Ran Xu, Rex Ying, Shuiwang Ji, Suman Jana, Tianlong Chen, Tianming Liu, Tianyi Zhou, William Yang Wang, Xiang Li, Xiangliang Zhang, Xiao Wang, Xing Xie, Xun Chen, Xuyu Wang, Yan Liu, Yanfang Ye, Yinzhi Cao, Yong Chen, and Yue Zhao. Trustllm: Trustworthiness in large language models. In *Forty-first International Conference on Machine Learning*, 2024.
- [28] Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*, 2023.
- [29] Aly M Kassem, Omar Mahmoud, Niloofar Miresghallah, Hyunwoo Kim, Yulia Tsvetkov, Yejin Choi, Sherif Saad, and Santu Rana. Alpaca against vicuna: Using llms to uncover memorization of llms. *arXiv preprint arXiv:2403.04801*, 2024.
- [30] Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations - democratizing large language model alignment. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 47669–47681. Curran Associates, Inc., 2023.

- [31] Kalpesh Krishna, Gaurav Singh Tomar, Ankur P. Parikh, Nicolas Papernot, and Mohit Iyyer. Thieves on sesame street! model extraction of bert-based apis. In *International Conference on Learning Representations*, 2020.
- [32] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*, 2023.
- [33] Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*, 2024.
- [34] Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In *Thirteenth international conference on the principles of knowledge representation and reasoning*, 2012.
- [35] Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for” mind” exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.
- [36] Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, et al. Synthetic data (almost) from scratch: Generalized instruction tuning for language models. *arXiv preprint arXiv:2402.13064*, 2024.
- [37] Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Banghua Zhu, Joseph E. Gonzalez, and Ion Stoica. From live data to high-quality benchmarks: The arena-hard pipeline, April 2024.

- [38] Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 2023.
- [39] Bill Yuchen Lin, Khyathi Chandu, Faeze Brahman, Yuntian Deng, Abhilasha Ravichander, Valentina Pyatkin, Ronan Le Bras, and Yejin Choi. Wildbench: Benchmarking language models with challenging tasks from real users in the wild, 2024.
- [40] Bill Yuchen Lin, Abhilasha Ravichander, Ximing Lu, Nouha Dziri, Melanie Sclar, Khyathi Chandu, Chandra Bhagavatula, and Yejin Choi. The unlocking spell on base llms: Rethinking alignment via in-context learning. *arXiv preprint arXiv:2312.01552*, 2023.
- [41] Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- [42] Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinneng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. Best practices and lessons learned on synthetic data for language models. *arXiv preprint arXiv:2404.07503*, 2024.
- [43] Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. In *The Twelfth International Conference on Learning Representations*, 2024.
- [44] Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- [45] Meta. Llama 3. <https://ai.meta.com/blog/meta-llama-3/>, 2024.
- [46] Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.

- [47] OpenAI. Tiktoken. <https://github.com/openai/tiktoken>, 2024.
- [48] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [49] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [50] Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022.
- [51] Zhiqing Sun, Yikang Shen, Qinlong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. Principle-driven self-alignment of language models from scratch with minimal human supervision. *Advances in Neural Information Processing Systems*, 36, 2023.
- [52] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [53] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin,

- Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [54] Llama Team. Meta llama guard 2. https://github.com/meta-llama/PurpleLlama/blob/main/Llama-Guard2/MODEL_CARD.md, 2024.
- [55] Teknium. Openhermes 2.5: An open dataset of synthetic data for generalist llm assistants, 2023.
- [56] Teknium. Openhermes dataset, 2023.
- [57] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [58] Hoang Tran, Chris Glaze, and Braden Hancock. Iterative DPO alignment. Technical report, Snorkel AI, 2023.
- [59] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [60] Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. In *EMNLP*, 2024.
- [61] Jeffrey G Wang, Jason Wang, Marvin Li, and Seth Neel. Pandora’s white-box: Increased training data leakage in open llms. *arXiv preprint arXiv:2402.17012*, 2024.
- [62] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. In *Proceedings of the 61st Annual Meeting of the Association*

- for Computational Linguistics (Volume 1: Long Papers)*, pages 13484–13508, Toronto, Canada, 2023. Association for Computational Linguistics.
- [63] 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, Krима 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. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.
- [64] Zifeng Wang, Chun-Liang Li, Vincent Perot, Long T Le, Jin Miao, Zizhao Zhang, Chen-Yu Lee, and Tomas Pfister. Codeclm: Aligning language models with tailored synthetic data. *arXiv preprint arXiv:2404.05875*, 2024.
- [65] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint, 2024.
- [66] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- [67] Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in*

- Natural Language Processing*, pages 6268–6278, Singapore, December 2023. Association for Computational Linguistics.
- [68] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- [69] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [70] Da Yin, Xiao Liu, Fan Yin, Ming Zhong, Hritik Bansal, Jiawei Han, and Kai-Wei Chang. Dynosaur: A dynamic growth paradigm for instruction-tuning data curation. *arXiv preprint arXiv:2305.14327*, 2023.
- [71] Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. Bag of tricks for training data extraction from language models. In *International Conference on Machine Learning*, pages 40306–40320. PMLR, 2023.
- [72] Zhaojian Yu, Xin Zhang, Ning Shang, Yangyu Huang, Can Xu, Yishujie Zhao, Wenxiang Hu, and Qiufeng Yin. Wavecoder: Widespread and versatile enhancement for code large language models by instruction tuning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5140–5153, 2024.

- [73] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- [74] Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. Wildchat: 1m chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*, 2024.
- [75] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. LMSYS-chat-1m: A large-scale real-world LLM conversation dataset. In *The Twelfth International Conference on Learning Representations*, 2024.
- [76] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [77] Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36, 2023.
- [78] Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*, 2023.
- [79] Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*, 2024.

Appendix A

STATISTICS OF INSTRUCTION DATASETS GENERATED BY Magpie COMPARED TO OTHER INSTRUCTION DATASETS.

MAGPIE can be readily deployed to state-of-the-art open-weight model families, including but not limited to Llama-3 [45], Llama-3.1/3.3 [21], Qwen2 [68], Qwen2.5 [69], Gemma-2 [53], and Phi-3 [1] models.

In what follows, we compare datasets generated by MAGPIE with the above model families compared to other state-of-the-art instruction datasets. The MAGPIE dataset family encompasses over 13.4 million diverse and high-quality instructions and corresponding responses generated from state-of-the-art open-source models. This corpus represents the largest alignment dataset for LLMs that does not rely on human-written questions or employ complex multi-stage pipelines.

Table A.1: Statistics of MAGPIE family compared to other instruction datasets. Tokens are counted using the `tiktoken` library [47]. Links to the MAGPIE datasets are provided in the text.

Instruction Source	Dataset Name	#Convs	#Turns	Human Effort	Response Generator	#Tokens / Turn	#Total Tokens
Synthetic	Alpaca [52]	52K	1	Low	text-davinci-003	67.38 \pm 54.88	3.5M
	Evol Instruct [66]	143K	1	Low	ChatGPT	473.33 \pm 330.13	68M
	UltraChat [19]	208K	3.16	Low	GhatGPT	376.58 \pm 177.81	238M
Human	Dolly [16]	15K	1	High	ChatGPT	94.61 \pm 135.84	1.42M
	ShareGPT [76]	112K	4.79	High	ChatGPT	465.38 \pm 368.37	201M
	WildChat [74]	652K	2.52	High	GPT-3.5 & GPT-4	727.09 \pm 818.84	852M
	LMSYS-Chat-1M [75]	1M	2.01	High	Mix	260.37 \pm 346.97	496M
Mixture	Deita [43]	9.5K	22.02	-	Mix	372.78 \pm 182.97	74M
	OpenHermes [56]	243K	1	-	Mix	297.86 \pm 258.45	72M
	Tulu V2 Mixture [28]	326K	2.31	-	Mix	411.94 \pm 447.48	285M
Magpie	MAGPIE-Llama-3-Air	3M	1	No	Llama-3-8B-Instruct	426.39 \pm 217.39	1.28B
	MAGPIE-Llama-3-Air-MT	300K	2	No	Llama-3-8B-Instruct	610.80 \pm 90.61	366M
	MAGPIE-Llama-3-Pro	1M	1	No	Llama-3-70B-Instruct	478.00 \pm 211.09	477M
	MAGPIE-Llama-3-Pro-MT	300K	2	No	Llama-3-70B-Instruct	554.53 \pm 133.64	333M
	MAGPIE-Llama-3.1-Pro	1M	1	No	Llama-3.1-70B-Instruct	482.35 \pm 378.45	482M
	MAGPIE-Llama-3.1-Pro-MT	300K	2	No	Llama-3.1-70B-Instruct	552.53 \pm 325.49	331M
	MAGPIE-Llama-3.3-Pro	1M	1	No	Llama-3.3-70B-Instruct	568.59 \pm 391.54	569M
	MAGPIE-Qwen2-Air	3M	1	No	Qwen2-7B-Instruct	577.87 \pm 416.10	1.73B
	MAGPIE-Qwen2-Pro	1M	1	No	Qwen2-72B-Instruct	424.87 \pm 339.71	424M
	MAGPIE-Qwen2.5-Pro	1M	1	No	Qwen2.5-72B-Instruct	693.31 \pm 271.45	693M
	MAGPIE-Gemma2-Pro	534K	1	No	Gemma-2-27b-it	483.90 \pm 237.80	259M
	MAGPIE-Phi3-Pro	1M	1	No	Phi-3-Medium-Instruct	391.38 \pm 414.32	391M

Appendix B

FILTER SETUPS

In this section, we explore potential filter configurations for selecting high-quality instructional data for fine-tuning purposes. We provide the following metrics to enable users to customize their filtered MAGPIE dataset:

1. **Input Length:** The total number of characters in the instructions.
2. **Output Length:** The total number of characters in the responses.
3. **Task Category:** The specific category of the instructions. See Appendix C.1 for details.
4. **Input Quality:** The clarity, specificity, and coherence of the instructions, rated as ‘very poor’, ‘poor’, ‘average’, ‘good’, and ‘excellent’.
5. **Input Difficulty:** The level of knowledge required to address the task described in the instruction, rated as ‘very easy’, ‘easy’, ‘medium’, ‘hard’, or ‘very hard’.
6. **Minimum Neighbor Distance:** The embedding distance to the nearest neighbor. Can be used for filtering out repetitive or similar instances.
7. **Reward:** Denoted as r^* . See Section 4 for details. This metric can be used to filter out low-quality responses, such as repetitions or refusals.
8. **Reward Difference:** Denoted as $r^* - r_{base}$. See Section 4 for details.

We provide several off-the-shelf configurations, as demonstrated in Table B.1. We defer the detailed performance analysis of each filter configuration for MAGPIE-Pro to Appendix E.5.

Table B.1: Different filter configurations we provide. We note that the Output Length filter is applied last. Specifically, this filter selects the k instances of the longest responses. In our experiments, we empirically set $\tau_1 = -12$, and $\tau_2 = 0$.

Source Dataset	Filter Name	#Convs	Input Length	Output Length	Task Category	Input Quality	Input Difficulty	Min Neighbor Distance	Reward	Reward Difference
MAGPIE-Air	Filter	300K	-	Longest	-	\geq good	\geq medium	> 0	-	$> \tau_2$
	Filter	300K	-	Longest	-	\geq average	-	> 0	$> \tau_1$	-
	Filter2	300K	-	Longest	-	\geq good	\geq easy	> 0	$> \tau_1$	-
MAGPIE-Pro	Filter3	300K	-	Longest	-	-	-	> 0	$> \tau_1$	-
	Filter4	300K	-	Longest	-	\geq good	\geq easy	> 0	-	$> \tau_2$
	Filter5	338K	-	-	-	\geq good	\geq easy	> 0	$> \tau_1$	-
	Filter6	200K	-	Longest	-	-	50% easy + 50% $>$ easy	> 0	$> \tau_1$	-

Appendix C

MORE DATASET ANALYSIS

This section provides additional dataset analysis, complementing the discussions in Section 4. Statistics including lengths of instructions and responses are illustrated in Figure C.1.

C.1 Additional Analysis on Dataset Coverage and Attributes.

Dataset Coverage Measured by T-SNE and UMAP. Figure C.2 presents the t-SNE and UMAP plots of MAGPIE, Alpaca, Evol Instruct, and UltraChat. Each t-SNE and UMAP plot is generated by randomly sampling 10,000 instructions from the associated dataset. We observe that the t-SNE and UMAP plot of MAGPIE encompasses the area covered by the plots of Alpaca, Evol Instruct, and UltraChat. This suggests that MAGPIE datasets provides a broader or more diverse range of topics, highlighting its extensive coverage across varied themes and subjects.

Task Categories of Magpie-Pro and Magpie-Air. Figure C.4 illustrates the task category distributions for MAGPIE-Pro and MAGPIE-Air, as labeled by Llama-3-Instruct. We observe that the task category distributions of these two datasets are largely similar, however, MAGPIE-Pro exhibits a higher percentage of creative writing tasks.

Topic Diversity of Magpie-Pro and Magpie-Air. To validate the diversity of generated instructions, we conducted additional analysis using Topic Diversity metric from UltraChat [19]. Our results are summarized in Table C.1. The results demonstrate that our generated instructions are indeed more diverse in topic compared with other baselines.

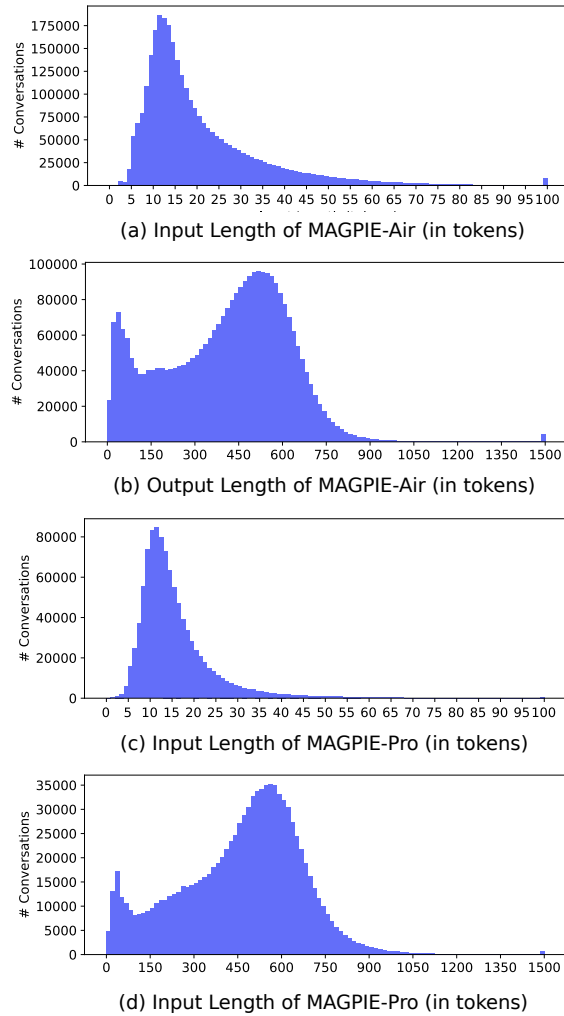


Figure C.1: Lengths of instructions and responses in MAGPIE-Air/Pro.

Table C.1: Comparison of Topic Diversity across Different Synthetic Datasets.

Dataset	Alpaca	Evol Instruct	UltraChat	Magpie-Air	Magpie-Pro
Topic Diversity (\downarrow)	0.13	0.09	0.10	0.05	0.06

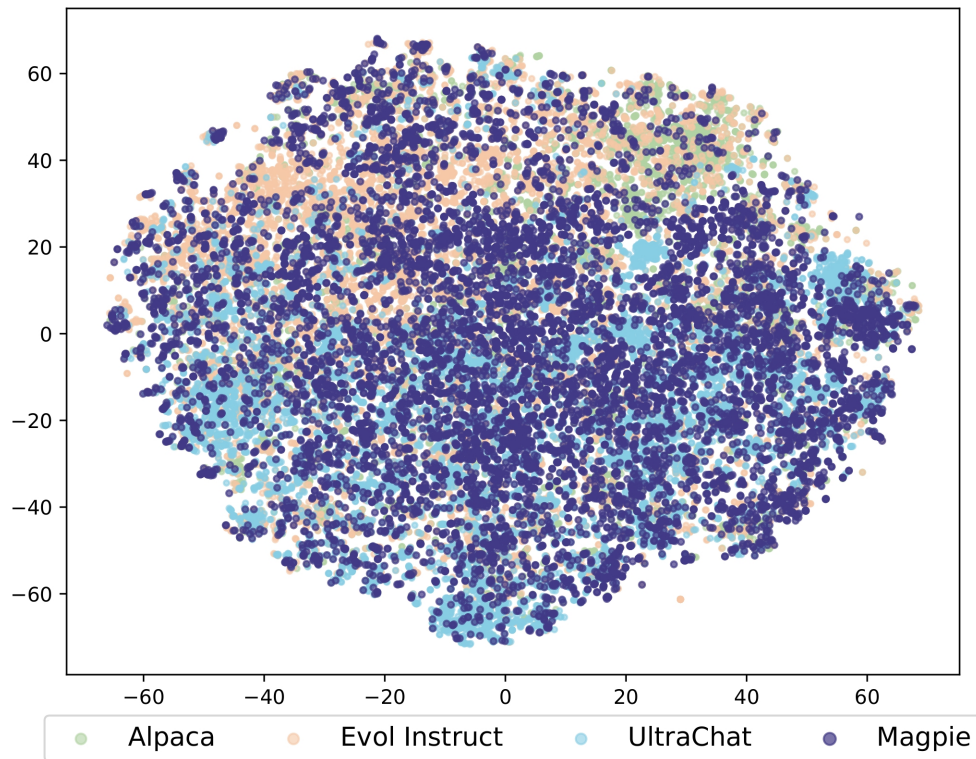


Figure C.2: This figure compares the t-SNE plot of MAGPIE-Pro with those of Alpaca, Evol Instruct, and UltraChat, each of which is sampled with 10,000 instructions. The t-SNE plot of MAGPIE-Pro encompasses the area covered by the other plots, demonstrating the comprehensive coverage of MAGPIE-Pro.

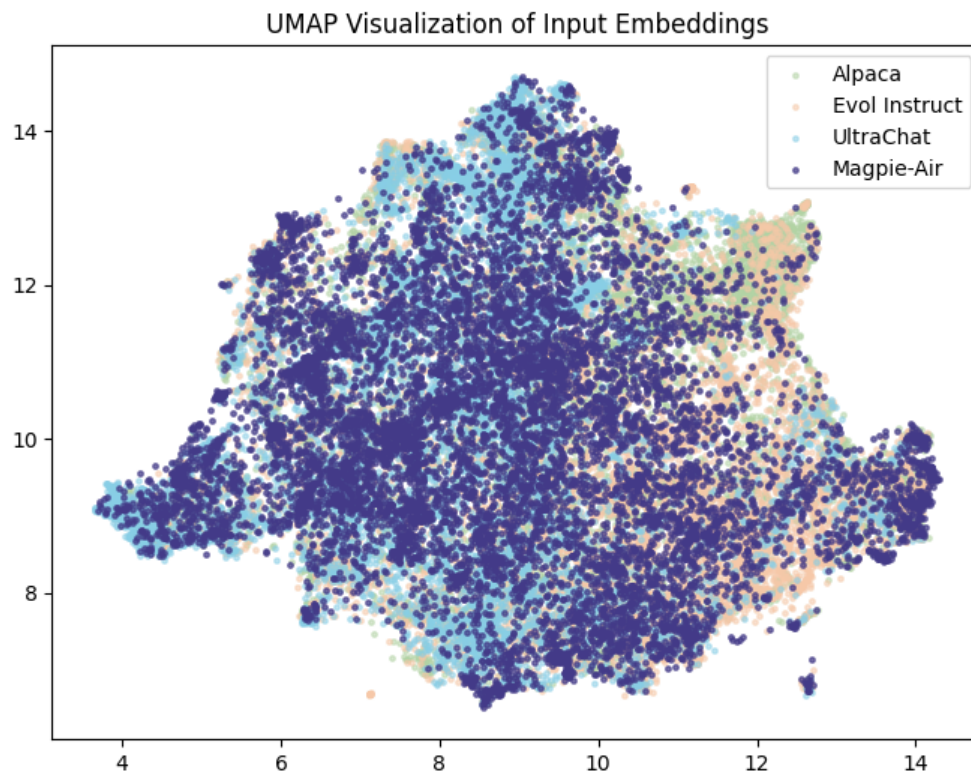


Figure C.3: This figure compares the UMAP plot of MAGPIE-Air with those of Alpaca, Evol Instruct, and UltraChat, each of which is sampled with 10,000 instructions. The UMAP plot of MAGPIE-Air encompasses the area covered by the other plots, demonstrating the comprehensive coverage of MAGPIE-Air.

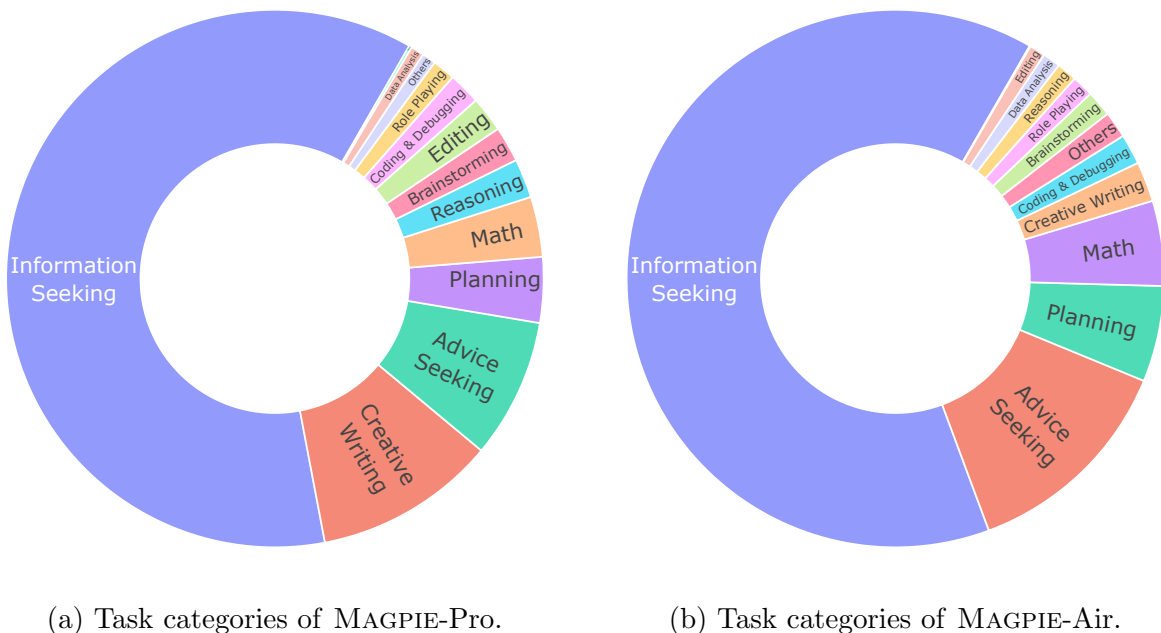


Figure C.4: This figure visualizes the task category of MAGPIE-Pro and MAGPIE-Air by topic tags.

Visualization of Root Verbs and Their Direct Noun Objects. Figure C.5 visualizes the top common root verbs and their direct noun objects of MAGPIE-Air dataset. This indicates the diverse topic coverage of MAGPIE-Air.

C.2 Additional Safety Analysis

Table C.2 illustrates the percentage of different unsafe categories of MAGPIE-Air and MAGPIE-Pro, as labeled by Llama-Guard-2 [54]. We have two key observations. First, the proportion of data containing potentially harmful queries is minimal, with less than 1% for both datasets. Second, the majority of unsafe responses fall into the category of specialized advice, which includes responses that may offer specialized financial, medical, or legal advice, or suggest that dangerous activities or objects are safe.

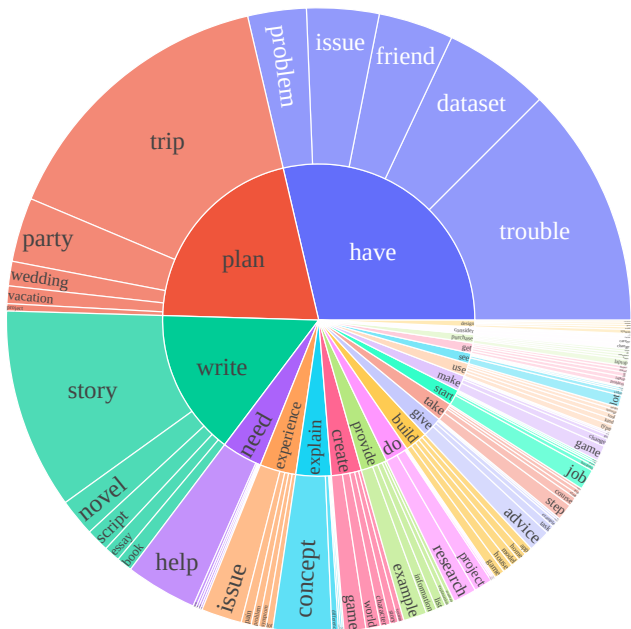


Figure C.5: This figure demonstrates the top 20 most common root verbs (shown in the inner circle) and their top 5 direct noun objects (shown in the outer circle) within the MAGPIE-Air dataset. This indicates that MAGPIE encompasses a broad range of topics.

C.3 Ablation Analysis on Generation Configurations

Ablation Analysis on Decoding Parameters. We conduct an ablation analysis on the decoding parameters used in generating instruction with MAGPIE. Specifically, we use three different temperatures (i.e., 1, 1.1, and 1.2) and top-p values (i.e., 1, 0.995, and 0.99) during Step 1 of MAGPIE. We use three metrics, **Average Quality Score**, **Average Difficulty Score** and **Average Minimum Neighbor Distance** to characterize the quality, difficulty, and diversity of instructions using different decoding parameters. The Average Quality Score is calculated by averaging the ratings of all data within a specific temperature-top-p pair, on a scale from 1 (‘very poor’) to 5 (‘excellent’). Similarly, the Average Difficulty Score is rated on a scale from 1 (‘very easy’) to 5 (‘very hard’). The Average Minimum Neighbor Distance is calculated by averaging the minimum neighbor distances, as defined in Section 4, for all data generated using the same decoding parameters.

Table C.2: This table shows the percentage of different unsafe categories of MAGPIE-Air and MAGPIE-Pro tagged by Llama-Guard-2 [54] model.

Dataset	Safe	Violent Crimes	Non-Violent Crimes	Sex-Related Crimes	Child Sexual Exploitation	Specialized Advice	Privacy	Intellectual Property	Indiscriminate Weapons	Hate	Suicide & Self-Harm	Sexual Content	Others
MAGPIE-Air	99.128%	0.001%	0.073%	0.003%	0.000%	0.636%	0.022%	0.026%	0.038%	0.001%	0.002%	0.009%	0.062%
MAGPIE-Pro	99.347%	0.001%	0.049%	0.002%	0.000%	0.446%	0.015%	0.074%	0.014%	0.001%	0.004%	0.011%	0.036%

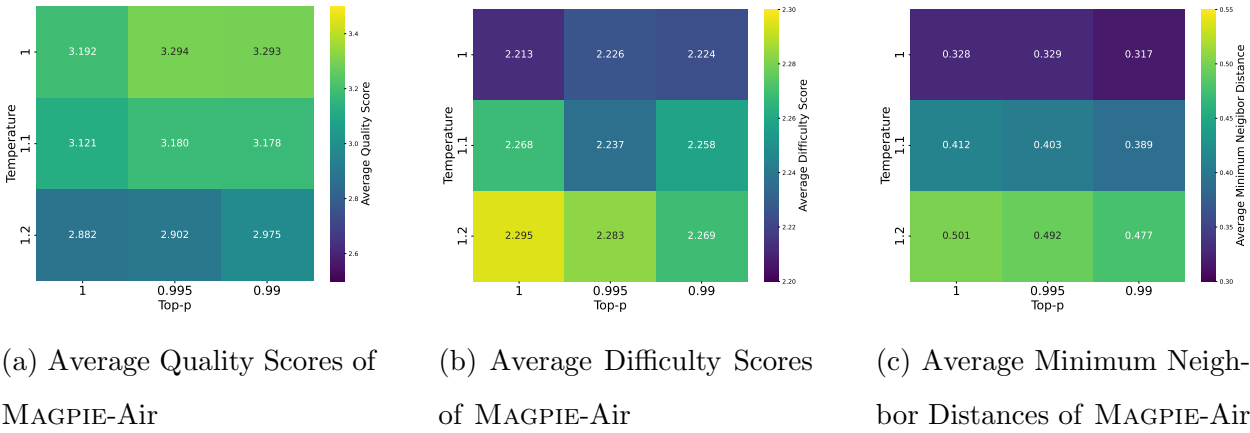


Figure C.6: This figure illustrates the impact of varying decoding parameters on the quality, difficulty, and diversity of generated instructions. We observe that while higher temperature and top-p values may decrease the overall quality, they tend to increase both the difficulty and diversity of the instructions.

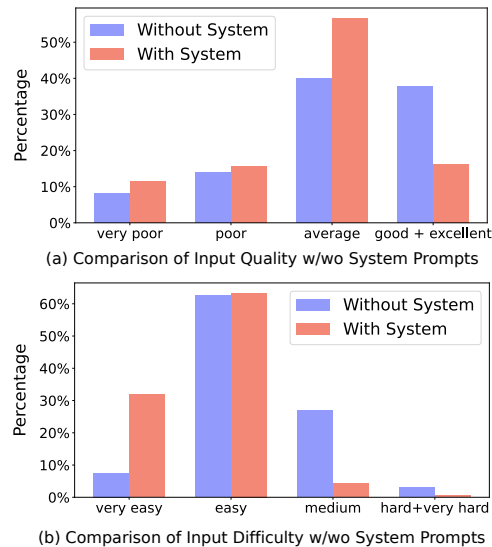


Figure C.7: This figure compares the input quality and difficulty with and without system prompts.

The findings are summarized in Figure C.6. We observe that higher temperature and top-p values may slightly decrease the overall quality of instructions, while simultaneously increasing the difficulty and remarkably enhancing the diversity of the instructions generated. The selection of these hyper-parameters should be tailored to the user’s specific requirements, balancing the trade-offs between quality, difficulty, and diversity.

Ablation Analysis on the System Prompt. Figure C.7 compares the use of system prompt compared with not using it in Step 1 of MAGPIE. Since the Llama-3 model does not have an official system prompt, we use the default system prompt from Vicuna [12]: A chat between a curious user and an artificial intelligence assistant. The assistant gives helpful, detailed, and polite answers to the user’s questions. We observe that using a system prompt generally results in a decrease in the overall quality of instructions, and the instructions are easier. Consequently, we recommend not appending system prompts in default settings.

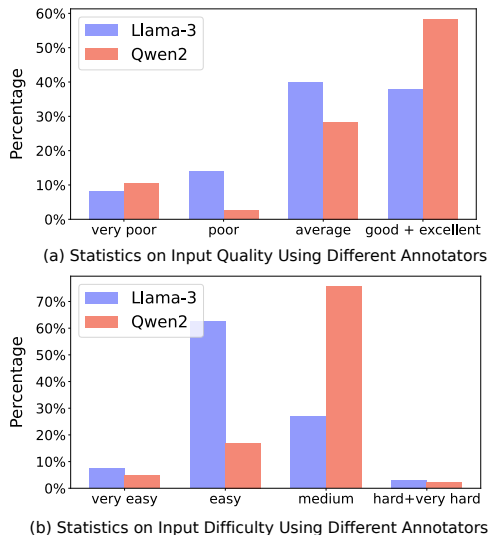


Figure C.8: This figure compares the impact of different annotators on evaluating the instruction quality and difficulty.

C.4 Impact of Annotating Models

We note that LLMs may occasionally favor its own response [17]. In what follows, we conduct experiments to evaluate the impact of annotating models when labeling quality and difficulty of the MAGPIE-Air dataset. We used the Qwen-2-7B-Instruct model (outside the Llama-3 family) to annotate the quality and difficulty of our MAGPIE-Air dataset. The statistics are summarized in Figure C.8.

Our findings show that even when evaluated by Qwen-2-7B-Instruct, the MAGPIE-Air dataset maintains high quality and difficulty, which is even higher than those originally annotated by Llama-3-8B-Instruct. This suggests that our dataset’s quality is robust across different annotators.

C.5 Contamination Analysis

We conduct a contamination analysis of MAGPIE-generated instructions using the Alpaca Eval 2 and Arena Hard benchmarks. Following [72], we use an embedding model to convert both MAGPIE instructions and benchmark questions into embeddings, calculating cosine similarity scores to assess potential overlap.

Our analysis reveals that Arena Hard shows no evidence of data contamination, while some similarities were found between Alpaca Eval 2 and MAGPIE datasets generated by Llama-3, with a few entries showing cosine similarity greater than 0.9. Notably, the affected questions account for at most 10 out of 805 benchmark questions. This is expected, as Alpaca Eval data is constructed from a mixture of existing instruction datasets, many of which contain common questions such as “What’s the capital of Australia?”, “Can you code?”, and “What’s your name?” We note that even with this small degree of overlap, MAGPIE consistently outperformed the baselines, demonstrating its robustness and showing no significant impact on benchmark evaluations.

Appendix D

DETAILED EXPERIMENTAL SETUPS

D.1 Experimental Setups for Generating Magpie-Air and Magpie-Pro

As detailed in Appendix C.3, varying decoding parameters in Step 1 can significantly influence the quality, difficulty, and diversity of the generated instructions. To optimize the trade-offs among these attributes, we employ diverse decoding parameters for the generation of MAGPIE-Air and MAGPIE-Pro. Table D.1 presents the configurations of MAGPIE-Air and MAGPIE-Pro, showcasing how diverse decoding parameters shape each dataset.

We employ greedy decoding to generate responses in Step 2 for MAGPIE-Air and MAGPIE-Pro. The intuition is that the word with the highest probability is more likely to originate from the model’s training dataset.

D.2 Experimental Setups for Instruction Tuning and Preference Tuning

Supervised Fine-Tuning Hyper-parameters. Table D.2 demonstrates the detailed supervised fine-tuning hyper-parameters. These experiments were conducted using Axolotl¹.

Preference Tuning Hyper-parameters. Table D.3 demonstrates the detailed DPO hyper-parameters for aligning Llama-3-8B using MAGPIE-Air-DPO and MAGPIE-Pro-DPO. These experiments were conducted using Alignment Handbook².

Decoding parameters for evaluation benchmarks. For Arena-Hard [37] and Wild-Bench [39], we follow its default setting and use greedy decoding for all settings. For AlpacaEval 2 [38] which allows the model provider to specify decoding parameters, we also employ greedy decoding in all experiments with a slightly increased repetition penalty ($RP = 1.2$)

¹<https://github.com/OpenAccess-AI-Collective/axolotl>

²<https://github.com/huggingface/alignment-handbook>

Table D.1: This table demonstrates the configurations of generating instructions of MAGPIE-Air and MAGPIE-Pro datasets with varying decoding parameters.

Dataset	Decoding Parameters			Total #Convs
	Temperature	Top-p	#Convs	
MAGPIE-Air	1.0	1.00	300K	3M
	1.0	0.995	300K	
	1.0	0.990	300K	
	1.1	1.00	300K	
	1.1	0.995	300K	
	1.1	0.990	300K	
	1.2	1.00	300K	
	1.2	0.995	300K	
	1.2	0.990	300K	
	1.25	1.00	100K	
	1.25	0.995	100K	
	1.25	0.990	100K	
MAGPIE-Pro	1.0	1.00	300K	1M
	1.1	0.995	300K	
	1.2	0.995	300K	
	1.25	0.990	100K	

Table D.2: This table shows the hyper-parameters for supervised fine-tuning.

Hyper-parameter	Value
Learning Rate	2×10^{-5}
Number of Epochs	2
Number of Devices	4
Per-device Batch Size	1
Gradient Accumulation Steps	8
Effective Batch Size	32
Optimizer	Adamw with $\beta s = (0.9, 0.999)$ and $\epsilon = 10^{-8}$
Learning Rate Scheduler	cosine
Warmup Steps	100
Max Sequence Length	8192

Table D.3: This table shows the hyper-parameters for direct preference optimization.

Hyper-parameter	Value
Learning Rate	5×10^{-7}
Number of Epochs	1
Number of Devices	4
Per-device Batch Size	2
Gradient Accumulation Steps	16
Effective Batch Size	128
Optimizer	Adamw with $\beta s = (0.9, 0.999)$ and $\epsilon = 10^{-8}$
Learning Rate Scheduler	cosine
Warmup Ratio	10%

to mitigate the potential repetitive outputs during the generation.

Appendix E

ADDITIONAL EXPERIMENTAL RESULTS***E.1 Performance of Magpie-MT***

Table E.1 compares the performance of MAGPIE-Air-MT and MAGPIE-Pro-MT with their respective single-turn counterparts. We observe that the multi-turn datasets have enhanced performance, particularly in the Arena-Hard benchmark.

Table E.1: This table compares the performance of the multi-turn versions, MAGPIE-Air-MT and MAGPIE-Pro-MT, with their single-turn counterparts. All models are instruction-tuned on the Llama-8B base models.

Dataset		AlpacaEval 2						Arena-Hard
		GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR (%)
		LC (%)	WR (%)	SD	LC (%)	WR (%)	SD	
MAGPIE-Air	Single-Turn	22.66	23.99	1.24	49.27	50.80	1.44	14.9
	MT	22.98	24.02	1.27	49.63	51.42	1.40	15.5
MAGPIE-Pro	Single-Turn	25.15	26.50	1.30	50.52	52.98	1.43	18.9
	MT	24.21	25.19	1.28	52.92	54.80	1.41	20.4

E.2 Compare Magpie and Self-Instruct using Llama-3-8B-Instruct

To compare the performance of MAGPIE and other synthetic dataset generation methods using the same model, we follow the official Self-Instruct [62] setup and generate a 100K supervised fine-tuning dataset using Llama-3-8B-Instruct. For a fair comparison, we select

the first 100K data samples from the MAGPIE-Air dataset generated by Llama-3-8B-Instruct. The performance of models fine-tuned with these two datasets is shown in the table E.2.

Table E.2: This table compares the performance of models fine-tuned using 100K instruction-following datasets generated by Self-Instruct and MAGPIE. All models are supervised-fine-tuned on the Llama-8B base models. We observe that MAGPIE significantly outperforms Self-Instruct across all benchmarks.

Dataset	#Convs	AlpacaEval 2						Arena-Hard
		GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR (%)
		LC (%)	WR (%)	SD	LC (%)	WR (%)	SD	
MAGPIE-Air-100K	100K	20.17	21.33	1.21	46.82	48.76	1.44	15.7
Self-Instruct (Llama-3)	100K	7.21	5.18	0.7	17.86	12.73	1.05	4.0

We observe a significant performance gap between models fine-tuned with datasets generated by Self-Instruct and our MAGPIE. Our analysis revealed that the instruction format in Self-Instruct-generated datasets is predominantly constrained by the patterns defined in the seed instructions, resulting in a lack of diversity. This comparison indicates the novelty of our MAGPIE in generating diverse high-quality instructions without any seed questions.

E.3 Performance of domain-specific and multilingual Magpie datasets

Domain Specific Data Evaluation. We choose code data as representative domain-specific data. We generate domain-specific data using the code instruction system prompt detailed in Appendix F. Using Qwen2.5-72B-Instruct as data generator, we create 100K synthetic code instructions via MAGPIE. We then fine-tune both Llama-3-8B base and Llama-3-8B-Instruct models using this dataset. The models are evaluated on HumanEval [11]. The results shown in Table E.3 demonstrate that our MAGPIE-generated code dataset effectively enhances Llama-3’s performance on code-related tasks, validating MAGPIE’s applicability to

domain-specific instruction tuning.

Table E.3: Performance Comparison on HumanEval.

Alignment Setup	Pass@1	Pass@10	Pass@100
Llama-3-8B-Instruct	0.5574	0.7174	0.8049
Llama-3-8B-base + MAGPIE-Code-100K	0.5327	0.7134	0.8293
Llama-3-8B-Instruct + MAGPIE-Code-100K	0.5768	0.7334	0.8232

Multilingual Data Evaluation. We evaluate MAGPIE’s multilingual capabilities using Chinese as our representative language case. Following the method described in Section 2.2, we use Qwen2-72B-Instruct to generate 200K Chinese synthetic instructions. We then fine-tuned the Llama-3-8B base model with this dataset and evaluated its performance using multilingual MT-Bench. The results are presented in Table E.4. The results demonstrate that models fine-tuned with our Chinese MAGPIE dataset outperform the official Llama-3-8B-Instruct on multilingual MT-Bench (zh-cn). This suggests MAGPIE’s applicability to generate high-quality multilingual datasets.

Table E.4: Performance Comparison on Chinese MT-Bench.

Alignment Setup	Zh MT-Bench
Meta-Llama-3-8B-Instruct	7.75
Llama-3-8B-base + MAGPIE-Chinese-200K	7.80
Llama-3-8B-base + MAGPIE-Chinese-200K + MAGPIE-Pro-MT	7.96

E.4 Ablation Analysis on Data Quantity and Quality

In what follows, we compare within the family of datasets generated by MAGPIE in Table E.5. These datasets differ in size, deployment of filtering, and models used to generate data. We observe that as the dataset’s size increases, the fine-tuned model’s performance improves, indicating that data quantity plays a critical role in enhancing instruction-following capabilities. Furthermore, the model fine-tuned with MAGPIE-Pro-300K-Filtered outperforms those fine-tuned with the same or even higher amounts of raw data. This demonstrates the effectiveness of our filtering technique, and underscores the importance of data quality. Finally, we observe that the models fine-tuned with MAGPIE-Pro consistently outperform those fine-tuned with MAGPIE-Air. The reason is that MAGPIE-Pro is generated by the more capable model, i.e., Llama-3-70B-Instruct.

Table E.5: This table compares MAGPIE datasets within its family that differ in size, deployment of filtering, and models used to generate data. All models are supervised-fine-tuned on the Llama-8B base models.

Dataset	#Convs	AlpacaEval 2						Arena-Hard	
		GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR(%)	
		LC (%)	WR (%)	SD	LC (%)	WR (%)	SD		
MAGPIE-Air	300K-Raw	300K	21.99	21.65	1.21	48.63	48.06	1.42	15.8
	3M-Raw	3M	22.96	21.09	1.20	50.57	48.40	1.42	16.1
	300K-Filtered	300K	22.66	23.99	1.24	49.27	50.8	1.44	14.9
MAGPIE-Pro	300K-Raw	300K	21.65	22.19	1.2	49.65	50.84	1.42	15.9
	1M-Raw	1M	24.16	23.93	1.25	49.97	50.34	1.43	16.7
	100K-Filtered	100K	20.47	24.52	1.25	47.92	52.75	1.43	17.2
	200K-Filtered	200K	22.11	26.02	1.26	51.17	56.76	1.41	15.9
	300K-Filtered	300K	25.08	29.47	1.35	52.12	53.43	1.44	18.9
MAGPIE-Air + MAGPIE-Pro	4M-Raw	4M	24.45	24.08	1.26	51.96	52.08	1.42	15.5

E.5 Ablation Analysis on Filter Designs

We conduct an ablation analysis on various filter designs within MAGPIE-Pro to assess their impact on the performance of supervised fine-tuned models. The results are presented in Table E.6. We observe that different filtering strategies yield optimal performance on different benchmarks, and no single filter consistently achieves the best performance across all benchmarks. Therefore, determining how to select instructional data to enhance the performance in supervised fine-tuning is an interesting topic for future research.

Table E.6: This table compares the performance of different filter designs within MAGPIE-Pro. All models are supervised-fine-tuned on the Llama-8B base models.

Dataset and Filter	AlpacaEval 2						Arena-Hard	
	GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR (%)	
	LC (%)	WR (%)	SD	LC (%)	WR (%)	SD		
MAGPIE-Pro	Filter	25.08	29.47	1.35	52.12	53.43	1.44	18.9
	Filter 2	25.15	26.50	1.30	50.52	52.98	1.43	18.9
	Filter 3	23.90	25.21	1.25	51.45	53.64	1.41	18.3
	Filter 4	24.20	25.33	1.27	52.43	54.34	1.43	17.9
	Filter 5	24.85	25.12	1.26	52.12	53.43	1.44	18.4
	Filter 6	23.20	28.43	1.26	51.34	57.29	1.41	17.9

E.6 Ablation Analysis on Response Generator

To investigate the impact of the response generator on the supervised fine-tuning performance using MAGPIE, we conduct an ablation study by replacing the response generator with Qwen-2-7B-Instruct [68] within MAGPIE-Air-300K-Filtered. We note that the performance of Qwen-2-7B-Instruct is comparable to, or slightly weaker than, Llama-3-8B-Instruct. The results are summarized in Table E.7.

We observe that although there is a slight performance degradation, the model fine-tuned using Qwen-2-7B-Instruct as the response generator still outperforms all baselines, including those using GPT-4 as the response generator. These findings indicate two key points: (1) The success of MAGPIE depends little on the specific response generator used, and (2) the instructions generated by MAGPIE are of high quality and diversity.

Table E.7: This table compares the impact of different response generators on the model performance. All models are supervised-fine-tuned on the Llama-8B base models.

Response Generator	AlpacaEval 2						Arena-Hard
	GPT-4-Turbo (1106)			Llama-3-8B-Instruct			WR (%)
	LC (%)	WR (%)	SD	LC (%)	WR (%)	SD	
Llama-3-8B-Instruct	22.66	23.99	1.24	49.27	50.80	1.44	14.9
Qwen2-7B-Instruct	15.01	15.60	1.05	41.09	42.07	1.47	13.7

E.7 Trustworthiness of Magpie-Aligned Models

In what follows, we conduct more experiments to compare MAGPIE model and Llama-3-8B-Instruct on the TrustLLM benchmark [27]. The results for safety, fairness, ethics, privacy, and robustness are summarized in Table E.8.

We observe that our supervised-fine-tuned model slightly underperforms Llama-3-8B-Instruct in terms of safety and fairness. However, it outperforms the official instruct model on ethics, privacy, and robustness. Considering that our fine-tuned model uses much fewer data samples (300K compared to over 10M), these results again highlight the high quality of data generated by MAGPIE.

E.8 IFEval Evaluations of Magpie-Aligned Models and Baselines

We compare the models fine-tuned with MAGPIE against baselines on IFEval [78] using the LM-Evaluation-Harness framework [24]. The results are presented in Table E.9.

Table E.8: This table compares the performance of model supervised-fine-tuned using MAGPIE-Pro-300K-Filtered and the official Llama-3-8B-Instruct on the TrustLLM benchmark [27].

TrustLLM	Evaluation/Dataset	Llama-3-8B-Instruct	MAGPIE-Pro-300K-Filtered
Safety	Jailbreak (RtA \uparrow)	0.93	0.80
	Misuse (RtA \uparrow)	0.85	0.80
	Exaggerated Safety (RtA \downarrow)	0.54	0.52
Fairness	Stereotype Recognition (Acc \uparrow)	0.49	0.40
	Stereotype Query Test (RtA \uparrow)	1.00	0.99
	Disparagement Sex (p-value \uparrow)	0.99	0.99
	Disparagement Race (p-value \uparrow)	0.55	0.47
Ethics	Social Chemistry 101 (Acc \uparrow)	0.94	0.63
	ETHICS (Acc \uparrow)	0.65	0.69
	MoralChoice (Acc \uparrow)	0.97	0.95
	MoralChoice (RtA \uparrow)	0.97	0.98
Privacy	Privacy Awareness-Normal (RtA \uparrow)	0.33	0.71
	Privacy Awareness-Augmented (RtA \uparrow)	1.00	0.98
	Privacy Leakage (RtA \uparrow)	0.66	0.87
Robustness	AdvGlue (RobustScore \uparrow)	0.42	0.58
	OOD Detection (RtA \uparrow)	0.37	0.26
	OOD Generalization (F1-Score \uparrow)	0.83	0.84

Table E.9: This table compares the performance of models fine-tuned using MAGPIE and other baseline datasets on the IFEval benchmark [78].

Alignment Data	prompt_level_strict	inst_level_strict	prompt_level_loose	inst_level_loose
Self-Instruct (Llama-3)	0.333	0.465	0.372	0.501
ShareGPT	0.331	0.454	0.372	0.492
Evol Instruct	0.344	0.463	0.377	0.494
OpenHermes 1	0.340	0.453	0.377	0.488
Tulu V2 Mix	0.338	0.458	0.370	0.499
WildChat	0.372	0.489	0.423	0.538
OpenHermes 2.5	0.381	0.493	0.436	0.536
GenQA	0.307	0.458	0.331	0.484
Ultrachat	0.298	0.421	0.346	0.466
MAGPIE-Air-300K-Raw	0.366	0.489	0.477	0.590
MAGPIE-Air-300K-Filtered	0.355	0.484	0.475	0.597
MAGPIE-Air-300K-MT	0.368	0.496	0.495	0.614
MAGPIE-Pro-300K-Raw	0.338	0.472	0.455	0.582
MAGPIE-Pro-300K-Filtered	0.298	0.432	0.401	0.529
MAGPIE-Pro-300K-MT	0.336	0.452	0.455	0.568

Our results demonstrate that MAGPIE-generated datasets achieve comparable prompt-level and instruction-level strict accuracy scores to existing baseline datasets. Moreover, MAGPIE exhibits significantly higher performance in both prompt-level and instruction-level loose accuracy metrics. These findings indicate the high quality of MAGPIE-generated datasets.

E.9 Ablation Analysis on the Impact of Reward Models on DPO Performance

To investigate how the choice of reward model influences DPO performance, we performed an ablation study using two reward models: *RLHFlow/ArmoRM-Llama3-8B-v0.1* [60] and *sfairXC/FsfairX-LLaMA3-RM-v0.1* [65]. Notably, ArmoRM exhibited stronger performance on RewardBench [33]. Following the procedure described in Section 3.2, we constructed DPO

datasets based on scores from both reward models and evaluated their performance using AlpacaEval2 and Arena Hard benchmarks.

As shown in Table E.10, our results demonstrate that employing a higher-performing reward model (ArmoRM) yields a better-performing LLM during DPO. This underscores the critical role of reward model quality in enhancing alignment and overall model performance.

Table E.10: Comparison of Llama-3-8B DPO-trained models using different reward models. Metrics are AE2 LC, AE2 WR, and AH Score.

Instruction Model	AE2 LC	AE2 WR	AH Score
Llama-3-8B-Air-DPO-ArmoRM	45.48	50.43	35.9
Llama-3-8B-Air-DPO-sfairXC	37.87	45.08	33.1
Llama-3-8B-Pro-DPO-ArmoRM	50.10	53.53	35.7
Llama-3-8B-Pro-DPO-sfairXC	43.49	50.75	34.2

Appendix F

PROMPT TEMPLATES

F.1 Prompt Templates for Magpie Extension

This section presents the prompt template used to generate MAGPIE-MT and control instruction tasks, as detailed in Figure F.1 and Figure F.2, respectively.

Prompt for generating MAGPIE-MT

```

<|begin_of_text|><|start_header_id|>system<|end_header_id|>

You are a helpful AI assistant. The user will engage in a multi-round conversation with you, asking
initial questions and following up with additional related questions. Your goal is to provide thorough,
relevant and insightful responses to help the user with their queries.<|eot_id|><|start_header_id|>
user<|end_header_id|>

{instruction}<|eot_id|><|start_header_id|>assistant<|end_header_id|>

{response}<|eot_id|><|start_header_id|>user<|end_header_id|>

```

Figure F.1: Prompt for generating MAGPIE-MT. We take Llama-3-8B-Instruct as an example. The placeholder `{instruction}` and `{response}` are from the first turn.

F.2 Prompt Templates for Evaluation

Here, we present the prompt template employed to generate task categories, quality, and difficulty, as detailed in Figure F.3, Figure F.4, and Figure F.5, respectively. The placeholder

System Prompt Template

```
<|begin_of_text|><|start_header_id|>system<|end_header_id|>

{System Prompt}<|eot_id|><|start_header_id|>user<|end_header_id|>
```

System prompt for controlling math instruction tasks

You are an AI assistant designed to provide helpful, step-by-step guidance on solving math problems. The user will ask you a wide range of complex mathematical questions. Your purpose is to assist users in understanding mathematical concepts, working through equations, and arriving at the correct solutions.

System prompt for controlling code instruction tasks

You are an AI assistant designed to provide helpful, step-by-step guidance on coding problems. The user will ask you a wide range of coding questions. Your purpose is to assist users in understanding coding concepts, working through code, and arriving at the correct solutions.

System prompt for controlling translation tasks

You are an AI assistant designed to provide accurate and contextually appropriate translations. Users will ask you to translate text between various languages. Your purpose is to assist users in understanding and conveying meaning across languages, maintaining the original context and nuances.

System prompt for controlling multilingual instruction generation (Japanese + Math)

あなたはAIアシスタントで、数学のを解くために役立つ、ステップバイステップのガイダンスを提供するようにされています。

Figure F.2: Prompts for controlling instruction generation tasks. These examples illustrate how to guide Llama-3-8B-Instruct in generating instructions for specific domains: mathematics, coding, translation, and multilingual tasks. To adapt this approach for different instruction tasks, replace the **System Prompt** placeholder in the System Prompt Template with the appropriate domain-specific prompt.

`input` represents the instructions to be evaluated.

Prompt for generating task categories

```

# Instruction
Please label the task tags for the user query.

## User Query
""{input}""

## Tagging the user input
Please label the task tags for the user query. You will need to analyze the user query and select the most relevant task tag from the list below.

all_task_tags = [
    "Information seeking", # Users ask for specific information or facts about various topics.
    "Reasoning", # Queries require logical thinking, problem-solving, or processing of complex ideas.
    "Planning", # Users need assistance in creating plans or strategies for activities and projects.
    "Editing", # Involves editing, rephrasing, proofreading, or other tasks related to the composition of general written content.
    "Coding & Debugging", # Users seek help with writing, reviewing, or fixing code in programming.
    "Math", # Queries related to mathematical concepts, problems, and calculations.
    "Role playing", # Users engage in scenarios requiring ChatGPT to adopt a character or persona.
    "Data analysis", # Requests involve interpreting data, statistics, or performing analytical tasks.
    "Creative writing", # Users seek assistance with crafting stories, poems, or other creative texts.
    "Advice seeking", # Users ask for recommendations or guidance on various personal or professional issues.
    "Brainstorming", # Involves generating ideas, creative thinking, or exploring possibilities.
    "Others" # Any queries that do not fit into the above categories or are of a miscellaneous nature.
]

## Output Format:
Note that you can only select a single primary tag. Other applicable tags can be added to the list of other tags.
Now, please output your tags below in a json format by filling in the placeholders in <...>:
...
{{
  "primary_tag": "<primary tag>",
  "other_tags": ["<tag 1>", "<tag 2>", ... ]
}}
...

```

Figure F.3: Prompt for generating task categories

Prompt for generating quality of instructions

Instruction

You need to rate the quality of the user query based on its clarity, specificity, and coherence.

The rating scale is as follows:

- very poor: The query is unclear, vague, or incoherent. It lacks essential information and context.
- poor: The query is somewhat unclear or lacks important details. It requires significant clarification.
- average: The query is moderately clear and specific. It may require some additional information for a complete understanding.
- good: The query is clear, specific, and mostly well–formed. It provides sufficient context for understanding the user’s intent.
- excellent: The query is very clear, specific, and well–articulated. It contains all the necessary information and context for providing a comprehensive response.

User Query

““{input}””

Output Format

Given the user query, you first need to give an assessment, highlighting the strengths and/or weaknesses of the user query. Then, you need to output a rating from very poor to excellent by filling in the placeholders in [...]:

““

{{

 "explanation": "[...]",

 "input_quality": "[very poor/poor/average/good/excellent]"

}}

““

””

Figure F.4: Prompt for generating quality of instructions

Prompt for generating difficulty of instructions

Instruction

You first need to identify the given user intent and then label the difficulty level of the user query based on the content of the user query.

User Query

“{input}”

Output Format

Given the user query, in your output, you first need to identify the user intent and the knowledge needed to solve the task in the user query. Then, rate the difficulty level of the user query as ‘very easy’, ‘easy’, ‘medium’, ‘hard’, or ‘very hard’.

Now, please output the user intent and difficulty level below in a json format by filling in the placeholders in [...]:

“

{{

"intent": "The user wants to [...]",

"knowledge": "To solve this problem, the models need to know [...]",

"difficulty": "[very easy/easy/medium/hard/very hard]"

}}

“

Figure F.5: Prompt for generating difficulty of instructions

Appendix G

MAGPIE EXAMPLES

In this section, we demonstrate several representative instructions generated by MAGPIE with different task categories. The following instructions are generated by the Llama-3-70B-Instruct model.

Information Seeking

A few days ago, I was at a restaurant and I got a cup of coffee. However, when I went to take a sip, I realized it was a little too hot. So, I decided to let it cool down for a few minutes. As I waited, I noticed that coffee is actually two different colors. The part that was closest to the surface of the coffee is a lighter color, and the part that is deeper is a darker color. Have you ever observed this phenomenon before?

Reasoning

You have 3 switches but they are not labeled. Each switch corresponds to one of three light bulbs in a room. Each light bulb is either on or off. You can turn the switches on and off as many times as you want, but you can only enter the room one time to observe the bulbs. How can you figure out which switch corresponds to which light bulb?

Planning

You are the Founder of a Financial Planning Company. As a professional financial advisor, you are scheduled to meet a new client tomorrow. Specifically, what are you planning to do to prepare for this meeting?

Editing

What is the best way to re-write the sentence: "We call this the 'core' product or the 'core' offering" using proper quotation marks and avoiding the word "this"?

Coding & Debugging

Write a Python program that calculates the total cost of a customer's order. The program should ask for the customer's name, the number of items they want to purchase, and the price of each item. It should then calculate the total cost by multiplying the number of items by the price of each item and adding 8% sales tax. The program should display the customer's name, the number of items, the price of each item, and the total cost, including sales tax.

Math

In the following problem, please use integers to solve it. A water tank has 1000 L of water. On the first day, $\frac{1}{5}$ of the water was drained. On the second day, $\frac{3}{10}$ of the remaining water was drained. On the third day, $\frac{2}{5}$ of the remaining water was drained. On the fourth day, $\frac{3}{4}$ of the remaining water was drained. How many liters of water are left after the fourth day?

Role Playing

In this game, you will be the host, and I will be the contestant. You will ask me a series of questions, and I will try to answer them correctly. The questions will be multiple choice, and I will have a 25% chance of getting the correct answer if I just randomly guess. However, I am a clever contestant, and I will try to use logic and reasoning to increase my chances of getting the correct answer.

Data Analysis

The personnel manager at a company is tasked with finding the average salary of new hires. She has collected data on the salaries of 13 new hires. She wants to know if there is a statistical difference between the average salary of new hires and the national average salary. The national average salary is \$45,000. The sample of new hires has a mean salary of \$42,800 and a standard deviation of \$4,200.

Creative Writing

Write a paragraph about a mythical creature that you created. The creature is small, no larger than a house cat. It has shimmering scales that reflect light, and can emit a soft, pulsing glow from its body. It has large, round eyes that seem to see right through you, but with a gentle kindness. It has a soft, melodious voice, and can communicate with humans through a form of telepathy.

Advice Seeking

How do you handle stress and overwhelm?

Brainstorming

Can you give me some ideas for a spontaneous, fun and memorable birthday celebration for my partner?

Others

What does "sdrawkcaB" mean?

MAGPIE can also generate domain-specific instructions using models that are tailored to particular fields, as mentioned in Section 3.2. The following instructions are generated by DeepSeek-Coder-V2 [79] and Qwen2-Math-7B-Instruct [68], respectively.

DeepSeek-Coder-V2 (Code Instruction)

You are given a list of emails. You need to write a Python function that returns the domain, excluding the @ symbol, for each email.

Qwen2-Math-7B-Instruct (Math Instruction)

A rectangle with length 12 units and width 8 units is scaled by a factor of 2 to form a new rectangle. Determine the dimensions of the new rectangle and calculate its area. Compare the area of the new rectangle to the area of the original rectangle.

We note that MAGPIE’s capabilities extend beyond generating English datasets to producing diverse multilingual datasets. The following instructions are generated by the Qwen2-72B-Instruct model.

Chinese

从给定的两个整数中找到较大的一个。但是你不能使用任何比较操作符（如 $<$, $>$, $!=$ 等）或数学运算符（如 $+$, $-$, $*$, $/$ 等）来实现它。你只能使用位操作符和逻辑操作符。

German

Das Debate-System 'Oxford-Oberhaus' wird bei ersten Auseinandersetzungen verwendet. Bitte erklären sie, wie dieses System funktioniert.

Spanish

Según la encuesta anual de satisfacción al cliente que acabamos de realizar, parece que la satisfacción general de los clientes con nuestro rendimiento ha disminuido. ¿Podrías preparar una presentación detallada para la reunión del lunes que analice los resultados, identifique las áreas problemáticas y proporcione posibles soluciones basadas en los datos recogidos?

Portuguese

Ho comprato una nuova inchiostriera sulla quale è presente la scritta "Non manipolare". Cosa significa?

Italian

Crie um exemplo de uma conversa entre dois personagens, um MC de hip hop e um pianista clássico, discutindo sobre seus estilos favoritos de música.