



CROWDSELECT: Synthetic Instruction Data Selection with Multi-LLM Wisdom

Anonymous ACL submission

Abstract

Distilling advanced Large Language Models’ instruction-following capabilities into smaller models using a selected subset has become a mainstream approach in model training. While existing synthetic instruction data selection strategies can identify valuable subsets for distillation, they predominantly rely on single-dimensional signals (*i.e.*, reward scores, model perplexity). We argue that such narrow signals may overlook essential nuances of user instructions, especially when each instruction can be answered from multiple perspectives. Therefore, we investigate more diverse signals to capture comprehensive instruction-response pair characteristics and propose three foundation metrics that leverage Multi-LLM wisdom: (1) diverse responses across multiple LLMs and (2) reward model assessment. Based on these metrics, we propose CROWDSELECT, which combines all three metrics with diversity preservation through clustering. Our comprehensive experiments demonstrate that our foundation metrics consistently improve performance across 4 base models on MT-bench and Arena-Hard. Our CROWDSELECT, as an integrated metric, achieves *state-of-the-art* performance in both Full and LoRA fine-tuning, showing improvements of 4.81% on Arena-Hard and 11.1% on MT-bench with Llama-3.2-3b-instruct. We hope our findings will bring valuable insights for future research in this direction.

1 Introduction

In recent years, Large Language Models (LLMs) (Achiam et al., 2023; Jaech et al., 2024; Team et al., 2024; Guo et al., 2025) have demonstrated remarkable capability in following user instructions to generate coherent and contextually helpful responses (Jiang et al., 2023; Zheng et al., 2023b; Wen et al., 2024). Yet, the computational overhead for instruction tuning and massive parameter sizes of these models create a considerable barrier to practical deployment (Peng et al., 2023). To address this,

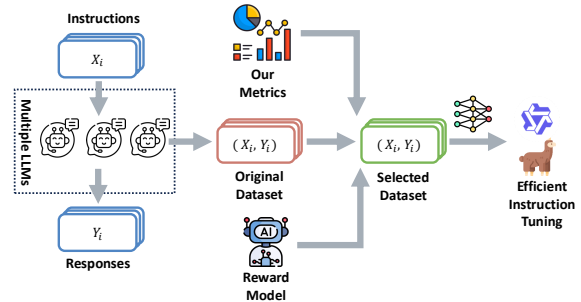


Figure 1: A demonstration of instruction tuning with selected synthetic instruction-response pairs.

many approaches distill the instruction-following prowess of advanced LLMs into smaller, more efficient models through a small-amount instruction tuning with synthetic responses (Xia et al., 2024; Zhou et al., 2024a).

A critical bottleneck, however, lies in selecting the right data for this distillation process. Most existing data selection methods rely on predefined rules (Chen et al., 2023a), automated single-dimensional signals — such as reward scores (Wu et al., 2024b; Lambert et al., 2024) or difficulty metrics (Li et al., 2023b, 2024b) — to identify valuable examples for fine-tuning. While effective to a point, such narrow signals may overlook essential nuances of user instructions, especially when each instruction can be answered from multiple perspectives (Händler, 2023; Feng et al., 2025). This raises a fundamental question: “Can we leverage multi-dimensional signals to better reflect the various facets of each sample for more effective instruction tuning data selection?”

Inspired by previous work that leverage Multi-LLMs collaboration, we take an explorative step toward more robust and comprehensive data selection by introducing CROWDSELECT, a framework that harnesses pre-collected Multiple LLMs’ responses and their reward scores, treating as different reflection of the instruction to leverage Multi-LLM Wis-

dom. Instead of treating each instruction–response pair in isolation—typically derived from just one model’s output—our method aggregates multiple responses for each instruction from a diverse set of LLMs. Crucially, we also factor in each response’s reward as provided by state-of-the-art reward models. This multi-view setup captures more “facets” of each instruction, illuminating subtle differences in how various models handle the same query. Based on these observations, we propose three base explorative metrics:

- **Difficulty** - Identifies instructions on which the majority of models struggle, surfacing challenging prompts critical to learning.
- **Separability** – Highlights instructions whose response quality exhibits high variance across models, making them especially useful for differentiating stronger from weaker capabilities.
- **Stability** – Measures how consistently model performance follows expected size-based ranking across families, ensuring the selected data helps reinforce well-grounded alignment signals.

Our exploratory experiments in FFT and low-rank adaptation (LoRA) (Hu et al., 2021) experiments on Llama-3.2-3b-base/instruct (Dubey et al., 2024) and Qwen-2.5-3b-base/instruct (Yang et al., 2024b) demonstrate the robustness and efficacy of our proposed metrics through significant performance gaps between *top-scored* and *bottom-scored* data subset fine-tuning, with potential further improvements through metric combination.

Subsequently, we propose CROWDSELECT that combines these metrics with a clustering strategy to preserve diversity and explore the upperbound of leveraging Multi-LLM wisdom to identify a compact yet high-impact subset of instruction–response data. Experimental results show that models fine-tuned on our selected subset significantly outperform baselines and previous state-of-the-art data selection methods, achieving improvements of 4.81% on Arena-Hard and 11.1% on MT-bench with Llama-3.2-3b-instruct. Furthermore, CROWDSELECT achieves *state-of-the-art* performance across four models on two benchmarks, demonstrating both the generalizability and robustness of our selected data and methodology, paving a new dimension for efficient instruction tuning.

Our contributions are summarized as follows:

- **Investigation of Multi-LLM Wisdom in Instruction Data Selection.** We propose a novel

approach that utilizes multiple synthesized responses from different LLMs for each instruction, enhancing the diversity and quality of data.

- **Novel Metrics and Methods.** We design three new explorative base metrics—*Difficulty*, *Separability*, and *Stability*—that leverage multi-LLM responses and reward scores as more comprehensive signals, and combine them into CROWDSELECT to explore upperbound in selecting high-quality data for instruction tuning.
- **State-of-the-art Performance.** We demonstrate that combining our metrics and clustering techniques for data selection leads to a new SOTA in efficient instruction tuning in both Llama-3.2-3b and Qwen-2.5-3b.

2 Related Work

Instruction Tuning Data Selection. Instruction Tuning stands out to be a method to solve the gap between pretrained knowledge and real-world user scenarios (Ouyang et al., 2022; Bai et al., 2022). Recent efforts like Vicuna (Peng et al., 2023) and LIMA (Zhou et al., 2024a) demonstrate high performance with a carefully selected small dataset, highlighting the growing importance of efficient instruction tuning. Three key metrics determine instruction data quality: *Difficulty*, *Quality*, and *Diversity*. *Difficulty*, focusing mainly on the question side, is considered more valuable for model learning (Li et al., 2023b, 2024b; Liu et al., 2024b; Lee et al., 2024; Wang et al., 2024b). *Quality*, mainly addressing the response side, measures the helpfulness and safety of model responses, typically assessed using LLM evaluators (Chen et al., 2023a, 2024b; Liu et al., 2024c; Ye et al., 2024), reward models (Son et al., 2024; Lambert et al., 2024), and gradient similarity search (Xia et al., 2024). *Diversity* also plays a crucial role in covering various instruction formats and world knowledge, primarily improving model robustness (Bukharin and Zhao, 2023; Wang et al., 2024d).

Data Synthesis for Instruction Tuning. While the development of LLMs initially relied on human-curated instruction datasets for instruction tuning (Zheng et al., 2023a; Zhao et al., 2024; Lightman et al., 2023), this approach proved time-consuming and labor-intensive, particularly as the complexity and scope of target tasks increased (Demrozi et al., 2023; Wang et al., 2021). Consequently, researchers began exploring the use of frontier LLMs

to generate synthetic instruction datasets, aiming to both address these scalability challenges (Ding et al., 2023; Chen et al., 2023b, 2024d) and leverage models’ advanced capabilities in developing next-generation foundation models (Burns et al., 2023; Charikar et al., 2024). Recent advancements streamline this process by utilizing instructions directly from pretrained LLMs with simple prompt templates (Xu et al., 2024a; Chen et al., 2024c; Zhang et al., 2024), significantly reducing the required custom design from human effort.

Deriving Crowded Wisdom from Multi-LLM.

Single LLM’s response to a question face limitations in its representation of data (particularly cutting-edge knowledge) (Lazaridou et al., 2021; Dhingra et al., 2022; Kasai et al., 2023), skills (as no single LLM is universally optimal *empirically*) (Sun et al., 2022; Liang et al., 2022; Chen et al., 2024a), and diverse perspectives (Feng et al., 2025). Previous work has demonstrated that *online* multi-LLM wisdom (also known as compositional agent frameworks (Gupta and Kembhavi, 2023)) tends to outperform single models across various domains, providing more comprehensive and reflective solution on complex downstream tasks (Wang et al., 2024c; Wu et al., 2023; Li et al., 2023a; Ouyang et al., 2025; Gui et al., 2025). *Offline* crowded wisdom, where data are pre-collected rather than real-time inference, also show potential in model alignment (Gallego, 2024; Rafailov et al., 2023; Meng et al., 2025) and benchmark construction (Ni et al., 2024b,a). In this paper, we pioneer the use of *offline* multi-LLM wisdom for instruction data selection by utilizing these LLMs’ responses and their reward score as *reflections* to measure instruction-response pairs’ *Difficulty* and *Quality*.

3 Methodology

In this section, we first define our synthetic data selection task and propose three foundational metrics that leverage responses and assessment scores from multiple advanced LLMs. We then introduce CROWDSELECT, which combines these metrics with diversity-preserving clustering to explore the upper bounds of Multi-LLM Wisdom.

3.1 Preliminaries

We formulate the instruction quality as the consensus among N LLMs. Given an instruction-tuning dataset, we extract all instructions from the dataset to form instruction dataset Q . For each instruction

$q_i \in Q$, a response set R_i is obtained by querying multiple LLMs. An assessment model then evaluates the response set R_i to form the score set C_i^M based on metrics M . The index of M is omitted for brevity in the following context unless specified. The top- k instruction subset of metric M is defined as

$$S_k^M = \arg \max_{S \subset \mathcal{S}, |S|=k} M(C_i^M) \quad (1)$$

where S_k^M consists of the k instructions that maximize the metric M .

The corresponding response r_i^M for each instruction q_i^M from the instruction subset S_k^M is subsequently obtained by

$$r_i^M = \text{Top}(R_i, C_i^M) \quad (2)$$

where $\text{Top}(R_i^S, C_i^M)$ denotes the best responses in r_i^S ranked by C_i^M . The produced instruction-answer subset $\tilde{Q} = \{(r_i^M, q_i^M)\}$ is then utilized for finetuning as an alternative of the original dataset.

3.2 Base Metrics

In this section, we introduce three new base metrics to leverage multiple LLMs’ responses and their reward scores as various “*facets*” to reflect the value of each sample.

Difficulty. The difficulty score C^{diff} is defined as the negative value of the average score, which is the mean score of all the model responses for a given instruction.

$$C^{diff} = -\frac{\sum C_i^M}{N} \quad (3)$$

Higher *difficulty* indicates more challenging instructions. This metric is particularly well-suited for fine-tuning on reasoning tasks, e.g. mathematics and planning, where the goal is often to improve performance on complex problems. By focusing on instructions with higher *difficulty*, we prioritize examples that are likely to be answered incorrectly by the majority of models. This ensures that the fine-tuning dataset includes a substantial proportion of challenging instructions, maximizing the model’s exposure to difficult material and potentially leading to greater improvements in performance.

Separability. The separability score C^{sep} is defined as the score variance, which is the variance of all the response scores for an instruction.

$$C^{sep} = \text{var}(C_i^M) \quad (4)$$

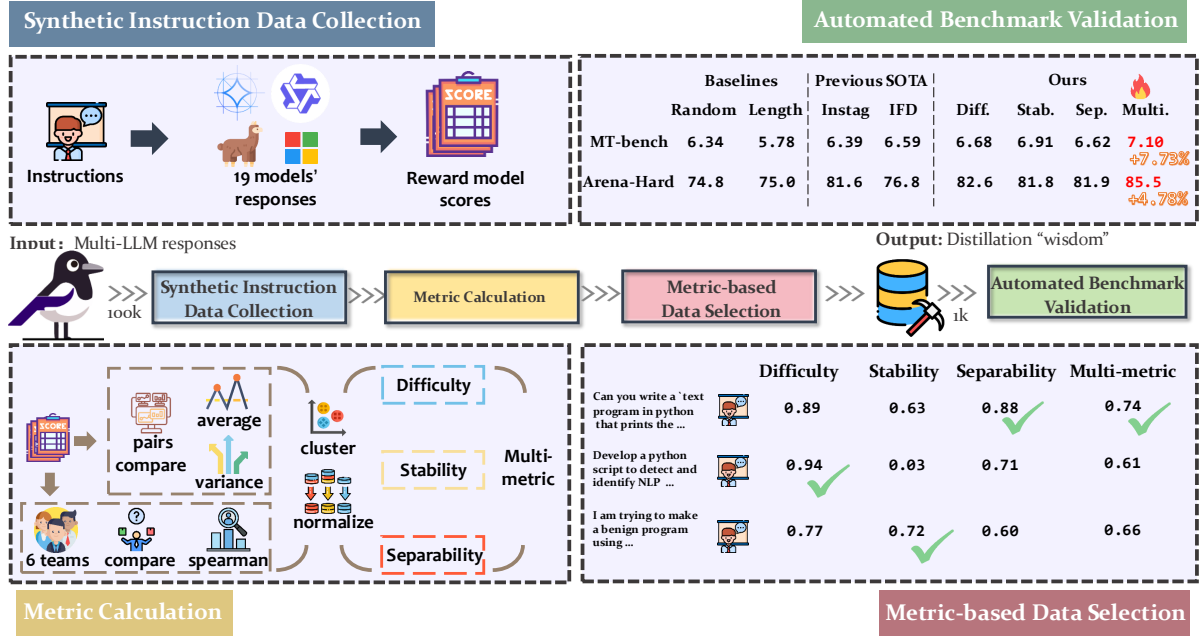


Figure 2: The overall pipeline of our CROWDSELECT, which innovatively leverages metrics calculated from multiple facets of instructions using pre-collected synthesized responses from various LLMs and their corresponding reward model scores. We enhance data selection through clustering for diversity and metric combination to explore the method’s potential. Finally, we evaluate the effectiveness of our selected instruction subset through FFT or LoRA fine-tuning (Hu et al., 2021) for efficient instruction tuning.

Higher *Separability* indicates that a considerable proportion of models cannot perform well on the instruction, thus this instruction is more effective in differentiating between models. This characteristic makes the *Separability* particularly well-suited for curating datasets of knowledge remembering or preference alignment. In such datasets, some models may exhibit strong performance while others struggle. By selecting instructions with high separability, we prioritize examples that effectively distinguish between these varying levels of competence. These “discriminatory” examples are valuable because they provide the fine-tuned model with opportunities to learn from the specific challenges that differentiate successful models from less successful ones. Focusing on these examples enforces the finetuned model to handle the nuances and complexities that separate high-performing models.

Stability. *Stability* is defined as the average spearman factor, which is the mean of five spearman factors, corresponding to five model families. The spearman factor is calculated based on r^a and r^b :

$$\frac{\frac{1}{n} \sum_{i=1}^n (r_i^a - \bar{r}^a) \cdot (r_i^b - \bar{r}^b)}{\sqrt{\left(\frac{1}{n} \sum_{i=1}^n (r_i^a - \bar{r}^a)^2\right) \cdot \left(\frac{1}{n} \sum_{i=1}^n (r_i^b - \bar{r}^b)^2\right)}} \quad (5)$$

- r^a refers to the original ranking within a model family, where models with larger pa-

rameters are theoretically ranked higher, naturally aligning with the performance rank.

- r^b is determined by the ranking of models based on their response scores (e.g., if LLaMA-3B has a response score of 90 and LLaMA-8B has a response score of 75, then 3B ranks higher than 8B within the LLaMA family).

Stability effectively captures how well performance rankings align with expected model size rankings using Spearman’s rank correlation (Schober et al., 2018), making it robust to variations in score scales and non-linear relationships. Averaging across model families further strengthens the robustness of the score, alleviating performance gaps among model families.

3.3 CROWDSELECT: Explore the Upperbound with Multi-LLM Wisdom

Diversity Preservation with Clustering. To facilitate clustering, all instructions were embedded into a fixed-dimensional latent space using a pre-trained embedding model. Within each cluster, instructions were then ranked with the given metric, and the highest-ranked instructions were selected. To avoid over-representing dominant clusters and neglecting potentially valuable information con-

tained within smaller or less frequent clusters, we draw equally from each cluster to form a more robust and generalizable subset.

Multi-metric Integration Building upon the cluster-based selection strategy, we introduce a multi-metric approach to leverage the diverse information captured by the difficulty, separability, and stability scores. Each instruction-response pair is thus characterized by a vector of associated scores, reflecting its various attributes. However, these metrics exhibit different distributions, ranges, and magnitudes. Therefore, we employ a three-stage normalization process to ensure equitable contribution from each metric.

Specifically, each metric score is standardized to standard normal distribution. The standardized scores are then normalized to $[0, 1]$ using a min-max scaling approach. Finally, to further refine the distribution and mitigate the impact of potential outliers, we apply a quantile transformation that maps the normalized scores to a uniform distribution between $[0, 1]$.

$$Z_i^M = \frac{(C_i^M - \mu^M)}{\sigma^M} \quad (6)$$

$$N_i^M = \frac{(Z_i^M - \min(Z^M))}{(\max(Z^M) - \min(Z^M))} \quad (7)$$

$$\rho_i^M = \text{quant}(N_i^M | N^M) \quad (8)$$

Following this normalization procedure, we aggregate the transformed scores into a single multi-metric score \hat{C} for each instruction-response pair. This aggregation is performed using a weighted sum of the proposed metrics:

$$\hat{C}_i = \sum_j w_i * \rho_i^{M_j} \quad (9)$$

where $\rho_i^{M_j}$ represent the quantile-transformed scores for metric j , and w_i are the corresponding weights assigned to each metric. This weighted multi-metric approach, combined with the preceding normalization steps, ensures a balanced and robust data selection process that leverages the complementary information provided by the different metrics.

4 Experiment

In this section, we first validate our base metrics through comparative experiments between top-scored and bottom-scored data subsets. We then

evaluate CROWDSELECT against existing baselines and *state-of-the-art* methods. Finally, we conduct an ablation study to analyze the contribution of each sub-module within CROWDSELECT.

4.1 Experiment Setups

Datasets. We conduct our experiments on Magpie-100K-Generator-Zoo¹ given that it directly matches our problem setting that contains answers from 19 models—Qwen2 (Yang et al., 2024a), Qwen2.5 (Yang et al., 2024b), Llama 3 (Dubey et al., 2024), Llama 3.1 (Dubey et al., 2024), Gemma 2 (Team et al., 2024), Phi-3 (Abdin et al., 2024) families and GPT-4 (Achiam et al., 2023)—and their reward scores from three state-of-the-art reward models from RewardBench (Lambert et al., 2024): ArmoRM-Llama3-8B-v0.1 (Wang et al., 2024a), Skywork-Reward-Llama-3.1-8B (Liu and Zeng, 2024), and Skywork-Reward-Gemma-2-27B (Liu and Zeng, 2024).

Evaluation. To evaluate the instruction-following capabilities, we use two widely-used instruction-following benchmarks: MT-Bench (Zheng et al., 2023b) and Arena-Hard (Li et al., 2024c). Both benchmarks mainly leverage LLM-as-a-Judge (Zheng et al., 2023b) for evaluation, while MT-Bench leverage 1-10 rating scoring and Arena-Hard leverage direct pairwise comparison and finally provide a leaderboard with one model as anchor-points. In our experiment, we set the base model (*i.e.*, LLaMA-3.2-3B-base) as the anchor point for models for arena battles. We unified the LLM-as-a-Judge model in both benchmarks as DeepSeek-V3 (Liu et al., 2024a) through official API² and Together API³ given its high performance on natural language generation tasks. Thanks to the unified judge model, we additionally report the **Average Performance (AP)** as a ranking computed by the ranking in MT-Bench and Arena-Hard. **Each experiment was conducted 3 times. The average results are reported to ensure the reliability and reproducibility.**

Base Models. Following Xu et al. (2024b), we consider four small models from different developers as student models, including base

¹<https://huggingface.co/datasets/Magpie-Align/Magpie-100K-Generator-Zoo>

²<https://platform.deepseek.com/>

³<https://api.together.ai/>

Table 1: Validation of our three foundation metrics on Full fine-tuning Llama-3b-instruct with *top-scored* (\uparrow) and *bottom-scored* (\downarrow) instruction selection and different response selection strategy. Best and second results for each metric is highlighted in **bold** and underline.

Strategy	DirectScore	Difficulty		Separability		Stability		Multi
		\downarrow	\uparrow	\downarrow	\uparrow	\downarrow	\uparrow	
MT-Bench								
Best-answer	4.406	4.506	<u>4.738</u>	4.731	5.056	4.675	5.088	5.125
Random	4.470	4.469	4.688	4.695	4.785	4.500	4.581	4.613
Top5-random	<u>4.435</u>	4.681	4.870	4.788	<u>5.008</u>	4.619	<u>4.956</u>	<u>5.048</u>
Arena-Hard								
Best-answer	75.3 _(-2.0, 1.6)	<u>78.6</u> _(-1.9, 2.1)	76.8 _(-1.6, 1.7)	81.8 _(-1.8, 1.2)	83.3 _(-1.8, 1.7)	<u>80.0</u> _(-1.5, 1.6)	82.3 _(-1.6, 2.2)	85.5 _(-0.8, 1.1)
Random	<u>74.5</u> _(-1.1, 1.2)	78.5 _(-1.6, 1.3)	80.4 _(-1.0, 1.5)	79.0 _(-1.3, 1.4)	80.6 _(-1.6, 1.6)	76.2 _(-0.8, 1.6)	77.0 _(-1.0, 1.8)	82.3 _(-1.2, 1.3)
Top5-random	73.7 _(-1.2, 1.8)	75.9 _(-1.6, 1.5)	76.8 _(-1.2, 1.4)	<u>82.0</u> _(-1.3, 1.2)	80.0 _(-0.7, 1.3)	75.0 _(-4.4, 5.8)	76.9 _(-1.4, 1.6)	<u>83.1</u> _(-1.4, 1.7)

and instruct models—Qwen-2.5-3B, Qwen-2.5-3B-Instruct (Yang et al., 2024b) and LLaMA-3.2-3B, LLaMA-3.2-3B-Instruct (Dubey et al., 2024). We use 10 clusters for diversity preservation, and the multimetric setting uses $w = (1, 1, 2)$ for metric intergration in the following experiments.

Baselines. We include 7 baselines in our experiments. *Random*, denotes a randomly selected instruction-answer set from the original dataset. We also compared two previous *state-of-the-art* data selection method: Instag (Lu et al., 2023), and IFD (Li et al., 2023b). For rule-based method, We include *Length* and *Reward Score* (Liu et al., 2023). More details are shown in Appendix B.3.

Instruction-Tuning Setups. We conduct our fine-tuning and evaluation on single A800 and A6000 servers. For fine-tuning, we use LLaMA-Factory (Zheng et al., 2024). For evaluation, we leverage the official codebase of MT-Bench⁴ and Arena-Hard⁵ for automatic assessments. See Appendix B for more details of experiment setups.

4.2 Experiment Results.

Three foundation metrics demonstrate effectiveness in selecting valuable samples. As shown in Table 1, our three foundation metrics consistently identify valuable instruction samples across all response selection strategies. Models fine-tuned on *Top-scored* samples consistently outperform *Bottom-scored* samples, with *Stability* exceed the most margin. We also explore the response selection strategies to build a foundation for following experiments. *Best-answer* setting outperforms

both *Random* and *Top5-random* approaches, indicating that responses with higher reward scores provide better quality data for distillation. This consistent performance across individual metrics establishes strong foundation for further improvements through integration. Therefore, we use *top-scored* as the instruction selection and *Best-answer* as the corresponding response for all experiments.

CROWDSELECT achieves new state-of-the-art performance on both benchmarks. As shown in Table 2, our approach significantly outperforms previous baselines across four models, demonstrating robust generalization. On Arena-Hard and MT-bench, CROWDSELECT with Llama-3.2-3b-instruct achieves scores of 85.5 and 7.103 respectively, surpassing the previous best results by 4.81% and 11.1%. For Qwen-2.5-3b-instruct, CROWDSELECT outperforms the strongest baseline by 3.90%, validating our approach of post-training with high-quality instructions and model distillation. Even for base models, our foundation metrics and CROWDSELECT prove effective, notably improving Llama-3.2-3b’s performance on MT-bench by 12.3%.

CROWDSELECT metrics perform robust on various finetuning methods. Beyond demonstrating superior performance on standard benchmarks, the proposed metrics were further evaluated for robustness across a range of fine-tuning methodologies. Table 1 revealed consistent and stable performance of the proposed metrics. This robustness across varying training paradigms highlights the generalizability of the metrics and suggests their applicability in a wider range of practical scenarios.

4.3 Ablation Studies

In this section, we conduct ablation studies for each module in CROWDSELECT to provide a compre-

⁴https://github.com/lm-sys/FastChat/tree/main/fastchat/llm_judge

⁵<https://github.com/lmarena/Arena-Hard-auto>

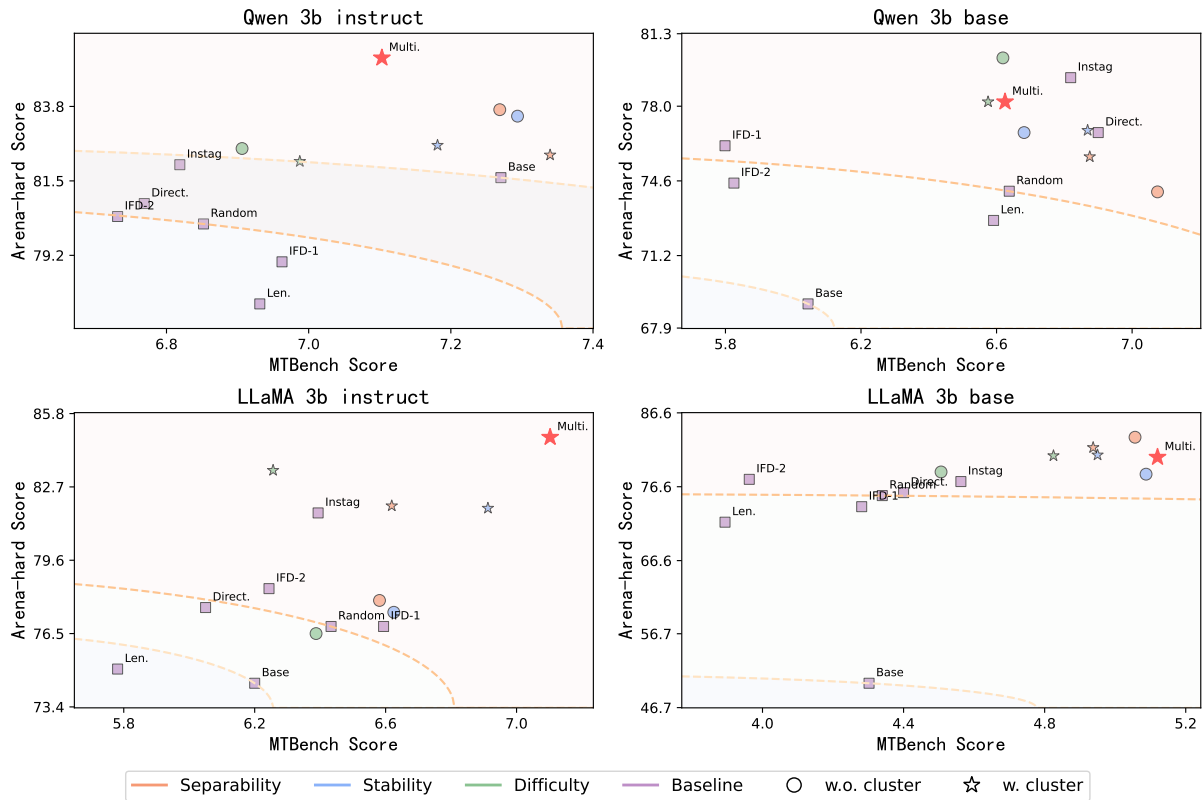


Figure 3: Overall results demonstrate that our foundation metrics and CROWDSELECT consistently outperform baseline methods by a significant margin across FFT settings of four models, with particularly strong performance improvements on Llama-3b-instruct.

hensive analysis of our approach.

Dataset Size. Cao et al. (2023) suggests that selecting concise subsets from all datasets yield competitive results. Following this finding, we collect 1k instruction-response pairs overall in our main experiments. Further experiments on various dataset sizes also support this finding. From the results in 4, small elite datasets behaves on par with a large dataset. This highlights the importance of data quality over sheer quantity in instruction tuning.

Metric Coefficient Combination Our experiments explored various coefficient combinations to determine the optimal balance for creating high-quality, robust datasets. Table 3 details the process of optimizing the weights assigned to different metrics when evaluating dataset quality. As shown in the table, the coefficient combination $w = (1, 1, 2)$ consistently yielded superior results compared to other tested combinations.

Number of Clusters Clustering’s impact on dataset quality was investigated by varying the number of clusters during dataset construction (see Table 4). While no strong positive correlation was

observed between cluster count and quality, all clustered datasets outperformed those constructed without clustering. The results highlight the importance and robustness of the clustering process.

Response Generation Strategy The response generation strategy largely affects the generation quality of the finetuned LLM. Table 1 shows that the best-answer strategy substantially outperforms other strategies, highlighting the importance of dataset response quality. We argue the reason for strategy independence of the difficulty metric is that the core challenge in these instructions is inherently tied to the complexity of the task, not in the method of response formulation. For instance, a highly challenging instruction may require the model to synthesize information from multiple domains, reason through abstract concepts, or produce detailed, contextually rich outputs. These demands remain consistent, regardless of response generation strategies.

Further experiments on finetuning with LoRA and reward model selection are also presented in Appendix C.

Table 2: Performance comparison of full finetuned Llama3.2-3b-base/instruct and Qwen2.5-3b-base/instruct models with different data selection strategies. The best and second results are in **bold** and underline.

Benchmark	Base	Baselines			Our Metrics			
		Random	Tags	IFD	Difficulty	Separability	Stability	Multi
Llama3.2-3b-base								
MT-Bench	4.302	4.406	4.562	3.962	4.738	5.056	<u>5.088</u>	5.125
Arena-Hard	50.0(-0.0, 0.0)	75.3(-2.0, 1.6)	77.3(-1.1, 1.2)	77.6(-1.6, 1.6)	76.8(-1.6, 1.7)	83.3 (-1.8, 1.7)	78.3(-1.6, 2.2)	<u>80.6</u> (-2.4, 1.6)
Llama3.2-3b-instruct								
MT-Bench	6.200	6.356	6.393	6.243	<u>6.648</u>	6.581	6.625	7.103
Arena-Hard	74.4(-1.0, 1.5)	74.8(-1.5, 1.6)	<u>81.6</u> (-0.2, 0.2)	78.4(-1.7, 1.5)	80.5(-0.9, 1.3)	77.9(-1.5, 1.7)	77.4(-1.5, 1.1)	85.5 (-0.8, 1.1)
Qwen2.5-3b-base								
MT-Bench	6.043	6.500	<u>6.818</u>	5.825	6.613	7.075	6.681	6.625
Arena-Hard	69.0(-2.2, 1.6)	72.9(-2.2, 1.9)	<u>79.3</u> (-2.2, 1.9)	74.5(-1.5, 1.5)	73.8(-2.5, 1.8)	74.1(-1.6, 2.4)	76.8(-1.8, 1.8)	79.9 (-1.6, 1.8)
Qwen2.5-3b-instruct								
MT-Bench	7.138	6.793	6.818	6.731	7.182	<u>7.269</u>	7.294	7.131
Arena-Hard	81.6(-1.8, 1.4)	78.2(-1.7, 2.0)	82.0(-2.4, 1.6)	80.4(-1.3, 1.0)	81.8(-1.6, 1.3)	<u>83.7</u> (-1.4, 1.2)	83.5(-1.4, 1.4)	85.2 (-1.2, 1.1)

Table 3: Hyperparameter comparison of CROWDSELECT using Llama-3b-instruct models with varying cluster numbers. The sequence represents (*Difficulty, Separability, Stability*).

Hyperparameter	MT-Bench	Arena-Hard
1_1_1	6.913	81.8(-0.5, 0.8)
1_-1_1	6.625	84.2(-0.7, 1.0)
1_1_2	7.103	85.5 (-0.8, 1.1)
1_1_-1	6.650	82.7(-1.5, 1.4)
1_1_1.5	6.850	84.7(-1.6, 1.3)
1_-1_1.5	6.781	83.0(-1.4, 1.4)
-1_-1_1	6.781	81.9(-1.5, 1.3)
-1_-1_2	6.838	84.8(-1.3, 1.2)
-1_-1_1.5	6.638	81.8(-1.3, 1.3)

Table 4: Performance comparison of FFT-version of Llama-3b-instruct on different coefficient combinations for multiple metrics with clustering.

Benchmark	Random	Difficulty	Separability	Stability
10 clusters				
MT-Bench	6.443	6.675	6.619	6.913
Arena-Hard	80.9	82.6	81.9	81.8
Arena-Hard-95%CI	(-1.3, 1.4)	(-1.2, 1.8)	(-1.7, 1.7)	(-1.5, 1.7)
20 clusters				
MT-Bench	6.607	6.615	6.591	6.686
Arena-Hard	82.8	83.1	85.2	82.8
Arena-Hard-95%CI	(-1.2, 1.4)	(-1.1, 1.7)	(-1.3, 1.1)	(-1.4, 1.1)
30 clusters				
MT-Bench	6.721	6.737	6.725	6.562
Arena-Hard	83.2	84.9	83.3	83.8
Arena-Hard-95%CI	(-1.3, 1.1)	(-1.0, 1.1)	(-1.4, 1.4)	(-1.4, 1.2)

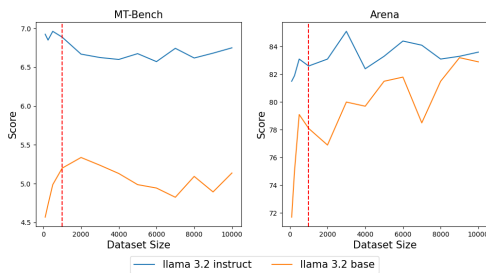


Figure 4: Results show that small elite datasets behaves on par with a large dataset. Our implementation (line in red) achieves reasonably good results on all scenarios.

5 Conclusion

This paper presents novel metrics for synthetic instruction data selection based on Multi-LLM Wisdom, capturing the *difficulty* of instructions from multiple perspectives through various LLMs' responses and their corresponding reward scores. We validate our hypothesis through the strong perfor-

mance of individual metrics on both MT-Bench and Arena-Hard using FFT and LoRA fine-tuning on Llama-3.2-3b and Qwen-2.5-3b. By combining diversity enhancement through clustering with our proposed metrics, CROWDSELECT consistently outperforms *state-of-the-art* data selection methods, establishing both new perspectives and a robust baseline for instruction tuning data selection.

Limitations

While CROWDSELECT demonstrates significant improvement in synthetic data selection tasks, we acknowledge several limitations. Our approach computes data selection metrics by leveraging responses from multiple models across different model families and their corresponding reward from model scores. However, this methodology may be susceptible to reward model biases, including potential reward hacking issues. Although a more

543	organic integration of multiple reward scores could	Dongping Chen, Ruoxi Chen, Shu Pu, Zhaoyi Liu,	594
544	potentially enhance robustness, the computation	Yanru Wu, Caixi Chen, Benlin Liu, Yue Huang, Yao	595
545	of these scores requires additional computational	Wan, Pan Zhou, et al. 2024a. Interleaved scene graph	596
546	resources. Furthermore, our experiments were con-	for interleaved text-and-image generation assessment.	597
547	ducted on both A800 and A6000 GPUs, and the	<i>arXiv preprint arXiv:2411.17188.</i>	598
548	variation in hardware environments may introduce	Dongping Chen, Ruoxi Chen, Shilin Zhang, Yinuo Liu,	599
549	some instability and affect experimental results, po-	Yaochen Wang, Huichi Zhou, Qihui Zhang, Yao	600
550	tentially impacting the reproducibility of our find-	Wan, Pan Zhou, and Lichao Sun. 2024b. Mllm-	601
551	ings.	as-a-judge: Assessing multimodal llm-as-a-judge	602
		with vision-language benchmark. <i>arXiv preprint</i>	603
		<i>arXiv:2402.04788.</i>	604
552	References	Jiuhai Chen, Rifaa Qadri, Yuxin Wen, Neel Jain, John	605
553	Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed	Kirchenbauer, Tianyi Zhou, and Tom Goldstein.	606
554	Awadallah, Ammar Ahmad Awan, Nguyen Bach,	2024c. Genqa: Generating millions of instructions	607
555	Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat	from a handful of prompts. <i>ArXiv</i> , abs/2406.10323.	608
556	Behl, et al. 2024. Phi-3 technical report: A highly ca-	Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa	609
557	pable language model locally on your phone. <i>arXiv</i>	Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniv-	610
558	<i>preprint arXiv:2404.14219.</i>	asan, Tianyi Zhou, Heng Huang, et al. 2023a. Al-	611
559	Josh Achiam, Steven Adler, Sandhini Agarwal, Lama	pagasus: Training a better alpaca with fewer data.	612
560	Ahmad, Ilge Akkaya, Florencia Leoni Aleman,	<i>arXiv preprint arXiv:2307.08701.</i>	613
561	Diogo Almeida, Janko Altenschmidt, Sam Altman,	Lin Chen, Jinsong Li, Xiao wen Dong, Pan Zhang, Cong-	614
562	Shyamal Anadkat, et al. 2023. Gpt-4 technical report.	hui He, Jiaqi Wang, Feng Zhao, and Dahua Lin.	615
563	<i>arXiv preprint arXiv:2303.08774.</i>	2023b. Sharegpt4v: Improving large multi-modal	616
564	Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne	models with better captions. In <i>European Confer-</i>	617
565	Longpre, Nathan Lambert, Xinyi Wang, Niklas	<i>ence on Computer Vision.</i>	618
566	Muennighoff, Bairu Hou, Liangming Pan, Hae-	Lin Chen, Xilin Wei, Jinsong Li, Xiao wen Dong, Pan	619
567	won Jeong, Colin Raffel, Shiyu Chang, Tatsunori	Zhang, Yuhang Zang, Zehui Chen, Haodong Duan,	620
568	Hashimoto, and William Yang Wang. 2024. A sur-	Bin Lin, Zhenyu Tang, Li Yuan, Yu Qiao, Dahua Lin,	621
569	vey on data selection for language models. <i>ArXiv</i> ,	Feng Zhao, and Jiaqi Wang. 2024d. Sharegpt4video:	622
570	abs/2402.16827.	Improving video understanding and generation with	623
571	Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda	better captions. <i>ArXiv</i> , abs/2406.04325.	624
572	Askeff, Anna Chen, Nova DasSarma, Dawn Drain,	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao,	625
573	Stanislav Fort, Deep Ganguli, Tom Henighan, et al.	Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and	626
574	2022. Training a helpful and harmless assistant with	Maosong Sun. 2023. Ultrafeedback: Boosting lan-	627
575	reinforcement learning from human feedback. <i>arXiv</i>	guage models with high-quality feedback. <i>ArXiv</i> ,	628
576	<i>preprint arXiv:2204.05862.</i>	abs/2310.01377.	629
577	Alexander Bukharin and Tuo Zhao. 2023. Data diversity	Florenc Demrozi, Cristian Turetta, Fadi Al Machot,	630
578	matters for robust instruction tuning. <i>arXiv preprint</i>	Graziano Pravadelli, and Philipp H. Kindt. 2023. A	631
579	<i>arXiv:2311.14736.</i>	comprehensive review of automated data annotation	632
580	Collin Burns, Pavel Izmailov, Jan Hendrik Kirchner,	techniques in human activity recognition. <i>ArXiv</i> ,	633
581	Bowen Baker, Leo Gao, Leopold Aschenbrenner,	abs/2307.05988.	634
582	Yining Chen, Adrien Ecoffet, Manas R. Joglekar,	Bhuwan Dhingra, Jeremy R Cole, Julian Martin	635
583	Jan Leike, Ilya Sutskever, Jeff Wu, and OpenAI.	Eisenschlos, Daniel Gillick, Jacob Eisenstein, and	636
584	2023. Weak-to-strong generalization: Eliciting	William W Cohen. 2022. Time-aware language mod-	637
585	strong capabilities with weak supervision. <i>ArXiv</i> ,	els as temporal knowledge bases. <i>Transactions of the</i>	638
586	abs/2312.09390.	<i>Association for Computational Linguistics</i> , 10:257–	639
587	Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun.	273.	640
588	2023. Instruction mining: Instruction data selection	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi	641
589	for tuning large language models. <i>arXiv preprint</i>	Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun,	642
590	<i>arXiv:2307.06290.</i>	and Bowen Zhou. 2023. Enhancing chat language	643
591	Moses Charikar, Chirag Pabbaraju, and Kirankumar	models by scaling high-quality instructional conver-	644
592	Shiragur. 2024. Quantifying the gain in weak-to-	sations. <i>ArXiv</i> , abs/2305.14233.	645
593	strong generalization. <i>ArXiv</i> , abs/2405.15116.	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	646
		Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	647
		Akhil Mathur, Alan Schelten, Amy Yang, Angela	648
		Fan, et al. 2024. The llama 3 herd of models. <i>arXiv</i>	649
		<i>preprint arXiv:2407.21783.</i>	650

651	Yann Dubois, Bal'azs Galambosi, Percy Liang, and Tatsunori Hashimoto. 2024. Length-controlled alpaca-eval: A simple way to debias automatic evaluators. <i>ArXiv</i> , abs/2404.04475.	
652		
653		
654		
655	Shangbin Feng, Wenxuan Ding, Alisa Liu, Zifeng Wang, Weijia Shi, Yike Wang, Zejiang Shen, Xiaochuang Han, Hunter Lang, Chen-Yu Lee, Tomas Pfister, Yejin Choi, and Yulia Tsvetkov. 2025. When one llm drools, multi-llm collaboration rules.	
656		
657		
658		
659		
660	Víctor Gallego. 2024. Refined direct preference optimization with synthetic data for behavioral alignment of llms. <i>arXiv preprint arXiv:2402.08005</i> .	
661		
662		
663	Chujie Gao, Qihui Zhang, Dongping Chen, Yue Huang, Siyuan Wu, Zhengyan Fu, Yao Wan, Xiangliang Zhang, and Lichao Sun. 2024. The best of both worlds: Toward an honest and helpful large language model. <i>arXiv preprint arXiv:2406.00380</i> .	
664		
665		
666		
667		
668	Yi Gui, Yao Wan, Zhen Li, Zhongyi Zhang, Dongping Chen, Hongyu Zhang, Yi Su, Bohua Chen, Xing Zhou, Wenbin Jiang, et al. 2025. Uicopilot: Automating ui synthesis via hierarchical code generation from webpage designs. In <i>THE WEB CONFERENCE 2025</i> .	
669		
670		
671		
672		
673		
674	Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. <i>arXiv preprint arXiv:2501.12948</i> .	
675		
676		
677		
678		
679	Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual programming: Compositional visual reasoning without training. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14953–14962.	
680		
681		
682		
683		
684	Thorsten Händler. 2023. Balancing autonomy and alignment: A multi-dimensional taxonomy for autonomous llm-powered multi-agent architectures. <i>ArXiv</i> , abs/2310.03659.	
685		
686		
687		
688	Jiwoo Hong, Noah Lee, and James Thorne. 2024. ORPO: Monolithic preference optimization without reference model. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 11170–11189, Miami, Florida, USA. Association for Computational Linguistics.	
689		
690		
691		
692		
693		
694	Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. 2023. Metagpt: Meta programming for multi-agent collaborative framework. <i>arXiv preprint arXiv:2308.00352</i> .	
695		
696		
697		
698		
699	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	
700		
701		
702		
703		
	Yue Huang, Siyuan Wu, Chujie Gao, Dongping Chen, Qihui Zhang, Yao Wan, Tianyi Zhou, Xiangliang Zhang, Jianfeng Gao, Chaowei Xiao, et al. 2024. Datagen: Unified synthetic dataset generation via large language models.	704 705 706 707 708
	Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. 2024. Openai o1 system card. <i>arXiv preprint arXiv:2412.16720</i> .	709 710 711 712 713
	Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. 2023. Follow-bench: A multi-level fine-grained constraints following benchmark for large language models. <i>arXiv preprint arXiv:2310.20410</i> .	714 715 716 717 718 719
	Jungo Kasai, Keisuke Sakaguchi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, Kentaro Inui, et al. 2023. Realtime qa: What's the answer right now? <i>Advances in neural information processing systems</i> , 36:49025–49043.	720 721 722 723 724
	Jingun Kwon, Hidetaka Kamigaito, Manabu Okumura, et al. 2024. Instructcmp: Length control in sentence compression through instruction-based large language models. <i>arXiv preprint arXiv:2406.11097</i> .	725 726 727 728
	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. 2024. Rewardbench: Evaluating reward models for language modeling. <i>arXiv preprint arXiv:2403.13787</i> .	729 730 731 732 733 734
	Angeliki Lazaridou, Adhi Kuncoro, Elena Gribovskaya, Devang Agrawal, Adam Liska, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d'Autume, Tomas Kocisky, Sebastian Ruder, et al. 2021. Mind the gap: Assessing temporal generalization in neural language models. <i>Advances in Neural Information Processing Systems</i> , 34:29348–29363.	735 736 737 738 739 740 741
	Changho Lee, Janghoon Han, Seonghyeon Ye, Stanley Jungkyu Choi, Honglak Lee, and Kyunghoon Bae. 2024. Instruction matters, a simple yet effective task selection approach in instruction tuning for specific tasks. <i>arXiv preprint arXiv:2404.16418</i> .	742 743 744 745 746
	Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023a. Camel: Communicative agents for "mind" exploration of large language model society. <i>Advances in Neural Information Processing Systems</i> , 36:51991–52008.	747 748 749 750 751
	Haoran Li, Qingxiu Dong, Zhengyang Tang, Chaojun Wang, Xingxing Zhang, Haoyang Huang, Shaohan Huang, Xiaolong Huang, Zeqiang Huang, Dongdong Zhang, Yuxian Gu, Xin Cheng, Xun Wang, Si-Qing Chen, Li Dong, Wei Lu, Zhifang Sui, Benyou Wang, Wai Lam, and Furu Wei. 2024a. Synthetic data (almost) from scratch: Generalized instruction tuning for language models. <i>ArXiv</i> , abs/2402.13064.	752 753 754 755 756 757 758 759

760	Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024b. Superfiltering: Weak-to-strong data filtering for fast instruction-tuning . In <i>Annual Meeting of the Association for Computational Linguistics</i> .	Zhou. 2023. # instag: Instruction tagging for analyzing supervised fine-tuning of large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	814 815 816 817
765	Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2023b. From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning. <i>arXiv preprint arXiv:2308.12032</i> .	Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. 2023. Wizardcoder: Empowering code large language models with evol-instruct . <i>ArXiv</i> , abs/2306.08568.	818 819 820 821 822
770	Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gonzalez, and Ion Stoica. 2024c. From crowdsourced data to high-quality benchmarks: Arena-hard and benchbuilder pipeline. <i>arXiv preprint arXiv:2406.11939</i> .	Yu Meng, Mengzhou Xia, and Danqi Chen. 2025. Simpo: Simple preference optimization with a reference-free reward. <i>Advances in Neural Information Processing Systems</i> , 37:124198–124235.	823 824 825 826
775	Zhuang Li, Yuncheng Hua, Thuy-Trang Vu, Haolan Zhan, Lizhen Qu, and Gholamreza Haffari. 2024d. Scar: Efficient instruction-tuning for large language models via style consistency-aware response ranking . <i>ArXiv</i> , abs/2406.10882.	Rudra Murthy, Prince Kumar, Praveen Venkateswaran, and Danish Contractor. 2024. Evaluating the instruction-following abilities of language models using knowledge tasks. <i>arXiv preprint arXiv:2410.12972</i> .	827 828 829 830 831
780	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. <i>arXiv preprint arXiv:2211.09110</i> .	Jinjie Ni, Yifan Song, Deepanway Ghosal, Bo Li, David Junhao Zhang, Xiang Yue, Fuzhao Xue, Zian Andy Zheng, Kaichen Zhang, Mahir Shah, Kabir Jain, Yang You, and Michael Shieh. 2024a. Mixeval-x: Any-to-any evaluations from real-world data mixtures . <i>ArXiv</i> , abs/2410.13754.	832 833 834 835 836 837
785	Hunter Lightman, Vineet Kosaraju, Yura Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step . <i>ArXiv</i> , abs/2305.20050.	Jinjie Ni, Fuzhao Xue, Xiang Yue, Yuntian Deng, Mahir Shah, Kabir Jain, Graham Neubig, and Yang You. 2024b. Mixeval: Deriving wisdom of the crowd from llm benchmark mixtures. <i>arXiv preprint arXiv:2406.06565</i> .	838 839 840 841 842
790	Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. 2024a. Deepseek-v3 technical report. <i>arXiv preprint arXiv:2412.19437</i> .	OpenAI. 2024. Hello gpt-4o . Accessed: 2024-06-06.	843
795	Chris Yuhao Liu and Liang Zeng. 2024. Skywork reward model series. https://huggingface.co/Skywork . Hugging Face model repository.	Geliang Ouyang, Jingyao Chen, Zhihe Nie, Yi Gui, Yao Wan, Hongyu Zhang, and Dongping Chen. 2025. nvagent: Automated data visualization from natural language via collaborative agent workflow. <i>arXiv preprint arXiv:2502.05036</i> .	844 845 846 847 848
798	Liangxin Liu, Xuebo Liu, Derek F Wong, Dongfang Li, Ziyi Wang, Baotian Hu, and Min Zhang. 2024b. Selectit: Selective instruction tuning for large language models via uncertainty-aware self-reflection. <i>arXiv preprint arXiv:2402.16705</i> .	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in neural information processing systems</i> , 35:27730–27744.	849 850 851 852 853 854
803	Liangxin Liu, Xuebo Liu, Derek F. Wong, Dongfang Li, Ziyi Wang, Baotian Hu, and Min Zhang. 2024c. Selectit: Selective instruction tuning for large language models via uncertainty-aware self-reflection . <i>ArXiv</i> , abs/2402.16705.	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4 . <i>ArXiv</i> , abs/2304.03277.	855 856 857
808	Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning . <i>ArXiv</i> , abs/2312.15685.	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. <i>Advances in Neural Information Processing Systems</i> , 36:53728–53741.	858 859 860 861 862 863
812	Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren	Patrick Schober, Christa Boer, and Lothar A Schwarte. 2018. Correlation coefficients: appropriate use and interpretation. <i>Anesthesia & analgesia</i> , 126(5):1763–1768.	864 865 866 867

868	Guijin Son, Hyunwoo Ko, Hoyoung Lee, Yewon Kim, and Seunghyeok Hong. 2024. Llm-as-a-judge & reward model: What they can and cannot do. <i>arXiv preprint arXiv:2409.11239</i> .	924
869		925
870		926
871		927
872	Tian-Xiang Sun, Xiang-Yang Liu, Xi-Peng Qiu, and Xuan-Jing Huang. 2022. Paradigm shift in natural language processing. <i>Machine Intelligence Research</i> , 19(3):169–183.	928
873		929
874		
875		
876	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. <i>arXiv preprint arXiv:2408.00118</i> .	
877		
878		
879		
880		
881		
882	Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. 2024a. Interpretable preferences via multi-objective reward modeling and mixture-of-experts. <i>arXiv preprint arXiv:2406.12845</i> .	
883		
884		
885		
886	Jiahao Wang, Bolin Zhang, Qianlong Du, Jiajun Zhang, and Dianhui Chu. 2024b. A survey on data selection for llm instruction tuning. <i>arXiv preprint arXiv:2402.05123</i> .	
887		
888		
889		
890	Junlin Wang, Jue Wang, Ben Athiwaratkun, Ce Zhang, and James Zou. 2024c. Mixture-of-agents enhances large language model capabilities. <i>arXiv preprint arXiv:2406.04692</i> .	
891		
892		
893		
894	Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. 2024d. Diversity measurement and subset selection for instruction tuning datasets. <i>arXiv preprint arXiv:2402.02318</i> .	
895		
896		
897		
898		
899	Shuohang Wang, Yang Liu, Yichong Xu, Chenguang Zhu, and Michael Zeng. 2021. Want to reduce labeling cost? gpt-3 can help . <i>ArXiv</i> , abs/2108.13487.	
900		
901		
902	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023. How far can camels go? exploring the state of instruction tuning on open resources . <i>ArXiv</i> , abs/2306.04751.	
903		
904		
905		
906		
907		
908	Zifeng Wang, Chun-Liang Li, Vincent Perot, Long T. Le, Jin Miao, Zizhao Zhang, Chen-Yu Lee, and Tomas Pfister. 2024e. Codeclm: Aligning language models with tailored synthetic data . In <i>NAACL-HLT</i> .	
909		
910		
911		
912	Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu, Wendy Gao, Jiaxin Xu, et al. 2024. Benchmarking complex instruction-following with multiple constraints composition. <i>arXiv preprint arXiv:2407.03978</i> .	
913		
914		
915		
916		
917		
918	Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Auto-gen: Enabling next-gen llm applications via multi-agent conversation framework. <i>arXiv preprint arXiv:2308.08155</i> .	
919		
920		
921		
922		
923		
	Siyuan Wu, Yue Huang, Chujie Gao, Dongping Chen, Qihui Zhang, Yao Wan, Tianyi Zhou, Xiangliang Zhang, Jianfeng Gao, Chaowei Xiao, et al. 2024a. Unigen: A unified framework for textual dataset generation using large language models. <i>arXiv preprint arXiv:2406.18966</i> .	924
		925
		926
		927
		928
		929
	Yang Wu, Huayi Zhang, Yizheng Jiao, Lin Ma, Xiaozhong Liu, Jinhong Yu, Dongyu Zhang, Dezhi Yu, and Wei Xu. 2024b. Rose: A reward-oriented data selection framework for llm task-specific instruction tuning. <i>arXiv preprint arXiv:2412.00631</i> .	930
		931
		932
		933
		934
	Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: Selecting influential data for targeted instruction tuning . <i>ArXiv</i> , abs/2402.04333.	935
		936
		937
		938
	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions . <i>ArXiv</i> , abs/2304.12244.	939
		940
		941
		942
		943
	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024a. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing . <i>ArXiv</i> , abs/2406.08464.	944
		945
		946
		947
		948
	Zhangchen Xu, Fengqing Jiang, Luyao Niu, Bill Yuchen Lin, and Radha Poovendran. 2024b. Stronger models are not stronger teachers for instruction tuning. <i>arXiv preprint arXiv:2411.07133</i> .	949
		950
		951
		952
	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jiahong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Ke-Yang Chen, Kexin Yang, Mei Li, Min Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yunyang Wan, Yunfei Chu, Zeyu Cui, Zhenru Zhang, and Zhi-Wei Fan. 2024a. Qwen2 technical report . <i>ArXiv</i> , abs/2407.10671.	953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
	An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024b. Qwen2. 5 technical report . <i>arXiv preprint arXiv:2412.15115</i> .	969
		970
		971
		972
	Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer, Chao Huang, Pin-Yu Chen, et al. 2024. Justice or prejudice? quantifying biases in llm-as-a-judge . <i>arXiv preprint arXiv:2410.02736</i> .	973
		974
		975
		976
		977
	Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2023. Large language model as	978
		979
		980

981 attributed training data generator: A tale of diversity
982 and bias. *Advances in Neural Information Processing*
983 *Systems*, 36:55734–55784.

984 Jieyu Zhang, Le Xue, Linxin Song, Jun Wang, Weikai
985 Huang, Manli Shu, An Yan, Zixian Ma, Juan Carlos
986 Niebles, Silvio Savarese, Caiming Xiong, Zeyuan
987 Chen, Ranjay Krishna, and Ran Xu. 2024. [Provi-](#)
988 [sion: Programmatically scaling vision-centric instruc-](#)
989 [tion data for multimodal language models](#). *ArXiv*,
990 [abs/2412.07012](#).

991 Wenting Zhao, Xiang Ren, John Frederick Hessel,
992 Claire Cardie, Yejin Choi, and Yuntian Deng. 2024.
993 [Wildchat: 1m chatgpt interaction logs in the wild](#).
994 *ArXiv*, [abs/2405.01470](#).

995 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle
996 Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
997 Zhuohan Li, Zi Lin, Eric P Xing, et al. 2023a. Lmsys-
998 chat-1m: A large-scale real-world llm conversation
999 dataset. *arXiv preprint arXiv:2309.11998*.

1000 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
1001 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
1002 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023b.
1003 Judging llm-as-a-judge with mt-bench and chatbot
1004 arena. *Advances in Neural Information Processing*
1005 *Systems*, 36:46595–46623.

1006 Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan
1007 Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma.
1008 2024. [Llamafactory: Unified efficient fine-tuning](#)
1009 [of 100+ language models](#). In *Proceedings of the*
1010 *62nd Annual Meeting of the Association for Compu-*
1011 *tational Linguistics (Volume 3: System Demonstra-*
1012 *tions)*, Bangkok, Thailand. Association for Computa-
1013 tional Linguistics.

1014 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer,
1015 Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping
1016 Yu, Lili Yu, et al. 2024a. Lima: Less is more for
1017 alignment. *Advances in Neural Information Process-*
1018 *ing Systems*, 36.

1019 Huichi Zhou, Zhaoyang Wang, Hongtao Wang, Dong-
1020 ping Chen, Wenhan Mu, and Fangyuan Zhang. 2024b.
1021 [Evaluating the validity of word-level adversarial at-](#)
1022 [tacks with large language models](#). In *Annual Meeting*
1023 *of the Association for Computational Linguistics*.

A Detailed Related Works

Instruction Tuning Data Selection. While Large Language Models (LLMs) like GPT-4 (Achiam et al., 2023; OpenAI, 2024) and Llama-3 (Dubey et al., 2024) excel in natural language understanding and generation, their pretraining objectives often misalign with user goals for instruction-following tasks (Murthy et al., 2024; Gao et al., 2024; Wen et al., 2024). Instruction tuning (or supervised fine-tuning) addresses this gap by refining LLMs on curated datasets of prompts and responses. Recent efforts like Vicuna (Peng et al., 2023) and LIMA (Zhou et al., 2024a) demonstrate high performance with a carefully selected small dataset, highlighting the growing importance of efficient instruction tuning and paving the way for aligning models with selected samples. This involves determining which instruction-response pairs to include in the training dataset and how to sample them effectively (Albalak et al., 2024).

Three key metrics determine instruction data quality: *Difficulty*, *Quality*, and *Diversity*. *Difficulty*, focusing mainly on the question side, is considered more valuable for model learning (Liu et al., 2024b; Lee et al., 2024; Wang et al., 2024b). IFD (Li et al., 2023b) pioneered the measurement of instruction-following difficulty for specific pairs, later enhanced by utilizing GPT-2 for efficient estimation in a weak-to-strong manner (Li et al., 2024b). *Quality*, mainly addressing the response side, measures the helpfulness and safety of model responses, typically assessed using LLM evaluators (Chen et al., 2023a, 2024b; Liu et al., 2024c; Ye et al., 2024), reward models (Son et al., 2024; Lambert et al., 2024), and gradient similarity search (Xia et al., 2024). *Diversity*, spanning both instruction and response aspects, plays a crucial role in covering various instruction formats and world knowledge, primarily improving model robustness (Bukharin and Zhao, 2023; Wang et al., 2024d). Our work stands out by addressing all three key components in data selection, introducing novel approaches to measuring difficulty from multiple LLMs’ responses and ultimately enhancing model performance.

Data Synthesis for Instruction Tuning. While the development of LLMs initially relied on human-curated instruction datasets for instruction tuning (Zheng et al., 2023a; Zhao et al., 2024; Lightman et al., 2023), this approach proved time-consuming and labor-intensive, particularly as the complex-

ity and scope of target tasks increased (Demrozi et al., 2023; Wang et al., 2021). Consequently, researchers began exploring the use of frontier LLMs to generate synthetic instruction datasets, aiming to both address these scalability challenges (Ding et al., 2023; Chen et al., 2023b, 2024d) and leverage models’ advanced capabilities in developing next-generation foundation models (Burns et al., 2023; Li et al., 2024b; Charikar et al., 2024). Early approaches (Xu et al., 2023; Wang et al., 2024e; Zhou et al., 2024b; Luo et al., 2023) focused on leveraging LLMs to generate synthetic instructions through a subset of human-annotated seed instructions (Chen et al., 2023a; Wang et al., 2023), and further enhanced by few-shot (Li et al., 2024a) and attribute-guided prompting (Yu et al., 2023; Wu et al., 2024a; Huang et al., 2024). A parallel line of research explored summarizing world knowledge to create more diverse synthetic datasets, aiming to maximize the coverage of different domains and task types (Cui et al., 2023; Li et al., 2024a). Recent advancements have further streamlined this process by utilizing instructions directly from pre-trained LLMs with simple prompt templates (Xu et al., 2024a; Chen et al., 2024c; Zhang et al., 2024), significantly reducing the required custom design from human effort. While existing work has primarily focused on generating extensive, diverse, and high-quality datasets—often scaling to 100,000 examples or more—this approach introduces challenges in terms of computational efficiency and training resource requirements (Li et al., 2024d; Dubois et al., 2024).

Deriving Crowded Wisdom from Multi-LLM. Single LLM’s response to a question face limitations in its representation of data (particularly cutting-edge knowledge) (Lazaridou et al., 2021; Dhingra et al., 2022; Kasai et al., 2023), skills (as no single LLM is universally optimal *empirically*) (Sun et al., 2022; Liang et al., 2022; Chen et al., 2024a), and diverse perspectives (Feng et al., 2025). Previous work has demonstrated that *online* multi-LLM wisdom (also known as compositional agent frameworks (Gupta and Kembhavi, 2023)) tends to outperform single models across various domains, providing more comprehensive and reflective solution on complex downstream tasks (Wang et al., 2024c; Hong et al., 2023; Wu et al., 2023; Li et al., 2023a; Ouyang et al., 2025; Gui et al., 2025). *Offline* crowded wisdom, where data are pre-collected rather than real-time inference, also show poten-

tial in model alignment (Gallego, 2024; Rafailov et al., 2023; Meng et al., 2025) and benchmark construction (Ni et al., 2024b,b). In this paper, we pioneer the use of *offline* multi-LLM wisdom for instruction data selection by utilizing these LLMs’ responses and their reward Score as *reflections* to measure instruction-response pairs’ *Difficulty* and *Quality*.

B Detailed Experiment Setups

B.1 Models & Benchmarks & Datasets Introduction

Models. In our study, the synthetic instruction dataset used for data selection consists of 19 response generators across 6 model families. These families include Qwen2 (Yang et al., 2024a), Qwen2.5 (Yang et al., 2024b), LLaMA 3 (Dubey et al., 2024), LLaMA 3.1 (Dubey et al., 2024), Gemma 2 (Team et al., 2024), and Phi-3 (Abdin et al., 2024). In our experiments, we perform supervised fine-tuning on the LLaMA3.2-3B-base/instruct (Dubey et al., 2024) and Qwen-2.5-3b-base/instruct (Yang et al., 2024b) models using the selected 1K datasets. A comprehensive overview of the models used in our study is presented in Table 5.

Benchmarks. In order to evaluate the instruction-following capabilities of the models, we use two widely-used instruction-following benchmarks: MT-Bench(Zheng et al., 2023b) and Arena-Hard (Li et al., 2024c) in our study.

MT-Bench (Zheng et al., 2023b). MT-bench is a collection of open-ended questions designed to evaluate a chatbot’s performance in multi-turn conversations and its ability to follow instructions—two critical factors in aligning with human preferences. It consists of 80 high-quality multi-turn questions, which are divided into 8 categories: writing, roleplay, extraction, reasoning, mathematics, coding, knowledge I (STEM), and knowledge II (humanities/social sciences). Each category contains 10 questions. This framework provides a robust tool for assessing the practical effectiveness of LLMs and their alignment with human preferences, through meticulously designed questions and evaluations conducted by human annotators.

Arena-Hard (Li et al., 2024c). Arena-Hard is a benchmark consisting 500 challenging prompts curated by BenchBuilder. It extracts high-quality

prompts from crowdsourced datasets like Chatbot Arena (Zheng et al., 2023b) and WildChat-1M (Zhao et al., 2024) without human intervention. The prompts are Scored and filtered based on seven key qualities, including specificity, domain knowledge, complexity, problem-solving, creativity, technical accuracy, and real-world applicability. This ensures that the prompts are challenging and capable of distinguishing between models. Unlike static benchmarks, Arena-Hard can be continuously updated to reflect the latest advancements in LLMs, avoiding the risk of becoming obsolete or leaking test data.

Datasets. In this paper, we conduct our experiments on Magpie-100K-Generator-Zoo(Xu et al., 2024b) because it provides a sufficiently large quantity of high-quality instruction fine-tuning data. It is a subset sampled from the MagpieAir-3M (Xu et al., 2024a) dataset, a large-scale instruction dataset. Magpie-100K contains 100,000 high-quality instructions, which are categorized into several types, including information seeking, mathematics, planning, coding and debugging, advice seeking, creative writing, reasoning, data analysis, brainstorming, editing, role-playing, and more. Each instruction has responses from 19 models across 6 model families—and their reward scores form 3 reward models. The diversity of these instructions ensures that the dataset covers a wide range of scenarios and tasks, making it suitable for instruction tuning of large language models (LLMs).

B.2 Model Training Details

Table 2 demonstrates the detailed supervised fine-tuning (SFT) hyper-parameters. We perform experiments on a server with eight NVIDIA A800-SXM4-80GB GPUs, two Intel Xeon Platinum 8358P 64-Core Processor, and 1024 GB of RAM. These experiments were conducted using LLaMA-Factory(Zheng et al., 2024).

B.3 Baseline Introduction

In this section, we present five baseline methods for comparison in our study. For each baseline, we describe its implementation details and rationale for inclusion.

Length-Based Filtering (Kwon et al., 2024). The Length method filters instructions based on their token count. We use the LLaMA 3.2 3B Instruction tokenizer to compute the number of to-

Table 5: Overview of 22 models used in our study.

Model Family	Release Date	Model ID	Size
Qwen2 (Yang et al., 2024a)	Jun, 2024	Qwen2-1.5B-Instruct	1.5B
		Qwen2-7B-Instruct	7B
		Qwen2-72B-Instruct	72B
Qwen2.5 (Yang et al., 2024b)	Sept, 2024	Qwen2.5-3B	3B
		Qwen2.5-3B-Instruct	3B
		Qwen2.5-7B-Instruct	7B
		Qwen2.5-14B-Instruct	14B
		Qwen2.5-32B-Instruct	32B
Llama 3 (Dubey et al., 2024)	Apr, 2024	Llama-3-8B-Instruct	8B
		Llama-3-70B-Instruct	70B
Llama 3.1 (Dubey et al., 2024)	Jul, 2024	Llama-3.1-8B-Instruct	8B
		Llama-3.1-70B-Instruct	70B
		Llama-3.1-405B-Instruct	405B
Llama 3.2 (Dubey et al., 2024)	Jul, 2024	Llama-3.2-3B	3B
		Llama-3.2-3B-Instruct	3B
Gemma 2 (Team et al., 2024)	Jun, 2024	Gemma-2-2B-it	2B
		Gemma-2-9B-it	9B
		Gemma-2-27B-it	27B
Phi-3 (Abdin et al., 2024)	Jun, 2024	Phi-3-mini-128k-instruct	3.5B
		Phi-3-small-128k-instruct	7B
		Phi-3-medium-128k-instruct	14B

Table 6: This table includes the hyper-parameters for supervised fine-tuning.

Hyper-parameter	Value
Learning Rate	1×10^{-5}
Number of Epochs	3
Per-device Batch Size	1
Gradient Accumulation Steps	2
Optimizer	Adamw
Learning Rate Scheduler	cosine
Warmup Steps	150
Max Sequence Length	2048

1223 kens in each instruction. Instructions that meet the
1224 predefined length criteria are selected for further
1225 processing.

1226 **Instag-Based Selection (Lu et al., 2023).** The
1227 Instag method introduces instruction tagging to ana-
1228 lyze the supervised fine-tuning process of large lan-
1229 guage models. Our implementation follows these
1230 steps:

1231 We use DeepSeek’s API to obtain the true labels
1232 of instructions. Instructions are grouped based on
1233 their assigned labels. The complexity and diversity
1234 of each group are computed. Finally, we select a

subset of instructions that exhibit the most desirable
characteristics.

Direct Score Filtering The Direct Score method
is inspired by the work of Chen et al. (2023) (Chen
et al., 2023a), which proposes a scoring mecha-
nism for instruction selection. We implement this
approach as follows:

We use the same prompt templates as the original
paper. Instead of the original scoring model, we
use DeepSeek for scoring, ensuring consistency
with our other experimental setups. We select the
top 1,000 instructions based on their scores.

Instruction Filtering by IFD The Instruction
Filtering by IFD method is based on the work of
Li et al. (2023) (Li et al., 2023b), which introduces
self-guided data selection to improve instruction
tuning. We directly use the open-source implemen-
tation from Cherry LLM and apply the following
three-step process:

Train a Pre-Experienced Model to establish prior
knowledge. Compute IFD (Instruction Filtering De-
gree) using the Pre-Experienced Model. Filter the
dataset based on IFD scores to retain high-quality
instructions. To evaluate the effectiveness of IFD,
we implement two versions:

1260 IFD (with pre): Uses a trained Pre-Experienced
1261 Model to compute IFD.

1262 IFD (no pre): Computes IFD directly using the
1263 model to be trained.

1264 **Random Sampling** The Random baseline selects
1265 a random subset of X instructions. Additionally,
1266 for each instruction, we randomly select one of
1267 its 19 possible responses, ensuring that instruction-
1268 response pairs are fully randomized.

1269 C Additional Experiment Results

1270 C.1 Dataset Size Ablation Details

1271 Tables 7 and 8 details the training loss, evaluation
1272 loss, and scores of Llama3.2-3b-base/instruct fine-
1273 tuned on different dataset sizes when selected with
1274 the difficulty metric. The data clearly shows a rapid
1275 increase in accuracy in when increasing the dataset
1276 sizes up to 0.5k to 1k, and marginal increases af-
1277 terwards. This highlights the importance of data
1278 quality over sheer quantity in instruction tuning.

1279 C.2 CROWDSELECT Performance on LoRA

1280 Tables 9 and 10 details the performance of CROWD-
1281 SELECT and various baselines combined with
1282 LoRA finetuning. CROWDSELECT generally out-
1283 performs the baseline dataset selection methods
1284 on LoRA. However, more instability is found in
1285 LoRA training due to its limited learning capability
1286 compared with full finetuning.

1287 C.3 CROWDSELECT Performance on Full 1288 Finetuning

1289 Tables 11 and 12 details the performance of
1290 CROWDSELECT and various baselines combined
1291 with Full finetuning.

1292 C.4 Foundation Metric with Clustering 1293 Performance

1294 Table 13 details the performance of our foundation
1295 metric combined with clustering strategy.

1296 C.5 CROWDSELECT Integrated Metric 1297 Performance on Different Coefficient 1298 Combinations

1299 Tables 14, 15, and 16 details the performance of
1300 our Integrated metric performace on 9 sets of co-
1301 efficients. $w = (1, 1, 2)$ stands out as stable coeffi-
1302 cients among all other combinations.

C.6 CROWDSELECT Performance on 1303 Different Finetuning Methods 1304

1305 Table 17 details the performance of CROWDSE-
1306 LECT on SFT, DPO, SimPO, and ORPO(Hong
1307 et al., 2024). Data reveals consistent and stable
1308 performance our proposed metrics, while SimPO
1309 performs best on all scenarios.

C.7 CROWDSELECT Performance on 1310 Different Reward Models 1311

1312 Table 18 details the performance of CROWDSE-
1313 LECT on various reward models. Data reveals the
1314 importance of reward models on finetuned model
1315 performance. However, the strong points of differ-
1316 ent reward models is scattered. While the results
1317 show that reward models are crucial for effective
1318 fine-tuning, they also reveal a nuanced landscape
1319 where the strengths of different reward models are
1320 distributed across different aspects of performance.
1321 Scattered performance highlights the need for care-
1322 ful consideration when selecting a reward model
1323 and reflects that current Large Language reward
1324 models are of high variance. Further research into
1325 more robust reward models for Large Language
1326 models is therefore essential.

Table 7: Performance comparison of Llama-3b-instruct with different sizes of difficulty-based selected data.

Data Size	Train Loss	Eval. Loss	MT-Bench		Arena-Hard		
			Score	Avg. Tokens	Score	95% CI	Avg. Tokens
0.25k	0.418	0.951	6.850	301	81.9	(−1.2, 1.5)	275
0.5k	0.406	1.004	6.962	276	83.1	(−1.0, 1.1)	275
1k	0.407	0.942	6.887	271	82.6	(−1.5, 1.2)	273
2k	0.405	0.929	6.668	301	83.1	(−1.0, 1.4)	273
3k	0.415	0.871	6.625	304	85.1	(−1.3, 1.3)	276
4k	0.413	0.869	6.600	279	82.4	(−1.1, 1.7)	268
5k	0.415	0.867	6.675	295	83.3	(−0.7, 1.4)	272
6k	0.414	0.857	6.572	282	84.4	(−1.1, 1.3)	265
7k	0.413	0.848	6.743	286	84.1	(−0.9, 1.2)	266
8k	0.411	0.836	6.618	275	83.1	(−1.1, 1.6)	268
9k	0.411	0.822	6.681	274	83.3	(−1.3, 1.5)	269
10k	0.409	0.828	6.750	279	83.6	(−0.8, 1.7)	266

Table 8: Performance comparison of Llama-3b with different sizes of difficulty-based selected data.

Data Size	Train Loss	Eval. Loss	MT-Bench		Arena-Hard		
			Score	Avg. Tokens	Score	95% CI	Avg. Tokens
0.25k	0.567	1.138	4.731	492	75.0	(−1.1, 2.1)	289
0.5k	0.544	1.161	4.987	392	79.1	(−1.0, 1.7)	289
1k	0.539	1.123	5.200	325	78.1	(−1.4, 1.5)	289
2k	0.534	1.094	5.337	309	76.9	(−1.4, 2.2)	290
3k	0.537	1.046	5.237	286	80.0	(−1.6, 1.6)	289
4k	0.535	1.031	5.131	287	79.7	(−1.3, 1.5)	289
5k	0.534	1.022	4.987	271	81.5	(−1.0, 1.5)	289
6k	0.531	1.019	4.943	251	81.8	(−1.3, 1.5)	290
7k	0.529	1.004	4.825	218	78.5	(−1.2, 1.7)	289
8k	0.526	0.990	5.093	278	81.5	(−1.1, 1.3)	289
9k	0.519	0.982	4.893	245	83.2	(−1.5, 1.2)	289
10k	0.517	0.983	5.137	270	82.9	(−1.0, 1.1)	289

Table 9: Performance comparison of lora-version of Llama-3b-base/instruct and Qwen-3b-base/instruct models with different data selection strategies.

Benchmark	Base	Difficulty		Separability		Stability	
		↓	↑	↓	↑	↓	↑
Llama3.2-3b-instruct							
MT-Bench	6.200	6.456	6.688	6.100	6.725	6.131	6.866
Arena-Hard	74.4	69.6	76.8	69.4	72.9	69.8	74.6
Arena-Hard-95%CI	(-1.0, 1.5)	(-1.8,1.4)	(-1.5,1.9)	(-2.5,1.2)	(-1.6,1.5)	(-1.7,1.7)	(-1.7,2.0)
Llama3.2-3b-base							
MT-Bench	4.302	4.626	4.651	4.631	5.040	3.538	4.369
Arena-Hard	50.0	73.1	68.0	73.8	73.2	60.8	73.2
Arena-Hard-95%CI	(0.0,0.0)	(-1.8,1.6)	(-1.2,1.9)	(-1.2,1.8)	(-2.0,1.1)	(-1.7,1.2)	(-1.2,1.2)
Qwen3.2-3b-instruct							
MT-Bench	7.138	6.906	7.068	7.025	6.937	7.018	7.037
Arena-Hard	81.6	77.2	79.1	80.3	78.8	76.2	78.0
Arena-Hard-95%CI	(-1.8, 1.4)	(-1.9, 1.5)	(-2.1, 1.8)	(-1.9, 1.4)	(-1.2, 1.2)	(-1.7, 1.6)	(-1.8, 1.7)
Qwen3.2-3b							
MT-Bench	6.043	5.137	6.612	6.368	6.343	5.800	6.525
Arena-Hard	69.0	76.9	70.7	74.1	74.2	73.7	74.2
Arena-Hard-95%CI	(-2.2, 1.6)	(-2.0, 1.8)	(-1.8, 2.4)	(-1.8, 1.5)	(-2.1, 1.5)	(-2.0, 1.3)	(-1.8, 1.9)

Table 10: Performance comparison of lora-version of Llama-3b-base/instruct and Qwen-3b-base/instruct models with pre data selection strategies as baselines.

Benchmark	Random	Tags	Direct-Score		Length		IFD	
			↓	↑	↓	↑	no_pre	pre
Llama3.2-3b-instruct								
MT-Bench	6.325	6.610	6.631	6.406	6.087	5.375	6.706	6.768
Arena-Hard	74.2	80.1	80.0	74.8	78.1	67.5	81.2	79.5
Arena-Hard-95%CI	(-1.7, 1.3)	(-0.7, 0.7)	(-1.4, 1.7)	(-1.1, 1.8)	(-3.4, 2.1)	(-1.4, 0.9)	(-0.8, 1.5)	(-1.6, 1.8)
Llama3.2-3b-base								
MT-Bench	4.637	4.575	4.962	4.675	4.062	4.243	4.512	4.418
Arena-Hard	76.0	76.8	76.9	75.6	67.1	70.3	73.7	77.5
Arena-Hard-95%CI	(-2.0, 1.6)	(-1.6, 1.8)	(-1.8, 1.7)	(-1.6, 1.4)	(-2.0, 2.0)	(-2.3, 2.2)	(-1.5, 1.5)	(-1.8, 1.4)
Qwen2.5-3b-instruct								
MT-Bench	6.950	7.125	7.131	7.175	7.037	7.006	6.918	6.868
Arena-Hard	78.2	83.0	77.7	81.7	75.8	76.4	78.8	83.1
Arena-Hard-95%CI	(-1.5, 1.8)	(-1.7, 2.1)	(-1.6, 2.0)	(-1.7, 1.9)	(-2.0, 2.0)	(-1.4, 1.7)	(-1.3, 1.2)	(-0.8, 1.0)
Qwen2.5-3b-base								
MT-Bench	5.887	5.616	5.417	5.750	3.981	5.637	6.427	5.861
Arena-Hard	76.6	83.8	79.3	76.5	74.3	70.4	79.7	82.2
Arena-Hard-95%CI	(-1.7, 1.5)	(-1.3, 1.2)	(-1.8, 1.2)	(-2.0, 1.7)	(-1.8, 1.6)	(-1.6, 1.9)	(-1.3, 1.0)	(-1.3, 1.0)

Table 11: Performance comparison of fft-version of Llama-3b-base/instruct and Qwen-3b-base/instruct models with different data selection strategies.

Benchmark	Base	Difficulty		Separability		Stability	
		↓	↑	↓	↑	↓	↑
Llama3.2-3b-instruct							
MT-Bench	6.200	6.388	6.648	5.937	6.581	6.225	6.625
Arena-Hard	74.4	76.5	80.5	80.0	77.9	75.8	77.4
Arena-Hard-95%CI	(-1.0, 1.5)	(-1.6, 1.5)	(-0.9, 1.3)	(-1.3, 1.2)	(-1.5, 1.7)	(-1.3, 0.9)	(-1.5, 1.1)
Llama3.2-3b-base							
MT-Bench	4.302	4.506	4.738	4.731	5.056	4.675	5.088
Arena-Hard	50.0	78.6	76.8	81.8	83.3	80.0	78.3
Arena-Hard-95%CI	(0.0, 0.0)	(-1.9, 2.1)	(-1.6, 1.7)	(-1.8, 1.2)	(-1.8, 1.7)	(-1.5, 1.6)	(-1.6, 2.2)
Qwen2.5-3b-instruct							
MT-Bench	7.138	6.906	7.182	6.919	7.269	7.056	7.294
Arena-Hard	81.6	82.5	81.8	81.4	83.7	78.1	83.5
Arena-Hard-95%CI	(-1.8, 1.4)	(-1.8, 1.5)	(-1.6, 1.3)	(-1.7, 1.6)	(-1.4, 1.2)	(-1.2, 2.0)	(-1.4, 1.4)
Qwen2.5-3b-base							
MT-Bench	6.043	6.619	6.613	6.575	7.075	6.763	6.681
Arena-Hard	69.0	80.2	73.8	76.5	74.1	74.4	76.8
Arena-Hard-95%CI	(-2.2, 1.6)	(-1.7, 1.6)	(-2.5, 1.8)	(-1.8, 1.8)	(-1.6, 2.4)	(-1.5, 1.8)	(-1.8, 1.8)

Table 12: Performance comparison of fft-version of Llama-3b-base/instruct and Qwen-3b-base/instruct models with pre data selection strategies as baselines.

Benchmark	Random	Tags	Direct-Score		Length		IFD	
			↓	↑	↓	↑	no_pre	pre
Llama3.2-3b-instruct								
MT-Bench	6.356	6.393	6.068	6.050	5.612	5.781	6.593	6.243
Arena-Hard	74.8	81.6	76.9	77.6	72.9	75.0	76.8	78.4
Arena-Hard-95%CI	(-1.5, 1.6)	(-0.2, -0.2)	(-1.5, 2.0)	(-1.7, 1.9)	(-1.9, 1.9)	(-2.4, 2.0)	(-1.2, 1.6)	(-1.7, 1.5)
Llama3.2-3b-base								
MT-Bench	4.406	4.562	4.131	4.400	3.393	3.893	4.281	3.962
Arena-Hard	75.3	77.3	72.7	75.8	59.4	71.8	73.9	77.6
Arena-Hard-95%CI	(-2.0, 1.6)	(-1.1, 1.2)	(-2.4, 1.9)	(-1.4, 1.2)	(-1.1, 1.3)	(-1.0, 1.2)	(-1.0, 1.6)	(-1.6, 1.6)
Qwen2.5-3b-instruct								
MT-Bench	6.793	6.818	6.506	6.768	5.881	6.931	6.962	6.731
Arena-Hard	78.2	82.0	81.2	80.8	75.6	77.7	79.0	80.4
Arena-Hard-95%CI	(-1.7, 2.0)	(-2.4, 1.6)	(-1.5, 1.8)	(-2.1, 1.7)	(-1.0, 1.2)	(-1.7, 1.7)	(-1.0, 1.5)	(-1.3, 1.0)
Qwen2.5-3b-base								
MT-Bench	6.500	6.818	6.325	6.900	4.925	6.591	5.798	5.825
Arena-Hard	72.9	79.3	75.6	76.8	71.2	72.8	76.2	74.5
Arena-Hard-95%CI	(-2.2, 1.9)	(-2.2, 1.9)	(-1.6, 2.1)	(-1.9, 1.9)	(-1.7, 1.4)	(-2.3, 1.9)	(-1.4, 1.3)	(-1.5, 1.5)

Table 13: Performance comparison of cluster-chosen-data-fft-version of Llama-3b-base/instruct and Qwen-3b-base/instruct models with different data selection strategies.

Benchmark	Base	Random	Difficulty		Separability		Stability	
			↓	↑	↓	↑	↓	↑
Llama3.2-3b-instruct								
MT-Bench	6.200	6.743	6.256	6.675	6.094	6.619	6.275	6.913
Arena-Hard	74.4	80.9	81.4	82.6	84.8	81.9	80.0	81.8
Arena-Hard-95%CI	(-1.0, 1.5)	(-1.3, 1.4)	(-1.5, 2.0)	(-1.2, 1.8)	(-1.7, 1.4)	(-1.7, 1.7)	(-2.0, 2.2)	(-1.5, 1.7)
Llama3.2-3b-base								
MT-Bench	4.302	4.869	4.825	5.000	4.813	4.938	4.800	4.950
Arena-Hard	50.0	79.2	80.8	79.5	80.8	81.9	80.6	80.9
Arena-Hard-95%CI	(0.0, 0.0)	(-0.9, 0.9)	(-1.2, 1.7)	(-1.7, 2.2)	(-2.0, 1.6)	(-1.5, 2.1)	(-1.9, 1.8)	(-2.0, 1.6)
Qwen-3b-instruct								
MT-Bench	7.138	7.006	6.988	7.150	7.238	7.340	7.019	7.181
Arena-Hard	81.6	82.3	82.1	82.6	82.5	82.3	80.3	82.6
Arena-Hard-95%CI	(-1.8, 1.4)	(-1.0, 0.9)	(-1.6, 1.3)	(-1.9, 1.7)	(-2.1, 1.3)	(-1.0, 1.4)	(-1.5, 1.4)	(-1.4, 2.0)
Qwen-3b-base								
MT-Bench	6.043	7.162	6.575	6.800	6.856	6.875	6.819	6.869
Arena-Hard	69.0	74.6	78.2	78.5	78.0	75.7	73.6	76.9
Arena-Hard-95%CI	(-2.2, 1.6)	(-0.7, 1.0)	(-1.9, 2.4)	(-1.6, 1.7)	(-1.7, 1.8)	(-2.2, 2.1)	(-1.8, 1.8)	(-2.1, 1.6)

Table 14: Performance comparison of fft-version of Llama-3b-instruct on different coefficient combinations for multiple metrics with clustering.

Hyperparameter	Train Loss	Eval. Loss	MT-Bench		Arena-Hard		
			Score	Avg. Tokens	Score	95% CI	Avg. Tokens
1_1_1	0.312	0.715	6.913	307	81.8	(-0.5, 0.8)	266
1_-1_1	0.368	0.803	6.625	292	84.2	(-0.7, 1.0)	269
1_1_2	0.325	0.717	7.103	328	85.5	(-0.8, 1.1)	271
1_1_-1	0.294	0.617	6.650	298	82.7	(-1.5, 1.4)	278
1_1_1.5	0.338	0.721	6.850	312	84.7	(-1.6, 1.3)	266
1_-1_1.5	0.391	0.795	6.781	286	83.0	(-1.4, 1.4)	270
-1_-1_1	0.354	0.707	6.781	308	81.9	(-1.5, 1.3)	275
-1_-1_2	0.355	0.742	6.838	297	84.8	(-1.3, 1.2)	275
-1_-1_1.5	0.351	0.754	6.638	289	81.8	(-1.3, 1.3)	276

Table 15: Performance comparison of fft-version of Qwen-3b-instruct with different coefficient combinations for multiple metrics.

Hyperparameter	Train Loss	Eval. Loss	MT-Bench		Arena-Hard		
			Score	Avg. Tokens	Score	95% CI	Avg. Tokens
1_1_1	0.354	0.776	6.856	359	83.6	(-1.7, 1.2)	259
1_-1_1	0.432	0.861	7.138	383	81.6	(-1.4, 1.5)	259
1_1_2	0.371	0.776	7.131	366	85.2	(-1.2, 1.1)	262
1_1_-1	0.310	0.645	7.231	376	82.3	(-1.6, 1.5)	261
1_1_1.5	0.369	0.755	6.981	387	83.6	(-2.0, 1.2)	260
1_-1_1.5	0.430	0.872	7.371	390	82.4	(-1.7, 1.5)	260
-1_-1_1	0.431	0.874	7.025	397	81.9	(-1.1, 1.9)	260
-1_-1_2	0.431	0.888	6.963	377	80.6	(-1.8, 1.5)	259
-1_-1_1.5	0.433	0.869	6.956	377	82.4	(-1.8, 1.3)	260

Table 16: Performance comparison of fft-version of Llama-3b with different coefficient combinations for multiple metrics.

Hyperparameter	Train Loss	Eval. Loss	MT-Bench		Arena-Hard		
			Score	Avg. Tokens	Score	95% CI	Avg. Tokens
1_1_1	0.437	0.901	4.800	306	80.8	(-1.3, 1.6)	289
1_-1_1	0.497	1.007	5.019	319	80.3	(-2.2, 2.1)	290
1_1_2	0.454	0.904	4.613	282	82.1	(-1.8, 1.8)	290
1_1_-1	0.416	0.786	4.669	283	83.0	(-1.6, 2.0)	289
1_1_1.5	0.449	0.908	4.731	276	75.7	(-1.9, 2.4)	290
1_-1_1.5	0.496	1.016	5.125	309	80.6	(-2.4, 1.6)	290
-1_-1_1	0.469	0.973	5.050	307	80.7	(-1.8, 1.2)	289
-1_-1_2	0.469	0.968	4.719	268	81.6	(-1.2, 1.1)	290
-1_-1_1.5	0.469	0.968	4.588	291	80.0	(-2.0, 1.8)	290

Table 17: Performance comparison of Llama-3b-instruct models with different finetuning methods

Benchmark	Random	Difficulty		Separability		Stability	
		↓	↑	↓	↑	↓	↑
SFT							
MT-Bench	6.200	6.388	6.648	5.937	6.581	6.225	6.625
Arena-Hard	74.4	76.5	80.5	77.9	80.0	75.8	77.4
Arena-Hard-95%CI	(-1.0, 1.5)	(-1.6, 1.5)	(-0.9, 1.3)	(-1.5, 1.7)	(-1.3, 1.2)	(-1.3, 0.9)	(-1.5, 1.1)
DPO							
MT-Bench	6.463	6.431	6.768	6.431	6.418	6.256	6.818
Arena-Hard	74.2	75.1	77.3	76.1	78.5	73.2	76.2
Arena-Hard-95%CI	(-1.8, 1.6)	(-1.6, 1.6)	(-1.6, 1.7)	(-1.9, 1.9)	(-1.5, 1.4)	(-1.4, 1.3)	(-1.9, 1.5)
SimPO							
MT-Bench	6.950	6.425	7.137	6.518	7.043	6.675	6.931
Arena-Hard	78.7	78.0	78.8	78.2	79.7	76.0	75.5
Arena-Hard-95%CI	(-2.5, 2.0)	(-2.5, 3.1)	(-0.9, 1.2)	(-1.6, 0.8)	(-5.4, 6.5)	(-1.3, 1.1)	(-5.7, 6.2)
ORPO							
MT-Bench	6.412	6.450	6.450	6.525	6.431	6.312	6.400
Arena-Hard	73.7	73.2	73.7	73.3	74.6	73.2	75.6
Arena-Hard-95%CI	(-2.1, 2.2)	(-2.2, 1.8)	(-1.5, 2.0)	(-1.9, 1.8)	(-2.0, 2.2)	(-2.1, 2.2)	(-1.8, 2.2)

Table 18: Performance comparison of lora-version of Llama-3b-instruct models with different reward-models

Benchmark	Difficulty		Separability		Stability		Reward-Score	
	↓	↑	↓	↑	↓	↑	↓	↑
ArmoRM-Llama3-8B-v0.1								
MT-Bench	6.625	6.687	6.468	6.493	6.375	6.431	4.037	6.512
Arena-Hard	81.7	78.6	74.3	75.6	77.3	80.0	57.8	83.2
Arena-Hard-95%CI	(-2.0, 1.8)	(-1.8, 1.8)	(-1.8, 2.1)	(-2.0, 1.6)	(-1.8, 2.0)	(-1.0, 1.8)	(-2.0, 1.9)	(-1.5, 1.9)
Skywork-Reward-Llama-3.1-8B								
MT-Bench	6.456	6.688	6.100	6.725	6.131	6.866	4.012	6.675
Arena-Hard	69.6	76.8	69.4	72.9	69.8	74.6	52.6	77.4
Arena-Hard-95%CI	(-1.5,1.9)	(-1.8,1.4)	(-2.5,1.2)	(-1.6,1.5)	(-1.7,1.7)	(-1.7,2.0)	(-2.4, 2.0)	(-1.8, 2.1)
Skywork-Reward-Gemma-2-27B								
MT-Bench	6.512	6.593	6.756	6.881	6.637	6.756	3.793	6.943
Arena-Hard	76.2	78.2	75.4	80.2	79.7	83.6	56.1	79.6
Arena-Hard-95%CI	(-1.6, 2.0)	(-1.6, 1.5)	(-2.1, 2.1)	(-1.7, 2.4)	(-1.4, 1.4)	(-1.9, 2.0)	(-2.1, 2.1)	(-1.6, 1.7)

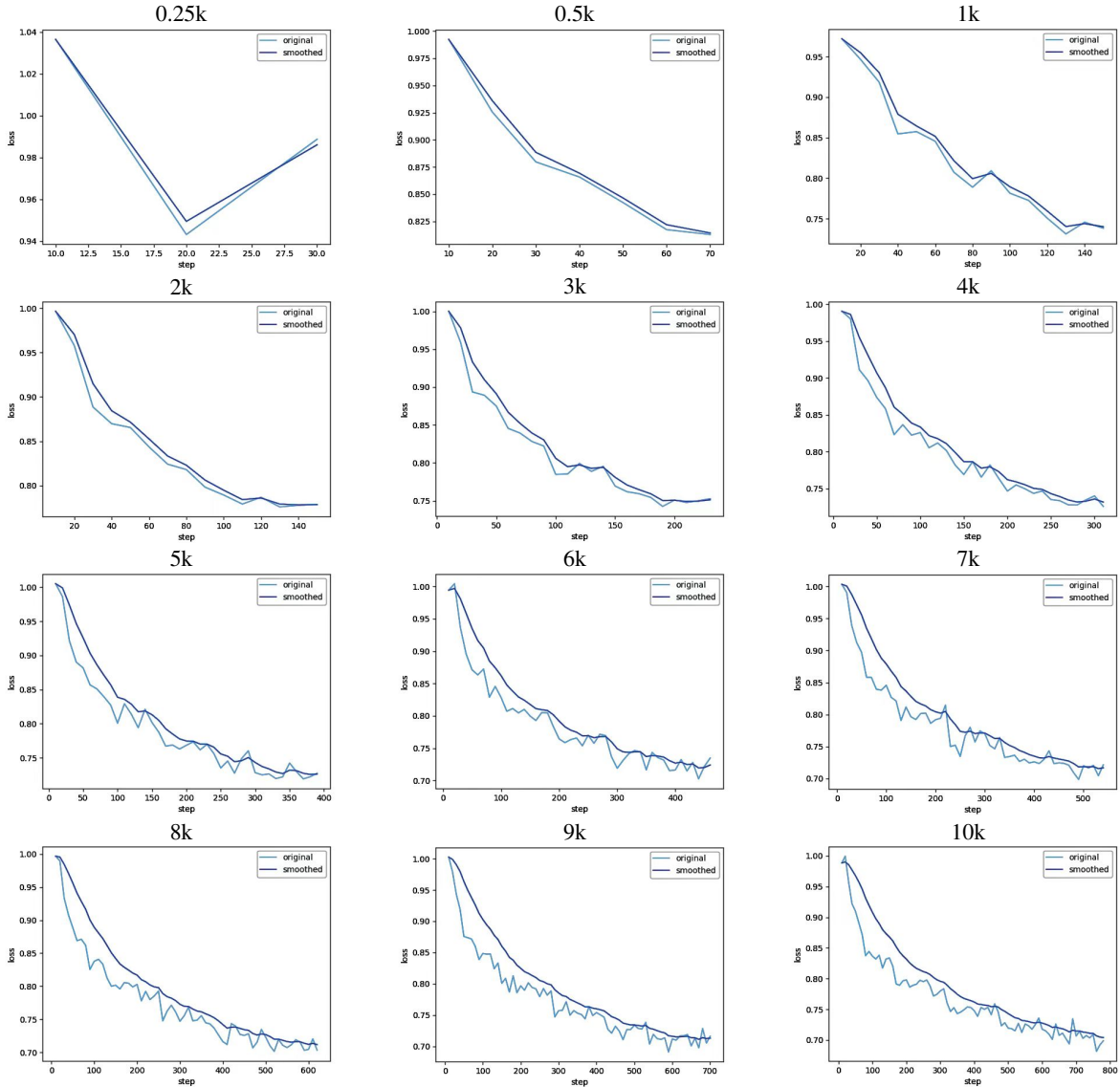


Figure 5: Lora train loss of training Llama-3b by using different sizes of randomly chosen data.