

The Art of Data Selection: A Survey on Data Selection for Fine-Tuning Large Language Models

Anonymous ACL submission

Abstract

Recently, Large Language Models (LLMs) have seen significant advancements, and supervised fine-tuning (SFT) plays a pivotal role in unleashing LLMs' potential to follow the users' instructions. As an emerging research field, data selection for fine-tuning LLMs aims to select a subset from a given candidate dataset for training selective-enhanced models to improve their performance and accelerate their training. Although some studies have already investigated these works, there is a lack of comprehensive analysis and comparison of them to provide potential research directions. To fill the gap, we first summarize a three-step scheme for data selection on existing works, including data preprocessing, data selector construction, and data selector evaluation, and comprehensively sort out the existing works according to this scheme. Then, we conduct an in-depth analysis of existing works from their efficiency and feasibility by making quantitative and qualitative comparisons and find that (1) the model-specific method who takes the loss output of the pending fine-tune model as an optimized goal is more effective; (2) increasing the complexity of the selector can improve the performance of the selective-enhanced model, but it needs more careful design to avoid introducing external factors. Finally, we summarize the trends in data selection and point out that the current main challenges are the lack of unified and efficient data quality measurement, as well as data selection for specific tasks and multiple turns of conversations.

1 Introduction

Large language models nowadays can generate natural and authentic human languages and complete many classic NLP challenges as well as real-world tasks (Naveed et al., 2023; Vaswani et al., 2023; Wang et al., 2022; Zhong et al., 2022). After the knowledge-based pretraining, the user-oriented supervised instruction fine-tuning endows LLMs

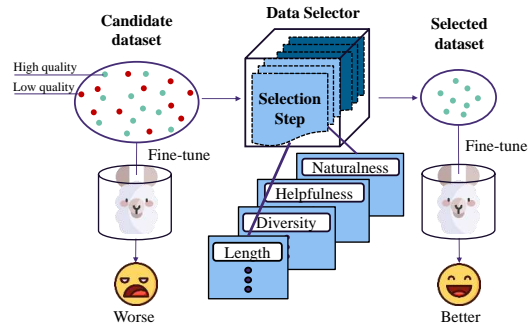


Figure 1: The Process of Data Selection.

with the most significant performance rise. With the success of LIMA (Zhou et al., 2023), data selection, that is, how to select a few high-quality samples from existing datasets to fine-tune better models in downstream tasks according to some prior indicators, has gradually become a research hotspot. It can improve fine-tuned LLMs' performance and accelerate their training simultaneously. Although recent works (Wang et al., 2024; Albalak et al., 2024) list most of the existing data selection methods for fine-tuning LLMs, there is a lack of in-depth analysis and comparison between each method for providing potential research directions.

Although recent works (Wang et al., 2024; Albalak et al., 2024) list most of the existing data selection methods for fine-tuning, there is a lack of in-depth analysis and comparison between each method.

To address these issues, we first summarize a three-step data selection scheme that can cover key parts of the entire data selection process, including data preprocessing, data selector construction, and data selector evaluation, after reviewing existing popular works. Then we conduct a comprehensive sort of the existing works based on the conversation format of the data to be selected, the indicator sources and calculation methods used by the selector, and the candidate datasets, models, and metrics used in the evaluation process.

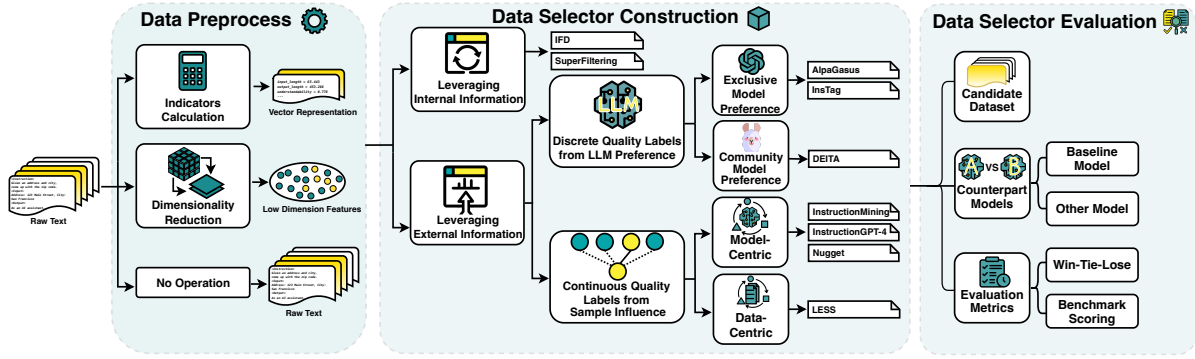


Figure 2: The Scheme of Data Selection for Fine-tuning LLMs.

072 Furthermore, we compared the existing work
 073 through quantitative and qualitative analysis.
 074 Specifically, we develop a unified efficiency mea-
 075 surement method based on the efficiency curve
 076 assumption to evaluate the performance of various
 077 models in selecting data, addressing hard compar-
 078 isons due to different experimental settings. We
 079 also qualitatively evaluate their feasibility by con-
 080 sidering simplicity and flexibility, including imple-
 081 mentation costs and reproducibility.

082 Through the analysis of the above two aspects,
 083 we not only obtained the technological develop-
 084 ment path of existing work, but also identified the
 085 following two findings that can help with data selec-
 086 tion in future work: (1) the selector taking the out-
 087 put of the pending fine-tune model as an optimized
 088 goal is more effective; (2) increasing the complex-
 089 ity of the selector can improve the performance of
 090 the selective-enhanced model, but it needs more
 091 careful design to avoid introducing external factors.
 092 Meanwhile, we point out that there are still many
 093 challenges to further research, including building
 094 unified and efficient data quality measurement and
 095 designing data selection for specific domains or
 096 multi-turn conversations.

097 2 Scheme of Data Selection

098 The data selection for fine-tuning LLM aims to
 099 select high-quality samples consisting of a subset
 100 from a given candidate dataset according to the
 101 data quality, resulting in the Selective-Enhanced
 102 Model (SEM) fine-tuned on the high-quality subset
 103 being better than the Baseline model (BM) trained
 104 on the full dataset.

105 Therefore, building a data selection method
 106 requires considering the following three aspects.
 107 Firstly, the form of data representation. For data
 108 selection, it is necessary to determine the selection

109 perspectives. In addition to considering the char-
 110 acteristics of the text itself, linguistic features or
 111 vector representation of the text are also commonly
 112 used. After the conversion of original samples, the
 113 key to building a data selector is determining the se-
 114 lection measurement, which includes two aspects:
 115 the source of data quality clues and how to obtain
 116 the quality label of a sample. On the one hand, it
 117 can be obtained by calculating the statistical charac-
 118 teristics of the data itself. On the other hand, it can
 119 also be obtained through external information such
 120 as third-party scorer models or by comparing the
 121 performance with known good samples. Finally, it
 122 is necessary to verify the usefulness of the data se-
 123 lection method after obtaining a subset selected by
 124 the selector. It can pair-wise compare the response
 125 from the basic model and the selective-enhanced
 126 model directly or compare their scoring in some
 127 popular benchmarks.

128 By considering the above factors, we construct a
 129 full process scheme of data selection after review-
 130 ing existing popular works, which is divided into
 131 three steps: (1) data preprocessing (Section 3), (2)
 132 data selector construction (Section 4), and (3) data
 133 selector evaluation (Section 5).

134 3 Data Preprocessing

135 Data preprocessing converts raw texts into fea-
 136 ture representations of the data for selection. Ac-
 137 cording to the converted forms, it can be divided
 138 into human-preferred **explicit features** (Cao et al.,
 139 2023; Wei et al., 2023), such as the length of input,
 140 model-oriented **implicit features** (Xia et al., 2024),
 141 such as low-dimensional gradients from LoRA (Hu
 142 et al., 2021) or purely **original texts** (Li et al., 2023;
 143 Chen et al., 2024) where the raw texts can preserve
 144 the most information.

145 **Explicit Features.** Some works usually take a

series of indicators as explicit features manually based on the human’s prior linguistic knowledge that can be extracted from the sample. For example, **InstructionMining** (Cao et al., 2023) converts each sample into a bag of NLP indicators (such as the length of input, the length of output, understandability, etc.). Utilizing such interpretable indicators of linguistic knowledge to represent the original sample can guide future selection with respect to these human-preferred aspects of a sample.

Implicit Features. Other works convert original texts to vector representations as implicit features of the sample. For example, **LESS** (Xia et al., 2024) randomly projects all candidate samples into low-dimensional gradient features with a warm-up LoRA, which can reflect the influence of each sample on the optimization process of the loss function. These implicit representations, though less interpretable to humans, are more objective features indicating the quality of the data.

4 Data Selector Construction

The choice of data quality measurement is the primary concern in constructing the data selector, which can be divided into two branches. One branch leverages internal information like statistical features from the candidate dataset (Li et al., 2024a). Another branch leverages the external information from LLM or datasets rather than from the SEM and candidate datasets. The leveraging external information works can be further divided into the group that uses discrete quality labels from LLM’s preference, which refers to the score given by the external models like ChatGPT (Chen et al., 2024; Lu et al., 2023; Liu et al., 2023), and the group uses continuous quality label from sample influence (Cao et al., 2023; Wei et al., 2023; Li et al., 2024b; Xia et al., 2024).

4.1 Leveraging Internal Information

Leveraging internal information means using only the features of a given candidate dataset as quality clues to determine whether to select them. The pioneering work (Li et al., 2023) proposes Instruction Following Difficulty (**IFD**) as a quantified metric, which can be obtained by using only a candidate dataset and a backbone pre-trained model. A higher IFD score indicates a closer relationship between the sample’s instruction and output, which means more useful information is given in the instruction and, thus, a higher quality of that sam-

ple. To obtain the IFD score, they train a LLaMa-7b as a warm-up selector on only a small portion of the candidate dataset to give the model basic instruction-following ability. Then, the IFD score can be computed by the following equation:

$$r_{\theta}(Q, A) = \frac{s_{\theta}(A|Q)}{s_{\theta}(A)} \quad (1)$$

where $r_{\theta}(Q, A)$ is the IFD score of a (Q, A) sample pair, θ means the warm-up selector model while $s_{\theta}(A|Q)$ and $s_{\theta}(A)$ are the likelihood of generating the same answer with or without giving the question as instruction.

Another work, **SuperFiltering** (Li et al., 2024a) adopts a modified version of the IFD score by replacing the likelihood function with perplexity values and selects samples with lower scores. Moreover, they use GPT-2 to train a smaller selector to determine data quality compared the previous work (Li et al., 2023).

4.2 Leveraging External Information

Leveraging external information for data quality measurement uses knowledge that is not accessible from the given candidate datasets or Pending Fine-tune Models (PFMs). PFM is the target model to be fine-tuned on the selected subset, using the same backbone model as SEM. To compute quality labels, external information is used in the form of either **discrete quality labels** from other LLMs (Chen et al., 2024; Lu et al., 2023; Liu et al., 2023) or **continuous quality labels** from sample influence (Cao et al., 2023; Wei et al., 2023; Li et al., 2024b; Xia et al., 2024), where the sample influence is reflected in model’s performance gain induced by the sample.

4.2.1 Discrete Quality Labels from LLM Preference

Obtaining discrete quality labels relies on the use of external LLMs for their direct scoring or annotations on the candidate data. Such quality labels are discrete since other LLMs work as a blackbox that takes in the candidate data and outputs a response reflecting their preference. Based on whether the external LLM is trainable, one can obtain quality labels either **exclusive LLMs**, which are closed-source commercial models like GPT-4, or **community LLMs**, which are open-source trainable models like LLaMa.

Exclusive LLM Preference. Exclusive LLM’s preference can be utilized as data quality labels

because many commercial models have high agreement with human annotators when evaluating the quality of data. **AlpaGasus** (Chen et al., 2024) determines data quality entirely from ChatGPT’s preference reflected in its direct scoring on each sample. The score is obtained by prompting ChatGPT with a scoring template with common evaluation aspects, like helpfulness and accuracy. Then, they select higher quality samples with higher scores. Compared with AlpaGasus’ straightforward prompt, **InsTag** (Lu et al., 2023) specifies clearer evaluation dimensions as tags when prompting ChatGPT. They propose a Complexity-first Diverse Sampling procedure for data selection. To obtain quality labels, their measure first sample-level complexity (average number of tags for each sample in the candidate subset) and then dataset-level diversity (the total number of distinct tags in the subset), balancing the interplay between data quality and diversity.

Community Model Preference. Community models are open-source, trainable models, which can be tailored for specific evaluation tasks after aligning them with external commercial models. **DEITA** (Liu et al., 2023) relies on the preference of a community model LLaMa to measure data quality, where the LLaMa learns from ChatGPT for scoring. To train the scorer, they utilize the ideas of evolving from WizardLM (Xu et al., 2023) to evolve a small set of sample seeds into different levels of complexity and quality and then fine-tunes a LLaMa on ChatGPT’s scoring on these evolved samples. For selection, they propose the Score-First, Diversity-Aware selection similar to that proposed in InsTag.

4.2.2 Continuous Quality Labels from Sample Influence

In search of more direct and model-specific data selection methods, this research line obtains continuous quality labels from sample influence, which is quantified by the performance improvement a model gains when fine-tuned on a sample. These improvements are gauged by continuous outputs like **model-centric** evaluation scores (Cao et al., 2023; Wei et al., 2023; Li et al., 2024b) or **data-centric** gradient similarity (Xia et al., 2024).

Model-centric. When using one sample to fine-tune the PFM, performance improvement in SEM is expected to reflect the quality of that sample. Based on the assumption, InstructionMining (Cao et al., 2023) constructs the mapping between the

9-dimensional-indicator representations of the sample and the inference loss (Wang et al., 2023; Zheng et al., 2023). Then, they utilize BLENDSEARCH, effectively combining global and local optimizations with bayesian optimization and different local search threads, to determine the final selected dataset size. InstructionGPT-4 adopts the same logic on a multimodal model that can also process visual-caption features as input. To further avoid the fine-tuning cost, **Nugget** (Li et al., 2024b) measures the sample influence towards SEM by prompting the PFM to answer the same set of questions with or without that certain sample. Better question-answering results indicate a larger improvement the sample brings to PFM.

Data-centric. Compared with model-centric methods, data-centric approach compares the similarity between candidate data and known high-quality data’s ability to improve model performance. To make comparison between data, **LESS** (Xia et al., 2024) proposes the Low-rank gradiEnt Similarity Search method. They first perform a warm-up LoRA to obtain gradient representations of the candidate dataset and then compare the similarity with high-quality dataset.

5 Data Selector Evaluation

To evaluate the usefulness of selectors, the method is to select a subset from a candidate dataset through the selector and then fine-tune a model to be the selectively enhanced model (SEM) based on this subset to compare the performance with the same model fine-tuned on full data (Baseline model, BM) or other popular oracle LLMs. Table 1 shows the detailed evaluation setting, including the choice of candidate datasets, counterpart models used in the comparison, and evaluation metrics that provide the performance.

Candidate Datasets. Most of the works (Li et al., 2024a, 2023; Liu et al., 2023) use the popular open-sourced datasets as candidate datasets to push forward better performance of fine-tuned models by selecting higher-quality samples in them. The candidate dataset is further divided into the typical group, including Alpaca, Dolly, FLAN, etc., and the advanced group developed from the typical datasets to achieve higher quality, including WizardLM, UltraChart, etc.

Counterpart Models. To objectively evaluate the performance of the SEM, most works choose BM as the counterpart model for comparison. They

Method	Candidate Datasets	Evaluating SEMs	Counterpart Models		Evaluation Metrics	
			BM	Others	Wins-ties-losses	Benchmark Scoring
AlpaGasus	Alpaca	LLaMA-2 7B	✓	✓	Vicuna, Koala, WizardLM, self-Instruct	InstructEval
Instruction-Mining	ALPACA & OPEN ASSISTANT, STACK-EXCHANGE & WIKIHOW	LLaMA 7B	✓	✗	OPENORCA & DOLLY	OPENLLM
InstructionGPT-4 IFD	MiniGPT-4	LLaMA-2	✓	✗	LLaVA-Bench	MME, VQA, MM-Bench
	Alpaca & WizardLM	LLaMA-2 7B	✓	✗	Vicuna, Koala, WizardLM, self-Instruct, LIMA	OPENLLM
Superfiltering	Alpaca & Alpaca-GPT4 & WizardLM	LLaMA-2 7B/13B	✓	✗	WizardLM	OPENLLM, AlpacaEval
Nugget	Alpaca	LLaMA-2 7B	✓	✗	-	MT-Bench, AlpacaEval
LESS	FLAN V2 & CoT & DOLLY & OPEN ASSISTANT 1	LLaMA2-13B; Mistral 7B	✓	✗	-	MMLU, TYDIQA, BBH
InsTag	WizardLM & UltraChat & ShareGPT	LLaMA-1/-2	✗	✓	-	MT-Bench
DEITA	Alpaca & DOLLY & Oassit & FLAN 2022 & WizardLM & UltraChat & ShareGPT	LLaMA-1/-2 13B; Mistral 7B	✗	✓	-	OPENLLM, MT-Bench

Table 1: The candidate dataset, SEMs, counterpart models, and evaluation metrics used in each method. Some works, such as Instruction-Mining, use part of several datasets mentioned to form a candidate dataset. The "✓" under BM means the work uses the same BM as the evaluating SEM.; under Other Models, it means the work uses many models other than BM, including oracle LLM and other fine-tuned SEM; under wins-ties-losses, it means the work uses various methods to evaluate wins-ties-losses, such as AlpacaEval, and directly using GPT-4 .

tend to use the popular LLaMa series (Chen et al., 2024; Lu et al., 2023) as well as Mistral (Liu et al., 2023; Xia et al., 2024) models as backbones of the SEM and BM to obtain relative improvement evaluation, which directly shows the improvement effect of the selector. Other works (Xia et al., 2024; Chen et al., 2024) compare the SEM with SOTA models (such as GPT-4, Claude, and LLaMA-Chat 7B) to obtain absolute improvement evaluation, which indicates how good SEM achieves.

Evaluation Metrics. Similar to the counterpart models, the evaluation metric adopts the relative and absolute methods to comprehensively evaluate the selector. The absolute metric uses Wins-ties-losses pairing scored by GPT-4 to indicate the direct performance difference between the SEM and counterpart model, while the absolute metric uses benchmark scoring to directly score and rank the SEM. Benchmark scoring is separated into a traditional group, which examines the loss of response on test tasks (such as MLU, TYDIQA, and Mosaic Eval Gauntlet), and a group, which uses GPT-4 to score on various tasks, including MT-Bench, MM-Bench, AlpacaEval, and VicunaQA et al.

6 Analysis of Data Selection Method

To spot the common designs that lead to superior performance, we analyze the efficiency and feasibility of the existing data selection methods, distinguishing the superior and inferior work. Efficiency is quantified to examine the selection competence filtering, which is measured based on the overall consideration of the performance of SEM and the data size (selected dataset fraction), while feasibility uses a qualitative method to evaluate the difficulty of implementation, which entails both

simplicity and flexibility.

6.1 Efficiency of the Selector

We manage to compare the data selection competence across different works by using efficiency, which refers to the expectation of probability in selecting the ground truth high-quality data at each bet (Appendix A.1).

Performance Improvement Ratio. The competence is reflected in the performance of SEM, where higher-quality data leads to a larger performance improvement. Performance improvement is the ratio of SEM’s performance to that of the counterpart model, evaluated under various settings. Therefore, we first classify the evaluation settings into four categories and use Equation 4 to calculate the overall improvement rate for each category, and then further unify the four categories into the one that is the wins rate of SEM to BM which directly reflects the improvement effect brought by the selected subset (detailed information is in Appendix A.2).

Selected Dataset Fraction. However, the increase in data size also improves the performance (Kaplan et al., 2020). To evenly evaluate the impact of data size between different works, we use the selected dataset fraction, which refers to the fraction of the selected subset to the entire dataset. Though the size of the selected dataset is directly provided in the works, it is abandoned because it is heavily affected by the size of the entire dataset varying from 3,439 (Wei et al., 2023) to 306,044 (Lu et al., 2023).

To acquire efficiency from performance improvement, we further use the selected dataset fraction to eliminate the impact of data size. Figure 3 reflects the efficiency of the selected dataset fraction and

performance improvement ratio. In addition, as there are many pairs of dataset fractions and performances in a work, the one that achieves the optimal performance is adopted.

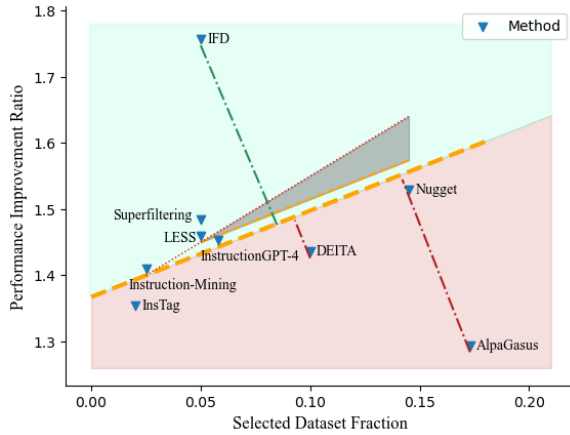


Figure 3: The comparison among popular data selection methods. The yellow line is the baseline based on Instruction-Mining and InstructionGPT-4. The signed distance between the method point and baseline is the efficiency difference between methods and baselines.

From the position of work in Figure 3, IFD has the highest efficiency because it achieves a colossal performance improvement with a small size of data. By contrast, AlpaGasus, which uses the largest selected dataset fraction but achieves the lowest performance improvement ratio, is the least efficient. Besides these two works, the superiority and inferiority of other works are elusive. Therefore, we develop the efficiency curve assumption to make the rest of the works mutually comparable (details are in Appendix A.3).

Under this assumption, the yellow dashed line represents the efficiency curve of Instruction-Mining and InstructionGPT-4, which is the baseline for separating the superior and inferior works. The grey area illustrates an infeasible area that every work has, where the efficiency of other work is incomparable (Appendix A.4). The red and green line, respectively indicates the superior and inferior efficiency of the work in comparison to the baseline, which is the signed distance between the baseline and the work, as calculated by the following:

$$Eff_i = \frac{Ax_i + By_i + C}{\sqrt{A^2 + B^2}} \quad (2)$$

$$l_{base} : Ax + By + C = 0 \quad (3)$$

where (x_i, y_i) is the position of work i , and l_{base} is the mathematical expression of the baseline. We

present the comparison efficiency of each work in Figure 4, where the efficiency of DEITA is invisible because it is in the infeasible area of InsTag. We consider DEITA to be better than InsTag because it refines the method InsTag used.

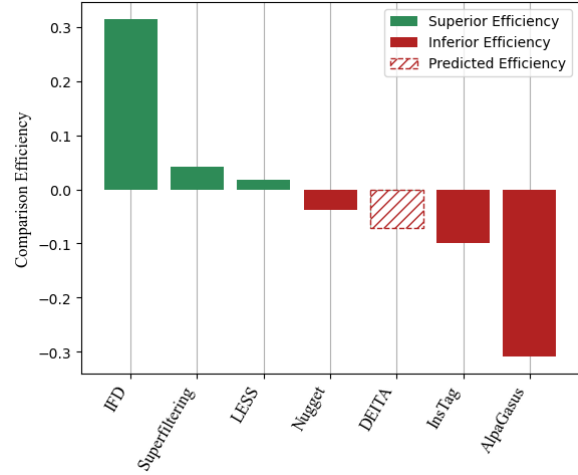


Figure 4: The comparison efficiency of each work relative to Instruction-Mining and InstructionGPT-4.

Figure 4 shows that the SEM-specific work exempted from indirect preference tends to achieve high efficiency. All the superior works use the BM LLM loss to measure the quality according to SEM preference instead of human or Oracle LLM preference, while the three least efficient works use the Oracle LLM score. Specifically, IFD achieves the highest efficiency by leveraging internal information, which is independent of all outer preferences and biases. LESS is far worse because it resorts to leveraging external information, which introduces human preferences. AlpaGasus is the worst because it solely relies on Oracle LLM to generate quality scores, whose preference deviates far from the SEM.

6.2 Feasibility of the Selector

Superior work should not only have high efficiency but also high feasibility. Feasibility employs the simplicity and the flexibility of the selector to respectively assess its implementation difficulty and competence in handling new selection tasks. We develop the feasibility rank of each work in Table 2, based on the consideration of these two aspects.

Simplicity. We evaluate the simplicity of method from its cost of implementation and reproducibility. The cost of implementing a method focuses on the training and inference cost of model,

Methods	Feasibility Rank	Simplicity	Flexibility
AlpaGasus	1	1	1
InsTag	2	3	1
Nugget	2	2	2
IFD	3	4	4
Superfilter	3	4	4
LESS	3	4	4
DEITA	4	6	3
Instruction-Mining	4	7	2
InstructionGPT-4	5	5	5

Table 2: The feasibility rank, and the corresponding simplicity and flexibility rank of each work. A smaller number indicates the work does better.

as the algorithm cost is so much smaller that it can be neglected. The cost of implementation considers the number of time involves model training and inference, and the actual number of models being trained. The reproducibility considers whether the implementation details is provided and code is open-sourced. Table 4 in Appendix A.5 shows the simplicity rank and the above considerations of each work. Specifically, the work uses oracle LLM score as quality measurement is with high simplicity, where AlpaGasus is the simplest, who involves 1 time of model inference.

Flexibility. Flexibility evaluates the extensibility and transferability of the selector. Extensibility examines the flexibility of reforming the model used in the selector, and transferability considers whether the work is dataset-independent or model-independent when dealing with migration tasks. Table 5 in Appendix A.5 shows the flexibility rank and the two considerations of each work, which are respectively derived from two questions: (1) Whether substitute the SOTA model used in the selector into open-sourced model defunct the selector; (2). Whether the model used in the selector needs to be retrained to maintain the optimal efficiency in handling new candidate datasets or SEM. Specifically, all the works that rely on oracle LLM to acquire quality scores are model inextensible, while six works using PLM loss are model dependent(Li et al., 2024b; Xia et al., 2024; Wei et al., 2023; Cao et al., 2023; Li et al., 2023).

6.3 Overall Consideration of the Selector

According to the overall performance of the works on efficiency and feasibility, we find out that: (1) The model-specific method achieves high efficiency without cost of feasibility; (2) Complex method can improve the efficiency, if it is deliberately designed to avoid the negative effect on model-specification, but it is always accompanied

with the feasibility loss.

Specifically, We divide the works into a group using oracle LLM score (Li et al., 2023, 2024a; Xia et al., 2024; Cao et al., 2023; Wei et al., 2023; Li et al., 2024b) and a group using PFM loss (Lu et al., 2023; Chen et al., 2024; Liu et al., 2023), based on the different choice of data quality measurement. The PFM loss is model-specific which examine the loss of PFM generated by data, while oracle LLM score is human-basic which leverages oracle LLM to imitate human scoring. The work in PFM loss group is always more efficient than the work in another group With the same feasibility rank, indicating that PFM loss is superior to oracle LLM score.

Moreover, these two groups have opposite relationship between efficiency and feasibility. The work using oracle LLM score has a trade-off relationship between these two aspects. AlpaGasus has the worst efficiency but highest feasibility, while DEITA and InsTag sacrifice feasibility to achieve better efficiency by using more complex methods. By contrast, in the group of BM LLM loss, the one with high efficiency typically has low feasibility. Among these works, IFD manages to achieve both high efficiency and feasibility, because it only resorts to the internal information without introducing redundant element during filtering. LESS, who leverages on an outer dataset to filter data, does worse on both aspects. The simple work tends to achieve high efficiency, because they are direct who truncate the irrelevant and detrimental information from data selection. The result shows that the complex method may be beneficial or detrimental to the efficiency, which depends on its impact on the directness of method.

7 Discussions

7.1 Trend

Figure 5 shows the trend of data selector. To achieve superior overall performance, the data selector becomes more model-specific which reflects on the goal designation and the choice of quality approximation method. Specifically, the goal is changed from selector motivation to task motivation. The early works aim to develop the selector with high transferability to handle any data selection task (Cao et al., 2023; Wei et al., 2023; Chen et al., 2024), while the later works focus on improving the performance of SEM of the specific task (Li et al., 2023; Xia et al., 2024).

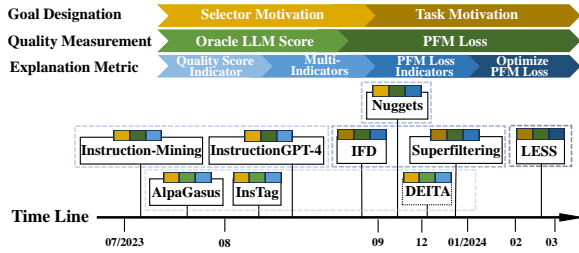


Figure 5: The timeline of the data selection methods.

The approximation method uses an explanation metric to predict the approximation object, which refers to the data quality label generated by the quality measurement. The quality measurement developed from relying on external oracle LLM score (Chen et al., 2024) to internal PFM loss (Xia et al., 2024) becomes more direct to the PFM, while the explanation metric becomes more complex from using concrete indicator to abstract indicator, and from using single to multiple indicators. The early works tend to use explicit indicators such as quality score and diversity, which introduce semantic factors to explain data quality, while the later works use abstract indicators developed from the PFM loss to reflect the data quality. On the other hand, the number of explicit indicators used in the selector increases. DEITA and InsTag employ more explicit indicators than AlpaGasus, which solely relies on quality scores from oracle LLM as quality indicators. Specifically, data diversity as the explicit indicator can largely improve the overall performance. DEITA and InsTag achieve far better overall performance than AlpaGasus because they take diversity into consideration.

7.2 Challenges

Through the above review and analysis, although there has been significant progress in data selection for fine-tuning LLMs, there are still three challenges.

(1) **Lack of unified and efficient metrics for high-quality data.** The existing data selection methods still have a vague definition of high quality. Although some methods consider explainable linguistic features or the complexity of data as well as the diversity of sampling, most of them focus on improving the performance of the model, that is, samples that can improve model performance are high-quality samples. Interestingly, recent work (Bai et al., 2024) has shown that some seemingly low-quality datasets considering their data sources also

have improvement effects on LLM fine-tuning. Our analysis has shown that different methods have their own strengths, and the main challenge is to organically combine their advantages to build a unified and efficient data selection method for fine-tuning.

(2) **Lack of data selection methods for specific domains.** Most data selection methods focus on overall performance improvement, but the contribution of selected data to different domains is not the same. The existing works (Cao et al., 2023; Wei et al., 2023; Chen et al., 2024; Lu et al., 2023; Li et al., 2024b) demonstrated that selected data can bring significant improvements in writing and role-playing but minor improvements in mathematics and reasoning. Therefore, future work needs to consider how to dynamically select data based on the shortcomings of the model in a specific domain to compensate for its performance.

(3) **Lack of data selection methods for multi-turn conversations.** Most existing data selection methods are aimed at single-turn conversations because their quality is easier to measure but lacks attention to multi-turn conversation data. Although current work (Lu et al., 2023; Li et al., 2024b; Liu et al., 2023) evaluates models in static multi-turn conversation scenarios such as MT-Bench, there is a lack of suitable metrics for measuring multi-turn conversation data quality for data selection.

8 Conclusion

In this paper, we conducted an extensive survey on data selection for fine-tuning large-scale language models. We first construct a three-stage data selection scheme for the entire process and review the current research progress of data selection based on it, including data preprocessing, data selector construction, and data selector evaluation. To address the issue of incompatibility caused by different experimental settings, we propose a quantitative evaluation based on the assumption of efficiency curves to compare the existing work. We also qualitatively analyzed the feasibility of existing work, including implementation costs and reproducibility. We find that the model-specific work achieves high efficiency, whereas complex method can improve the efficiency if it is deliberately designed to avoid the negative effect on model-specification, but it is always accompanied by feasibility loss. We have summarized the existing trends and provided insights for future research.

656 Limitation

657 (1) In analyzing the efficiency of the selector,
658 the two assumptions are proposed to make the effi-
659 ciency of the method comparable. Because of the
660 lack of information on the efficiency curve of the
661 method, we use distance to demonstrate the effi-
662 ciency of the methods, which generates the prob-
663 lem of infeasible area. Following the increase of
664 work included in a comparable group, the infeasible
665 area enlarges, which limits the number of works
666 in the group.

667 (2) This paper mainly focuses on data selec-
668 tion for instruction fine-tuning LLMs instead of
669 data rewriting or enhancement. Although we have
670 already comprehensively examined the existing
671 works, there may still be some works we neglected.

672 References

673 Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne
674 Longpre, Nathan Lambert, Xinyi Wang, Niklas
675 Muennighoff, Bairu Hou, Liangming Pan, Hae-
676 won Jeong, Colin Raffel, Shiyu Chang, Tatsunori
677 Hashimoto, and William Yang Wang. 2024. [A survey on data selection for language models](#).

679 Yuelin Bai, Xinrun Du, Yiming Liang, Yonggang Jin,
680 Ziqiang Liu, Juntong Zhou, Tianyu Zheng, Xincheng
681 Zhang, Nuo Ma, Zekun Wang, Ruibin Yuan, Haihong
682 Wu, Hongquan Lin, Wenhao Huang, Jiajun Zhang,
683 Wenhao Chen, Chenghua Lin, Jie Fu, Min Yang, Shi-
684 wen Ni, and Ge Zhang. 2024. [Coig-cqia: Quality is all you need for chinese instruction fine-tuning](#).

686 Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun.
687 2023. [Instruction mining: When data mining meets large language model finetuning](#).

689 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa
690 Gunaratna, Vikas Yadav, Zheng Tang, Vijay Sriniva-
691 san, Tianyi Zhou, Heng Huang, and Hongxia Jin.
692 2024. [Alpagasus: Training a better alpaca with fewer data](#).

694 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
695 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
696 Wang, Mostafa Dehghani, Siddhartha Brahma, Al-
697 bert Webson, Shixiang Shane Gu, Zhuyun Dai,
698 Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-
699 ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,
700 Dasha Valter, Sharan Narang, Gaurav Mishra, Adams
701 Yu, Vincent Zhao, Yanping Huang, Andrew Dai,
702 Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-
703 cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,
704 and Jason Wei. 2022. [Scaling instruction-finetuned language models](#).

706 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan
707 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and

Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). 708
709

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B
Brown, Benjamin Chess, Rewon Child, Scott Gray,
Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.
[Scaling laws for neural language models](#). *arXiv preprint arXiv:2001.08361*. 710
711
712
713
714

Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu
Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou.
2024a. [Superfiltering: Weak-to-strong data filtering for fast instruction-tuning](#). 715
716
717
718

Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang
Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and
Jing Xiao. 2023. [From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning](#). 719
720
721
722
723

Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiayi Yang,
Min Yang, Lei Zhang, Shuzheng Si, Junhao Liu,
Tongliang Liu, Fei Huang, and Yongbin Li. 2024b.
[One shot learning as instruction data prospector for large language models](#). 724
725
726
727
728

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](#). 729
730
731
732

Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Jun-
yang Lin, Chuanqi Tan, Chang Zhou, and Jingren
Zhou. 2023. [instag: Instruction tagging for analyzing supervised fine-tuning of large language models](#). 733
734
735
736

Elizaveta Moskvskaya, Olesya Chebotareva, Valeria
Efimova, and Sergey Muravyov. 2023. Predicting
dataset size for neural network fine-tuning with a
given quality in object detection task. *Procedia Com-
puter Science*, 229:158–167. 737
738
739
740
741

Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad
Saqib, Saeed Anwar, Muhammad Usman, Naveed
Akhtar, Nick Barnes, and Ajmal Mian. 2023. [A comprehensive overview of large language models](#). 742
743
744
745

Chen Sun, Abhinav Shrivastava, Saurabh Singh, and
Abhinav Gupta. 2017. Revisiting unreasonable effec-
tiveness of data in deep learning era. In *Proceedings
of the IEEE international conference on computer
vision*, pages 843–852. 746
747
748
749
750

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
Kaiser, and Illia Polosukhin. 2023. [Attention is all you need](#). 751
752
753
754

Jiahao Wang, Bolin Zhang, Qianlong Du, Jiajun Zhang,
and Dianhui Chu. 2024. [A survey on data selection for llm instruction tuning](#). 755
756
757

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa
Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh
Hajishirzi. 2023. [Self-instruct: Aligning language models with self-generated instructions](#). 758
759
760
761

762 Yizhong Wang, Swaroop Mishra, Pegah Alipoormo-
 763 labashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva
 764 Naik, Arjun Ashok, Arut Selvan Dhanasekaran,
 765 Anjana Arunkumar, David Stap, Eshaan Pathak,
 766 Giannis Karamanolakis, Haizhi Lai, Ishan Puro-
 767 hit, Ishani Mondal, Jacob Anderson, Kirby Kuznia,
 768 Krma Doshi, Kuntal Kumar Pal, Maitreya Patel,
 769 Mehrad Moradshahi, Mihir Parmar, Mirali Purohit,
 770 Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma,
 771 Ravsehaj Singh Puri, Rushang Karia, Savan Doshi,
 772 Shailaja Keyur Sampat, Siddhartha Mishra, Sujan
 773 Reddy A, Sumanta Patro, Tanay Dixit, and Xudong
 774 Shen. 2022. [Super-NaturalInstructions: Generaliza-
 775 tion via declarative instructions on 1600+ NLP tasks.](#)
 776 In *Proceedings of the 2022 Conference on Empirical
 777 Methods in Natural Language Processing*, pages
 778 5085–5109, Abu Dhabi, United Arab Emirates. As-
 779 sociation for Computational Linguistics.

780 Lai Wei, Zihao Jiang, Weiran Huang, and Lichao Sun.
 781 2023. [Instructiongpt-4: A 200-instruction paradigm
 782 for fine-tuning minigpt-4.](#)

783 Mengzhou Xia, Sadhika Malladi, Suchin Gururangan,
 784 Sanjeev Arora, and Danqi Chen. 2024. [Less: Select-
 785 ing influential data for targeted instruction tuning.](#)

786 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,
 787 Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin
 788 Jiang. 2023. [Wizardlm: Empowering large language
 789 models to follow complex instructions.](#)

790 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
 791 Zhuang, Zhonghao Wu, Yonghao Zhuang, Zi Lin,
 792 Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang,
 793 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging
 794 llm-as-a-judge with mt-bench and chatbot arena.](#)

795 Ming Zhong, Yang Liu, Da Yin, Yuning Mao, Yizhu
 796 Jiao, Pengfei Liu, Chenguang Zhu, Heng Ji, and
 797 Jiawei Han. 2022. [Towards a unified multi-
 798 dimensional evaluator for text generation.](#) In *Pro-
 799 ceedings of the 2022 Conference on Empirical Meth-
 800 ods in Natural Language Processing*, pages 2023–
 801 2038, Abu Dhabi, United Arab Emirates. Association
 802 for Computational Linguistics.

803 Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao
 804 Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu,
 805 Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis,
 806 Luke Zettlemoyer, and Omer Levy. 2023. [Lima: Less
 807 is more for alignment.](#)

808 A Appendix

809 A.1 Efficiency and High Quality Dataset

810 **Efficiency Definition** The efficiency of work
 811 is the expectation of probability in selecting the
 812 ground truth high quality data at each bet, which is
 813 derived from the consideration on the quality struc-
 814 ture of the selected subset. The quality structure
 815 refers to the proportion of high quality data to the
 816 size of dataset.

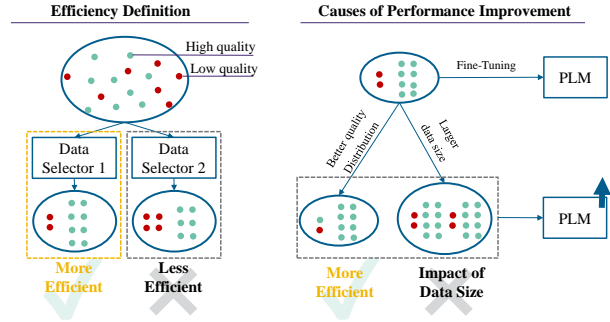


Figure 6: The left figure indicates the definition of efficiency is related to the quality distribution of selected subset. The right figure shows that the performance improvement of the corresponding fine-tuned PLM comes from both a better quality distribution and a larger data size

In Figure 6, data selector 1 and 2 selects two subsets with different quality structure from a candidate set. Subset 1 has better quality structure than subset 2, because data selector 1 has higher probability in selecting the high quality data at each bet, whose expectation is the proportion of high quality data to the size of subset.

Causes of Performance Improvement The quality of dataset decides the performance of PLM who is fine-tuned on it (Zhou et al., 2023). However, the performance also increases following the increase of data size, if the augmented data set has the same quality distribution with the original one (Chung et al., 2022).

831 A.2 Performance Improvement Ratio

832 Table 3 shows the performance improvement
 833 developed from the original performance values
 834 under 4 kinds of evaluation setting which we con-
 835 clude and extract from the primitive article. Wins
 836 rate and benchmark improvement is respectively de-
 837 veloped from the wine-ties-losses and benchmark
 838 scoring using the formula:

$$839 \frac{1}{n} \sum_{i=0}^n \frac{X_i}{Y_i} \quad (4)$$

840 where n is the total number of the evaluation
 841 settings using the same counterpart model with
 842 different evaluation metrics or candidate datasets,
 843 X_i and Y_i is respectively the score of the SEM and
 844 the counterpart model under the same evaluation
 845 setting i .

846 Then, the SEM Wins rate under Same Counter-
 847 part Model is chosen as the unified performance
 848 improvement ratio. Missing value is calculated by

Method	Same Counterpart Model (BM)		SEM	Other Counterpart Models (LLaMA chat 7B/13B)	
	Wins Rate	Benchmark Improvement		Wins Rate	Benchmark Improvement
AlpaGasus	1.284	0.949	LLaMA-2 7B	-	-
Superfiltering	1.475	0.962	LLaMA-2 7B	-	-
InsTag	1.344	-	LLaMA 13B	-	0.985
DEITA	1.467	-	LLaMA-2 13B	-	1.000
InstructionGPT-4	1.443	-	MiniGPT-4	-	-
Nuggets	1.519	-	LLaMA-2 7B	-	-
IFD	1.747	-	LLaMA-2 7B	-	-
LESS	-	0.973	LLaMA-2 13B	-	-
Instruction-Mining	-	-	LLaMA-2 7B	0.212	0.991

Table 3: The table shows the performance improvement under four evaluation settings which we deliberately choose and leverage on, to generate a unified performance improvement rate for each method. In the Other Counterpart Models, the parameter size of LLaMA chat is chosen as the same as the SEM to offset the impact of model size.

leveraging on the other method that bridges the value from other evaluation settings to the unified performance improvement ratio.

A.3 Assumption of Efficiency Curve

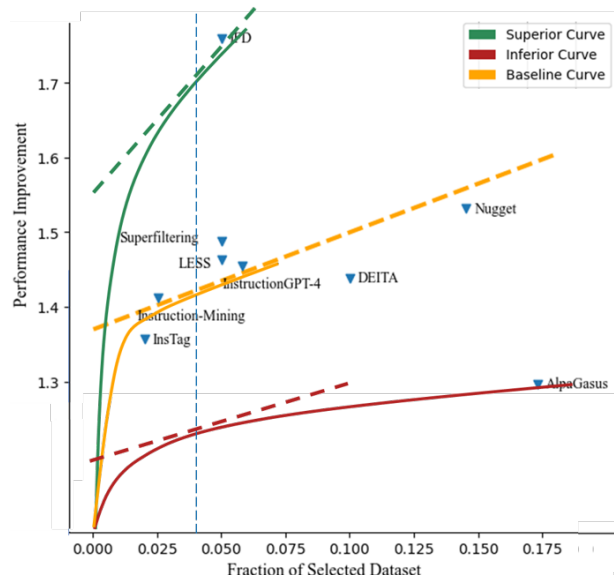


Figure 7: The demonstration of potential efficiency curve. The green, yellow, and red dashed lines represent the slopes of the methods at the same value of fraction.

We develop these assumptions to construct the efficiency curve of the work. The work on the same theoretical efficiency curve has the same efficiency, where the performance is purely affected by the impact of data size:

The first assumption. We assume for a dataset with fixed structure, its function of the performance improvement ratio and selected dataset fraction complies to the logarithm-like function which is upwarded, concaved, and approaching to linear after experiencing a rapid but short increasing.

The second assumption. We assume for two methods with different efficiency, the function's slope of the superior one is always larger than the

inferior one in the whole feasible domain of the selected dataset fraction which is between 0 and 1.

Figure 7 demonstrates the potential efficiency curves of three works. As IFD is the most efficient and AlpaGasus is the least efficiency, the slope of IFD is larger than AlpaGasus at the same selected dataset fraction.

For the first assumption, many articles suggest that the impact of logarithm data size on the loss is linear, if the augmented dataset maintains the same quality structure (Kaplan et al., 2020; Sun et al., 2017; Moskovskaya et al., 2023). Assumption 1 extends this relationship to the pair of the performance improvement ratio and selected dataset fraction. Moreover, the statement of rapid but short growth complies with the fact that the slope of baseline (≈ 1.303) is far less than the slope ($= 56$) between Instruction-Mining and the original point, which implies a rapidly growth of performance at early stage.

The second assumption can be intuitively deduced from the first assumption with the fact that high quality data leads to better performance of SEM (Zhou et al., 2023). Therefore, if the methods is superior which indicate its selected dataset is with good structure, its increasing on the performance improvement must be greater than the inferior one at every point selected dataset fraction.

A.4 Infeasible Area

Because of a lack of information, each work generates an infeasible area, which is in fact the possible area of its efficiency curve. Therefore, if other work is in the infeasible area, it is incomparable with the work that generates this infeasible area. The infeasible of the inferior work and superior work is generated differently, where Figure 8 shows respectively the by using LESS and InsTag. For both superior and inferior works, the yellow boundary of the infeasible area is parallel to the baseline.

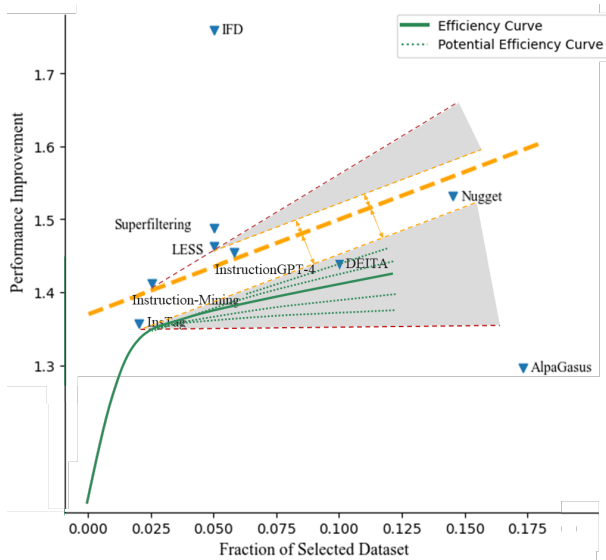


Figure 8: The demonstration of infeasible area. The green line is the efficiency curve of InsTag, where the dashed lines indicates its potential position.

906 For inferior work, the red boundary is horizontal
 907 because it cannot perform worse with a larger data
 908 size; For superior work, the red boundary is the line
 909 between the work and Instruction-Mining because
 910 the efficiency curve will never penetrate each other
 911 under the assumption.

912 A.5 Detailed Information of Feasibility

913 Table 4 shows the complexity rank and the corre-
 914 sponding considerations in structure difficulty and
 915 cost of running. Table 5 shows the flexibility rank
 916 and the corresponding considerations of each work.

Data Selection Method	Cost of Implementation			Reproducibility		Complexity Rank
	# Times involves Model Training	# Times involves Model Inference	# Models Being Trained	Implementation Details	Code Open Sourced	
AlpaGasus	0	0	1	✓	✗	1
Nuggets	0	2	0	✓	✓	2
InsTag	0	2	0	✓	✗	3
IFD	1	2	1	✓	✓	4
Superfilter	1	2	1	✓	✓	4
LESS	1	2	1	✓	✓	4
InstructionGPT-4	1	3	1	✓	✓	5
DEITA	2	3	1	✓	✓	6
Instruction-Mining	1	0	129	✓	✗	7

Table 4: The complexity rank and the corresponding considerations in structure difficulty and cost of running.

Method	Extensibility		Transferability		Flexibility Rank
	Model Independent	Dataset Independent	Model Independent	Model Independent	
InsTag	✗	✓	✓	✓	1
AlpaGasus	✗	✓	✓	✓	1
Nuggets	✓	✓	✗	✗	2
Instruction-Mining	✓	✓	✗	✗	2
DEITA	✗	✗	✓	✓	3
LESS	✓	✗	✗	✗	4
Superfiltering	✓	✗	✗	✗	4
IFD	✓	✗	✗	✗	4
InstructionGPT-4	✗	✗	✗	✗	5

Table 5: The feasibility rank and the corresponding considerations. In the real practicing, the indicator under Extensibility and Transferability has different priority. Therefore, we consider they contribute differently to the rank of feasibility, if the work has the same number of ✓. The priority from the highest to lowest is: Model Independent in Transferability, Dataset Independent, Model Independent in Extensibility.