



BRIEF: Bi-level Coreset Selection for Efficient Instruction Tuning in LLMs

Chaoyuan Shen*
Beijing Institute of Technology
Beijing, China
scy@bit.edu.cn

Chi Zhang*
Beijing Institute of Technology
Beijing, China
zc315@bit.edu.cn

Chengliang Chai[†]
Beijing Institute of Technology
Beijing, China
ccl@bit.edu.cn

Jiacheng Wang
Beijing Institute of Technology
Beijing, China
wangjc@bit.edu.cn

Jia Yuan
University of Arizona
Tucson, AZ, USA
jiayuan@arizona.edu

Yuping Wang
Beijing Institute of Technology
Beijing, China
wyp_cs@bit.edu.cn

Ye Yuan
Beijing Institute of Technology
Beijing, China
yuan-ye@bit.edu.cn

Guoren Wang
Beijing Institute of Technology
Beijing, China
wanggr@bit.edu.cn

Lei Cao
Massachusetts Institute of Technology
Cambridge, MA, USA
lcao@csail.mit.edu

ABSTRACT

Instruction tuning is a key step in adapting large language models (LLMs) to effectively understand and follow human instructions. It enables LLMs to transform general knowledge into task-specific responses that align with user intent. Although many high-quality instruction tuning datasets have been released, efficiently utilizing these data sources during supervised fine-tuning (SFT) is important, as training on the full high-quality corpus can be computationally expensive. To address this inefficiency, we explore whether a compact, high-quality subset of instruction data can achieve comparable performance to full-dataset SFT, thereby reducing training cost without sacrificing effectiveness. To this end, this work proposes to select such a subset (a.k.a., coreset) of instruction examples that maintains comparable downstream performance while improving training efficiency. The key idea is inspired by our discovered decomposition that in instruction tuning, the training loss can be decomposed into two components that effectively quantify the contribution of an instruction to the two fundamental capabilities of LLMs, namely knowledge-related capability and instruction following capability. We then revisit the objective of the classical coreset approaches to balance the two capabilities when selecting instruction examples. Based on a bi-level formulation and a composite gradient distance that makes the objective submodular, we design an effective algorithm to achieve a bounded approximation error. Experiments on 4 datasets across 9 downstream tasks demonstrate that BRIEF reduces computational costs by 3× while improving accuracy by 5% on Llama-3.1-8B, Qwen3-4B and Mistral-Nemo-12B.

PVLDB Reference Format:

Chaoyuan Shen, Chi Zhang, Chengliang Chai, Jiacheng Wang, Jia Yuan, Yuping Wang, Ye Yuan, Guoren Wang, and Lei Cao. BRIEF: Bi-level Coreset Selection for Efficient Instruction Tuning in LLMs. PVLDB, 19(6): 1264-1277, 2026.

*Chaoyuan Shen and Chi Zhang are co-first authors.

[†]Chengliang Chai is the corresponding author.

doi:10.14778/3797919.3797933

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/HR10108/BRIEF>.

1 INTRODUCTION

Instruction tuning, also known as supervised fine-tuning (SFT), improves generic Large Language Models (LLMs) [19, 39, 72] in following human instructions [49], and enables adaptation to specific domains [32, 61].

Recent work [30, 35, 61, 67] shows that domain SFT crucially depends on selecting instruction examples, since many instruction-tuning datasets contain irrelevant or even harmful data for the target domain [28, 53, 61]. In contrast, when SFT is used to improve *general* instruction following [11, 18, 65], data selection is largely overlooked, despite the high cost of training on massive corpora [22, 46, 70]. This makes large-scale SFT unaffordable for many small organizations, motivating data selection in this setting.

Specifically, if we can select a high-quality, representative subset (i.e., a *coreset* [14, 44]) such that fine-tuning on it yields performance competitive with training on the full set, we can greatly reduce SFT cost.

In deep learning, coreset methods [40, 41, 51] construct weighted subsets whose aggregated gradients approximate the full gradient. Their guarantees are typically based on bounding the gradient approximation (**GA**) error: the difference between the full-data gradient and the weighted sum of coreset gradients.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 19, No. 6 ISSN 2150-8097.

doi:10.14778/3797919.3797933

However, unlike traditional deep learning, SFT simultaneously improves two capabilities [60]: knowledge-related capability (KN) and instruction following capability (IF). KN measures factual correctness (e.g., Q: “What is the orbital period of the Moon around the Earth?” A: “Approximately 27.3 days”), while IF measures executing diverse instructions (e.g., Q: “Prepare a report following these specified steps.” A: “[Detailed step-by-step report]”). By decomposing the SFT loss into KN and IF, we expose why existing coreset objectives can be insufficient, and design BRIEF to preserve both capabilities.

Decomposition of the SFT Loss. First, in Section 4.1, we decompose SFT loss into KN and IF in measurable forms, making explicit that SFT jointly optimizes both terms and enabling precise definitions.

Accordingly, an ideal SFT coreset should preserve both KN and IF, i.e., simultaneously constrain GA errors for the gradients of each term. Yet directly applying traditional coreset selection to SFT typically minimizes GA error only for the *total* SFT loss, which can yield a coreset that matches the overall gradient while failing to control the two components.

Consider samples whose gradients align differently with KN and IF. Some samples are strong for one capability because they have large norms and directions close to KN or IF, and they still contribute to the other capability. Other samples are “average”, they have small norms and directions between the two, so they are less informative. Methods that match the *average* full gradient often prefer many average samples because their small and balanced gradients lie near the global mean. Intuitively, repeatedly training on such mediocre examples is suboptimal. We would rather select strong samples, possibly pointing to different directions, and make their *combined* gradient match the full gradient.

This issue can be understood via error vectors. The SFT error vector is the difference between the full gradient and the coreset approximation, and it equals the sum of the KN and IF error vectors. Even if the SFT error magnitude is small, the individual GA errors for KN and IF can be arbitrarily large when their error vectors nearly oppose each other, meaning the angle is close to 180° . By the triangle inequality, a small $\|\mathbf{e}_{\text{SFT}}\| = \|\mathbf{e}_{\text{KN}} + \mathbf{e}_{\text{IF}}\|$ does not imply small $\|\mathbf{e}_{\text{KN}}\|$ or $\|\mathbf{e}_{\text{IF}}\|$. Figure 1 illustrates this effect.

New Objective Function. To address this, BRIEF minimizes the *sum* of GA errors for KN and IF, instead of only the GA error of their gradient sum. This directly bounds both individual errors, and also bounds the SFT GA error because $\|\mathbf{e}_{\text{KN}} + \mathbf{e}_{\text{IF}}\| \leq \|\mathbf{e}_{\text{KN}}\| + \|\mathbf{e}_{\text{IF}}\|$ (Figure 1(c)). Thus the selected subset retains both knowledge-related and instruction following capabilities.

Optimization Methodology. We prove the reformulated problem is NP-hard and propose a bi-level framework BRIEF. Using an auxiliary variable to partition the solution space, we obtain an equivalent bi-level problem. At the lower level, given a partition, we form a unified objective by weighting GA errors of KN and IF. Via a gradient distance composition technique, this becomes a submodular maximization problem solvable by an efficient greedy algorithm with an approximation guarantee. The upper level searches over partitions: we show the objective upper bound behaves monotonically at both ends of the auxiliary-variable range, enabling efficient exploration by ternary search and avoiding exhaustive enumeration.

Contributions. We summarize our main contributions as follows:

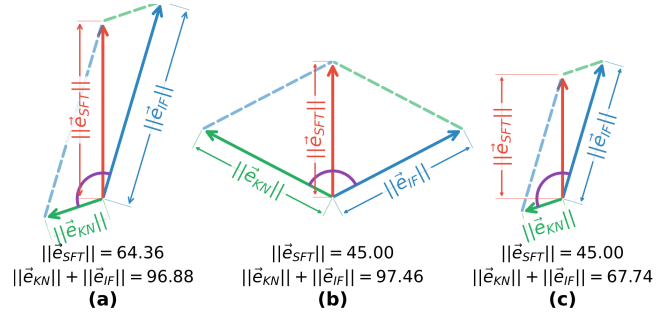


Figure 1: Different GA-error objectives: (a) original vectors, (b) minimizing GA error of the gradient sum, (c) minimizing sum of individual GA errors.

- By decomposing the SFT loss, for the first time we quantify the contribution of each instruction–response pair to the model’s knowledge-related capability and instruction following capability, offering a new loss-level perspective on instruction tuning.
- Building on this decomposition, we derive a new optimization objective for coreset selection in SFT that minimizes the sum of individual GA errors, rather than only the GA error of the gradient sum, thereby simultaneously retaining both knowledge-related and instruction following capabilities.
- By formalizing and solving a bi-level optimization problem driven by our composite gradient distance, which jointly controls the GA errors of KN and IF, we obtain an efficient coreset selection algorithm that significantly reduces the complexity of coreset selection while providing a theoretical guarantee on the resulting GA error bound.
- Experiments on several advanced LLMs and real world datasets show that BRIEF selects a well-performing coreset, reducing computational costs by $3\times$ while improving accuracy by 5% on Llama-3.1-8B, Qwen3-4B and Mistral-Nemo-12B.

2 RELATED WORK

We first review quality-based and task aware data selection methods for instruction tuning, then discuss gradient-based coreset methods originally developed for deep learning and their extensions to instruction-tuning settings, and finally connect these approaches to the bilevel optimization framework that underlies many recent coreset-based methods.

Quality-based Filtering for Instruction Tuning. Early work focused on removing low-quality data using simple rules [50, 56] or deduplication methods like SemDedup [1]. However, these heuristics often fail to generalize. Other researchers employ powerful LLMs (e.g., GPT-4) to score data quality, though this relies heavily on human-like intuition [17, 59, 71]. Perplexity is another common metric [6, 38, 43, 58], but it tends to favor simple and redundant sentences. To address this, Li et al. [31] introduced the Instruction Following Difficulty (IFD) score, which measures difficulty by comparing response generation with and without context.

Domain and Task aware Data Selection. Many methods select data to match specific downstream domains. Some use n -gram similarity to align corpora with validation sets [16, 63]. Others use

Table 1: Symbols used throughout the paper.

Symbol	Description
θ, Θ	Trainable parameters; Feasible parameter space
D, N	Full training dataset; Number of examples
(x_i, y_i)	i -th instruction–response pair from D
$ y_i , y_i^t, y_i^{<t}$	Token length; t -th token; Prefix before t
$p_\theta(\cdot)$	Model conditional probability under θ
$\mathcal{L}_{\text{SFT}}, \mathcal{L}_{\text{KN}}, \mathcal{L}_{\text{IF}}$	SFT loss; Knowledge loss; Instruction following loss
$\nabla \mathcal{L}_i(\theta)$	Per-example gradient of the loss at parameters θ
C, W, ω_j	Selected Coreset; Weight set; Weight of element j
$\gamma(j)$	Index map linking coreset item j back to D
K, ϵ	Size budget; Allowed GA error
$\text{IFD}(y x; \theta), \text{PPL}(\cdot)$	Instruction following difficulty; Perplexity
$d_{ij}^{\text{KN}}, d_{ij}^{\text{IF}}$	Pairwise gradient distances in KN and IF spaces
d_{ij}^{SFT}	Composite distance combining KN and IF terms
$\mathcal{B}_{\text{KN}}(C), \mathcal{B}_{\text{IF}}(C)$	Upper bounds on KN and IF GA errors for C
α, δ	Error budget split auxiliary parameter; Search tolerance
$F(C)$	Facility-location objective
$K_{\text{KN}}(\tau), K_{\text{IF}}(\tau)$	Minimum coreset size achieving error budget τ

influence functions to measure the impact of training data [68], as seen in LESS [62] and MATES [66]. Another approach involves training surrogate models to predict data relevance. For instance, DeepSeekMath [53] uses active learning, while RHO-1 [34] employs a high-quality model for token-level filtering. However, training these specialized classifiers or surrogate models often requires significant computational resources and may not adapt well to different domains.

Gradient-based Coreset Methods. Gradient-based methods select a small weighted subset of data to approximate the gradients of the full dataset. CRAIG [40] achieves this by maximizing a submodular objective, while GRAD-MATCH [23] uses a greedy algorithm to match gradients directly. Extensions include Camel [33] for streaming data, and Goodcore [4], which handles incomplete relational data while jointly cleaning and training. Additionally, cluster-based methods [5] perform selection over feature-space clusters rather than individual points. Recently, TAGCOS [69] applied similar gradient clustering ideas to instruction tuning to create diverse datasets. Beyond direct gradient matching, coreset selection can also be formulated as a bi-level optimization problem. GLISTER [24] and RETRIEVE [25] adopt this strategy: the inner loop trains on a subset, while the outer loop greedily selects examples to minimize validation loss using hypergradient approximations.

3 PRELIMINARY

In this section, we first introduce the basic notation and recall supervised fine-tuning (SFT) and classical coreset selection, which our method builds upon. Table 1 summarizes all the symbols used throughout this paper.

Supervised Fine-tuning (SFT). Suppose that θ denotes the model parameters and D denotes the full training dataset. During the SFT stage, the full training dataset D with N data points can be represented as an instruction–response corpus $D = \{(x_i, y_i)\}_{i=1}^N$. For each data point, the instruction sequence is x_i and the response sequence is y_i , where $|y_i|$ denotes the length of the response sequence y_i (i.e., how many tokens it contains), y_i^t denotes the t -th token in y_i and $y_i^{<t} = (y_i^1, \dots, y_i^{t-1})$ denotes the prefix sequence before the

t -th token in sequence y_i . Thus, the SFT loss can be written as:

$$\mathcal{L}_{\text{SFT}}((x_i, y_i); \theta) = -\frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \log p_\theta(y_i^t | (x_i, y_i^{<t})). \quad (1)$$

Coreset. The state-of-the-art coreset selection framework [14, 44] aims to select a subset of data points $C \subseteq D$ (with non-negative weights $\{\omega_j\}_{j=1}^{|C|}$) such that the coreset gradient well approximates the full gradient, i.e., we keep only a smaller weighted subset of instruction–response pairs while still matching the training effect produced by the original dataset.

Intuitively, each data point (x_i, y_i) contributes a gradient vector $\nabla \mathcal{L}_i(\theta)$, and the goal is to select a smaller weighted subset whose combined gradients approximate those of the entire dataset. This naturally involves a trade-off between subset compactness and approximation fidelity, which we formalize through two equivalent formulations below.

The most direct way is to minimize the gradient approximation error (GA error) subject to a size constraint:

$$\min_{C \subseteq D, \omega_j \geq 0} \max_{\theta \in \Theta} \left\| \underbrace{\sum_{i=1}^{|D|} \nabla \mathcal{L}_i(\theta)}_{\text{full gradient}} - \underbrace{\sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\gamma(j)}(\theta)}_{\text{coreset gradient}} \right\| \quad \text{s.t.} \quad |C| \leq K \quad (2)$$

gradient approximation error

Here, K denotes the coreset size budget and Θ represents the feasible parameter space.

Equation (2) aims to minimize the GA error, which measures the gap between the full dataset’s gradient and the weighted coreset’s gradient. However, directly solving this formulation is computationally intractable. To overcome this difficulty, many studies [5, 40, 41, 51] adopt a dual formulation by interchanging the optimization objective and the constraint, as shown below:

$$\min_{C \subseteq D, \omega_j \geq 0} |C| \quad \text{s.t.} \quad \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_i(\theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\gamma(j)}(\theta) \right\| \leq \epsilon \quad (3)$$

Equation (3) captures the essential trade-off between approximation accuracy (ϵ) and subset size ($|C|$). Modern coreset methods [5, 40, 41, 51] typically solve the fixed-budget problem (Eq. (2)) by optimizing this dual formulation—finding the minimal subset for a specific error tolerance—and then adapting the solution to the budget K . We build on the same dual-formulation framework, but explicitly split the SFT GA error into knowledge-related and instruction following terms, letting these two parts guide the data selection process.

Example. Consider a toy dataset of 8 data points forming three groups with identical gradients (group sizes: 3, 3, and 2). To satisfy $K = 3$, we can simply pick one point from each group and weight them by their group size (i.e., weights 3, 3, 2). This achieves near-zero gradient error, demonstrating how the method removes data redundancy without losing the training effect.

4 CORESET SELECTION FOR SFT

In this section, we decompose the SFT loss into a knowledge-related term and an instruction following term (Section 4.1), reformulate coreset selection to preserve both (Section 4.2), and present the BRIEF framework and algorithmic details (Sections 4.3–5).

4.1 SFT Loss Decomposition

By adding $\log p_\theta(y^t | y^{<t})$, Eq. (1) can be written as

$$\mathcal{L}_{\text{SFT}}(y | x; \theta) = -\frac{1}{|y|} \sum_{t=1}^{|y|} \left[\log p_\theta(y^t | y^{<t}) + \log \frac{p_\theta(y^t | x, y^{<t})}{p_\theta(y^t | y^{<t})} \right]. \quad (4)$$

This leads to the decomposition:

$$\mathcal{L}_{\text{SFT}}(y | x; \theta) = \mathcal{L}_{\text{KN}}(y; \theta) + \mathcal{L}_{\text{IF}}(y | x; \theta). \quad (5)$$

where

$$\mathcal{L}_{\text{KN}}(y; \theta) = -\frac{1}{|y|} \sum_{t=1}^{|y|} \log p_\theta(y^t | y^{<t}), \quad (6)$$

and

$$\mathcal{L}_{\text{IF}}(y | x; \theta) = -\frac{1}{|y|} \sum_{t=1}^{|y|} \log \frac{p_\theta(y^t | x, y^{<t})}{p_\theta(y^t | y^{<t})}. \quad (7)$$

The first component $\mathcal{L}_{\text{KN}}(y; \theta)$ denotes **knowledge-related loss**, which is in the same format as the pretraining loss of LLMs [46]. It represents the negative log-likelihood of predicting the next token y^t given only the previous tokens $y^{<t}$. This component measures how well the model predicts the response tokens without using the instruction x . Therefore, it mainly reflects the knowledge inherently stored in the answers.

Next, we show that the second component represents the **instruction following loss**. It measures the log-probability ratio of generating y with instruction x versus without it. This captures the additional information provided by the instruction x that helps produce the correct response y . Therefore, it effectively quantifies to what extent the instruction improves the model's ability to generate the desired output.

To prove this, we first introduce a metric called *instruction following Difficulty (IFD)* [29–31], which, given an instruction (x, y) , identifies discrepancies between the expected responses of a model and its generation capability.

$$\text{IFD}(y | x; \theta) = \frac{\text{PPL}(y | x; \theta)}{\text{PPL}(y; \theta)}, \quad (8)$$

where

$$\text{PPL}(y | x; \theta) = \exp(\mathcal{L}_{\text{SFT}}((x, y); \theta)), \quad (9)$$

$$\text{PPL}(y; \theta) = \exp(\mathcal{L}_{\text{KN}}(y; \theta)). \quad (10)$$

Intuitively, IFD measures the potential of this instruction to improve the instruction following capability of a model. The key observation here is that taking the log on $\text{IFD}(y | x; \theta)$ will derive the exact form of $\mathcal{L}_{\text{IF}}(y | x; \theta)$, i.e., the second component in Eq. (4).

$$\mathcal{L}_{\text{IF}}(y | x; \theta) = \log \text{IFD}(y | x; \theta). \quad (11)$$

This shows that \mathcal{L}_{IF} primarily focuses on exploring the training examples to improve the instruction following capacity of a model, thus representing the *instruction following loss*.

Notably, \mathcal{L}_{IF} is derived directly from the SFT loss, namely the actual objective used during training. This connection provides a principled explanation for the effectiveness of *instruction following difficulty (IFD)* and, in turn, proves the critical role of training examples in improving the instruction following capability of a model.

In summary, this decomposition highlights that the SFT loss reflects two core model capabilities: knowledge-related capability, driven by \mathcal{L}_{KN} , and instruction following capability, driven by \mathcal{L}_{IF} . At the level of individual training examples, each instruction-response pair influences the model through a knowledge-related term and an instruction following term, which are exactly the components we later control when selecting and weighting examples in the coreset.

4.2 Problem Definition

In this section, we reformulate the coreset selection problem for instruction tuning in LLMs. Given a full training dataset $D = \{(x_i, y_i)\}_{i=1}^{|D|}$, our goal is to select a subset $C \subseteq D$ with weights $\{\omega_j\}_{j=1}^{|C|}$ that preserves both knowledge-related and instruction following capabilities of D simultaneously during the instruction tuning period.

SFT Gradient Approximation. Incorporating the SFT loss from Eq. (1) into the traditional coreset objective (Eq. (2)), we obtain:

$$\max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{SFT}}((x_i, y_i); \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{SFT}}((x_{\gamma(j)}, y_{\gamma(j)}); \theta) \right\| \quad (12)$$

subject to $|C| \leq K$, where $\gamma(j)$ denotes the index mapping from coreset to the full dataset.

Triangle Inequality Analysis. Building on the SFT loss decomposition established in Eq. (5), we apply the triangle inequality to the gradient approximation error in Eq. (12):

$$\begin{aligned} & \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{SFT}}((x_i, y_i); \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{SFT}}((x_{\gamma(j)}, y_{\gamma(j)}); \theta) \right\| \\ & \leq \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{KN}}(y_{\gamma(j)}; \theta) \right\| \\ & \quad + \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{IF}}(y_{\gamma(j)} | x_{\gamma(j)}; \theta) \right\|. \end{aligned} \quad (13)$$

which implies that the SFT gradient approximation error is upper bounded by the sum of the two component errors. However, minimizing the SFT GA error alone does not necessarily control the knowledge-related and instruction following errors separately. Hence, we explicitly constrain both components and reformulate the objective as follows.

$$\begin{aligned} & \min_{C \subseteq D, \omega_j \geq 0} \left[\underbrace{\max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{KN}}(y_{\gamma(j)}; \theta) \right\|}_{\text{Knowledge-related GA Error}} \right. \\ & \quad \left. + \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{IF}}(y_{\gamma(j)} | x_{\gamma(j)}; \theta) \right\| \right] \\ & \text{s.t. } |C| \leq K \end{aligned} \quad (14)$$

Reformulated Optimization Objective. The reformulation differs from traditional approaches by explicitly constraining both component errors. Specifically, through the triangle inequality, any solution to Eq. (14), where the total error does not exceed a certain

upper bound, guarantees that the SFT GA error is also within that bound. This reformulated objective ensures that both knowledge-related (KN) and instruction following (IF) capabilities are preserved in the selected coreset, addressing the unique requirements of SFT that traditional coreset methods overlook. In the following sections, we present our BRIEF framework to efficiently solve this problem.

4.3 BRIEF Overview

The BRIEF framework addresses the coreset selection problem (Eq. (14)) through a bi-level optimization approach that simultaneously retains both knowledge-related and instruction following capabilities.

As illustrated in Figure 2, BRIEF starts with a warm-up phase to compute gradient-based distances for both knowledge-related and instruction following. Next, it employs a bi-level optimization: the upper level adjusts a parameter α to explore trade-offs, while the lower level selects the coreset using a combined metric. This structure transforms the original NP-hard problem into two manageable sub-tasks. Intuitively, we decouple the search by selecting a coreset for a fixed α , and then updating α via a simple one-dimensional search. This allows BRIEF to efficiently identify and weight the most representative instruction–response pairs from the SFT corpus.

Warm-up Phase. BRIEF starts by training the base model on a subset $\mathcal{D}_{\text{warm}}$. Then, it calculates gradients for both knowledge-related ($\nabla \mathcal{L}_{KN}$) and instruction following ($\nabla \mathcal{L}_{IF}$) capabilities to compute pairwise distances between data points. This step organizes training examples in the gradient spaces, enabling precise measurement of data similarity and coverage for the subsequent selection.

Bi-level Optimization Algorithm. The core of BRIEF lies in its bi-level optimization structure (see Section 5 for details). Specifically, the upper level employs ternary search over the auxiliary variable $\alpha \in (0, 1)$, which serves as a partition parameter that divides the solution space into different regions based on capability trade-offs. For each partition defined by α , the lower level performs coreset selection by formulating it as a submodular facility location problem for efficient calculation. In each iteration, this means we first fix how much error budget to allocate to KN versus IF and then choose a coreset that best represents the full dataset under this split.

Upper Optimization. The upper optimization employs ternary search, which iteratively narrows the search interval by evaluating two internal points (m_1 and m_2) and discarding the subinterval with worse performance. The algorithm exploits the monotonicity at both extremes of the α spectrum: the objective decreases for small α (where KN constraints dominate) and increases for large α (where IF constraints dominate). This property ensures that at least one local minimum exists in the interval, and the ternary search converges to a locally optimal partition parameter α^* in $O(\log(1/\delta))$ iterations, where δ is the convergence threshold. Once α^* is determined, the final coreset can be selected through the lower-level optimization, thereby preserving both model KN and IF capabilities.

Lower Optimization. For each iteration of the lower optimization (Algorithm 2), given the partition parameter α from the upper level, we construct a composite distance metric $d^{\text{SFT}} = \frac{1}{\alpha} d^{KN} + \frac{1}{1-\alpha} d^{IF}$ that integrates both capability measures. The algorithm then formulates coreset selection as a submodular facility location problem, where the objective measures the total distance from all

Algorithm 1 BRIEF Framework

Input: Candidate data pool \mathcal{D} , selection ratio γ , base model θ

Output: Coreset $C \subseteq \mathcal{D}$, weights W

```

1:  $\theta \leftarrow \text{WARMUPTRAINING}(\mathcal{D}, \theta)$ 
2: Compute  $\nabla \mathcal{L}_{KN}(y_i; \theta), \nabla \mathcal{L}_{IF}(y_i|x_i; \theta)$  for all  $(x_i, y_i) \in \mathcal{D}$ 
3: Compute pairwise distances  $d_{ij}^{KN}, d_{ij}^{IF}$  for all  $i, j \in [N]$ 
4:  $l \leftarrow 0, r \leftarrow 1$ 
5: while  $|r - l| > \delta$  do ▷ ternary search over  $\alpha, \delta \rightarrow 0$ 
6:    $m_1 \leftarrow l + \frac{r-l}{3}, m_2 \leftarrow r - \frac{r-l}{3}$ 
7:    $(C_1, W_1) \leftarrow \text{GREEDYCORESETSELECTION}(\mathcal{D}, \gamma, m_1, d^{KN}, d^{IF})$ 
8:    $(C_2, W_2) \leftarrow \text{GREEDYCORESETSELECTION}(\mathcal{D}, \gamma, m_2, d^{KN}, d^{IF})$ 
9:    $E_1 \leftarrow \text{COMPUTEGAERRORBOUND}(C_1, W_1)$ 
10:   $E_2 \leftarrow \text{COMPUTEGAERRORBOUND}(C_2, W_2)$ 
11:  if  $E_1 \leq E_2$  then
12:     $r \leftarrow m_2$ 
13:  else
14:     $l \leftarrow m_1$ 
15:  end if
16: end while
17:  $\alpha^* \leftarrow \frac{l+r}{2}$ 
18: return  $\text{GREEDYCORESETSELECTION}(\mathcal{D}, \gamma, \alpha^*, d^{KN}, d^{IF})$ 

```

data points to their nearest representatives in the coreset. Through greedy selection, it iteratively chooses $K = \lceil \gamma |\mathcal{D}| \rceil$ data points that maximally reduce the sum of distances from unselected points to their nearest selected representatives. This greedy approach provides a $(1 - 1/e)$ -approximation guarantee for the submodular maximization problem [47]. Finally, it computes weights for each selected point based on the number of data points assigned to it.

5 THE BRIEF APPROACH

In this section, we present the technical details that solve the optimization problem formulated in Eq. (14). This problem is proven to be prohibitively expensive (Section 5.1). We first demonstrate that directly optimizing this problem is computationally prohibitive, while transforming it into a bi-level optimization structure makes it possible to design an efficient solution (Section 5.2). We then detail the lower-level optimization, which solves the coreset selection problem w.r.t. a specific auxiliary variable α that partitions the solution space by allocating error budgets between the two capability components (Section 5.3). Finally, we present the upper-level optimization that efficiently searches for a locally optimal α value by exploiting the monotonic properties of the objective function (Section 5.4).

5.1 Problem Complexity

We now establish the computational intractability of the optimization problem in Eq. (14).

THEOREM 5.1. *The optimization problem in Eq. (14) is NP-hard.*

PROOF. We prove the NP-hardness by reduction from the classical coreset selection problem for gradient approximation, which is known to be NP-hard [40]. Specifically, we show that Eq. (14) subsumes the classical gradient-based coreset selection as a special case when instructions are set to be empty. The detailed proof is provided in our technical report [54]. \square

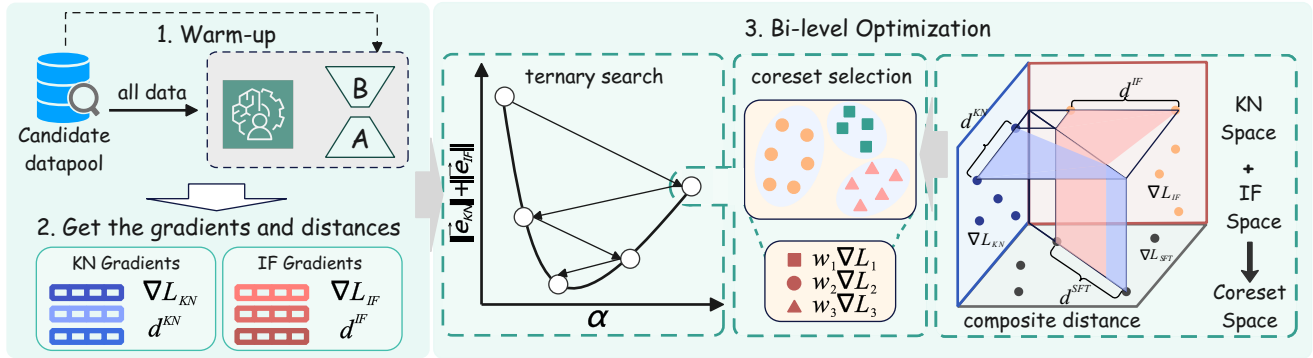


Figure 2: The Overall Framework of BRIEF.

5.2 Bi-level Optimization Problem

To efficiently solve the NP-hard problem formulated in Eq. (14), we transform it into a bi-level optimization framework. Following popular coreset selection approaches [40, 41, 51], we first derive the dual formulation of our coreset selection problem and then establish its equivalence to the original objective:

$$\begin{aligned}
 & \min_{C \subseteq D, \omega_j \geq 0} |C| \\
 \text{s.t.} \quad & \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{KN}}(y_{Y(j)}; \theta) \right\| \\
 & + \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{IF}}(y_{Y(j)} | x_{Y(j)}; \theta) \right\| \leq \varepsilon
 \end{aligned} \tag{15}$$

where ε represents the budget of the overall GA error.

Specifically, Eq.(15) serves as the dual formulation of Eq.(14), analogous to the relationship between Eqs.(2) and (3). We adopt this tolerance-constrained view to facilitate our analysis, while the derived algorithm effectively addresses the original size-constrained optimization. The challenge of directly solving Eq. (15) lies in the coupled constraint that involves the sum of two gradient approximation errors. To decouple this constraint and enable more efficient optimization, we introduce an auxiliary variable $\alpha \in (0, 1)$ that partitions the error budget between the KN and IF capabilities. Importantly, α is not a hyperparameter requiring manual tuning but is automatically optimized through our bi-level framework, as detailed below.

Specifically, we can transform our problem into the following equivalent bi-level optimization formulation:

$$\begin{aligned}
 & \min_{\alpha \in (0,1)} \min_{C \subseteq D, \omega_j \geq 0} |C| \quad \text{s.t.} \\
 & \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{KN}}(y_{Y(j)}; \theta) \right\| \leq \alpha \varepsilon \\
 & \max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \sum_{j=1}^{|C|} \omega_j \nabla \mathcal{L}_{\text{IF}}(y_{Y(j)} | x_{Y(j)}; \theta) \right\| \leq (1 - \alpha) \varepsilon
 \end{aligned} \tag{16}$$

We now show the equivalence of Eq. (15) and Eq. (16).

THEOREM 5.2. *The bi-level optimization problem in Eq. (16) is equivalent to the formulation in Eq. (15). In particular, any optimal*

solution (C^, W^*) to Eq. (15) corresponds to an optimal solution of Eq. (16) for some $\alpha^* \in (0, 1)$, and vice versa.*

PROOF. We provide the full proof in our technical report [54] and offer a simplified intuition here. The equivalence holds because the variable α simply acts as a slider that partitions the total error budget ε between the two capabilities. Specifically, any solution satisfying the combined constraint in Eq. (15) naturally corresponds to a specific split α in Eq. (16). Conversely, summing the two split constraints in Eq. (16) recovers the original total budget ε . Since both formulations minimize the same subset size $|C|$ under the same total error allowance, they yield identical optimal solutions. \square

This bi-level structure offers significant computational advantages. At the upper level, we only tune the continuous parameter α , which splits the GA error budget between KN and IF. Because the objective is unimodal in α , we can search for a good value using ternary search. For each fixed α , the lower level reduces to a standard submodular facility-location coreset selection problem with approximation guarantees.

5.3 Lower-level Optimization Problem

Given a specific parameter α learned by the upper-level optimization, the lower-level problem focuses on finding the minimal coreset C and corresponding weights W that satisfy the decoupled gradient approximation constraints in Eq. (16).

This lower-level problem represents a constrained coreset selection task, where we must simultaneously bound the gradient approximation errors for both knowledge-related capability (with error budget $\alpha\varepsilon$) and instruction following capability (with error budget $(1 - \alpha)\varepsilon$). While the parameter α partitions the total error budget between the two capabilities, solving this optimization problem remains computationally challenging due to its combinatorial nature and the need to evaluate gradient approximations over the parameter space Θ . However, we show that the aforementioned dual-constraint optimization problem can be reformulated as a submodular set cover problem, which admits efficient approximation algorithms.

Formulation as a Submodular Set Cover Problem. To precisely formulate our coreset selection task as a submodular set cover problem,

we define deterministic bounds $\mathcal{B}_{\text{KN}}(C)$ and $\mathcal{B}_{\text{IF}}(C)$ for any subset $C \subseteq D$ and parameter $\theta \in \Theta$.

$$\mathcal{B}_{\text{KN}}(C) \stackrel{\text{def}}{=} \sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{KN}}, \quad \mathcal{B}_{\text{IF}}(C) \stackrel{\text{def}}{=} \sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{IF}}. \quad (17)$$

Here, pairwise distances d_{ij}^{KN} and d_{ij}^{IF} are used to measure the normed difference between the gradient of data point $c_i \in D$ and data point $c_j \in C$, which are defined separately as:

$$d_{ij}^{\text{KN}} \stackrel{\text{def}}{=} \max_{\theta \in \Theta} \|\nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \nabla \mathcal{L}_{\text{KN}}(y_j; \theta)\| \quad (18)$$

$$d_{ij}^{\text{IF}} \stackrel{\text{def}}{=} \max_{\theta \in \Theta} \|\nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \nabla \mathcal{L}_{\text{IF}}(y_j | x_j; \theta)\| \quad (19)$$

Following the theoretical framework established in Mirzasoleiman et al. [40], we can derive upper bounds for the gradient approximation errors. Specifically, the maximum gradient approximation error for any subset C can be bounded by the sum of minimum pairwise distances between the gradient of each point in D and its closest representative in C . Intuitively, if every point in D has a close neighbour in C in the gradient space, then the sum of gradients over D can be well-approximated by reweighting the gradients of points in C ; the quantities $\mathcal{B}_{\text{KN}}(C)$ and $\mathcal{B}_{\text{IF}}(C)$ measure how well C covers D under the KN and IF gradient distances.

Formally, this gives us:

$$\max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{KN}}(y_i; \theta) - \sum_{j=1}^{|C^*|} \omega_j^* \nabla \mathcal{L}_{\text{KN}}(y_{y^*(j)}; \theta) \right\| \leq \mathcal{B}_{\text{KN}}(C), \quad (20)$$

$$\max_{\theta \in \Theta} \left\| \sum_{i=1}^{|D|} \nabla \mathcal{L}_{\text{IF}}(y_i | x_i; \theta) - \sum_{j=1}^{|C^*|} \omega_j^* \nabla \mathcal{L}_{\text{IF}}(y_{y^*(j)} | x_{y^*(j)}; \theta) \right\| \leq \mathcal{B}_{\text{IF}}(C). \quad (21)$$

This allows us to transform the intractable optimization problem in Eq. (16) into a more manageable form with scalar constraints.

Substituting $\mathcal{B}_{\text{KN}}(C)$ and $\mathcal{B}_{\text{IF}}(C)$ into Eq. (16), we obtain a simplified optimization problem with scalar constraints:

$$\min_{C \subseteq D} |C| \quad \text{s.t.} \quad \mathcal{B}_{\text{KN}}(C) \leq \alpha \varepsilon, \quad \mathcal{B}_{\text{IF}}(C) \leq (1 - \alpha) \varepsilon. \quad (22)$$

Unifying Dual Constraints into a Single Constraint. While the dual constraints in Eq. (22) provide explicit control over both capabilities, solving this bi-constraint optimization problem remains computationally challenging. To address this, we propose to transform the dual constraints into a single unified constraint that can be efficiently solved using existing submodular optimization techniques.

Our key insight is to construct a composite distance metric that effectively combines the KN and IF distance components, which is defined as follows:

$$d_{ij}^{\text{SFT}} \stackrel{\text{def}}{=} \frac{1}{\alpha} d_{ij}^{\text{KN}} + \frac{1}{1 - \alpha} d_{ij}^{\text{IF}}. \quad (23)$$

By weighting d_{ij}^{KN} with $\frac{1}{\alpha}$ and d_{ij}^{IF} with $\frac{1}{1 - \alpha}$, we make each component's contribution inversely proportional to its allocated error budget. We now show that any coreset satisfying a constraint on d_{ij}^{SFT} automatically satisfies both original constraints.

THEOREM 5.3. *Let $C \subseteq D$ be a coreset and d_{ij}^{SFT} be defined as in Eq. (23). If $\sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{SFT}} \leq \varepsilon$, then C satisfies both constraints in Eq. (22):*

$$\mathcal{B}_{\text{KN}}(C) \leq \alpha \varepsilon, \quad \mathcal{B}_{\text{IF}}(C) \leq (1 - \alpha) \varepsilon. \quad (24)$$

PROOF. By construction, d_{ij}^{SFT} is a single budget that is split between d_{ij}^{KN} and d_{ij}^{IF} according to the ratios α and $1 - \alpha$. The scaling by $1/\alpha$ and $1/(1 - \alpha)$ ensures that each component can never exceed its own share of the composite budget, i.e., $d_{ij}^{\text{KN}} \leq \alpha d_{ij}^{\text{SFT}}$ and $d_{ij}^{\text{IF}} \leq (1 - \alpha) d_{ij}^{\text{SFT}}$. Summing over all data points immediately yields $\mathcal{B}_{\text{KN}}(C) \leq \alpha \varepsilon$ and $\mathcal{B}_{\text{IF}}(C) \leq (1 - \alpha) \varepsilon$ whenever $\sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{SFT}} \leq \varepsilon$. A detailed proof is provided in our technical report [54]. \square

Therefore, the scaling by $1/\alpha$ and $1/(1 - \alpha)$ directly reflects the allocated error budgets, so components with tighter budgets induce larger distances and are prioritized by the selection.

This theoretical result allows us to simplify the original dual constraint optimization problem (Eq. (22)) into a single scalar constraint:

$$\min_{C \subseteq D} |C| \quad \text{s.t.} \quad \sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{SFT}} \leq \varepsilon, \quad (25)$$

Since d_{ij}^{SFT} is a non-negative linear combination of the distances d_{ij}^{KN} and d_{ij}^{IF} , and submodular functions are closed under non-negative linear combinations [15, 26], the function $\sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{SFT}}$ inherits the submodularity property. Previous work [40, 41, 51] has demonstrated that the optimization problems of the form in Eq. (25), where the objective involves minimizing sums of minimum distances, can be equivalently formulated as submodular set cover problems, allowing the use of efficient approximation algorithms. Formulation with Fixed Coreset Size. While Eq. (25) provides a theoretically elegant formulation with an error bound ε , in practice, directly specifying such a bound is challenging without prior knowledge of the dataset's gradient structure. Following standard practice in coreset selection literature [40, 41] and consistent with the dual formulation in Eq. (2), we adopt a more practical approach: we minimize the gradient approximation error with a coreset size constraint K . Specifically, our practical optimization objective becomes:

$$\min_{C \subseteq D} \sum_{i=1}^{|D|} \min_{c_j \in C} d_{ij}^{\text{SFT}} \quad \text{s.t.} \quad |C| \leq K \quad (26)$$

This formulation transforms the problem into a submodular maximization problem through the facility location framework. It is equivalent to solving the constraint problem in Eq. (25) for an appropriately chosen ε value, but eliminates the need to specify ε a priori. The greedy algorithm naturally adapts to minimize the total gradient approximation error given the budget of K data points, making it more practical for real-world applications where the appropriate error tolerance is unknown beforehand.

Solving the Unified Constraint via Greedy Framework. We solve the unified optimization problem in Eq. (26) by adapting the standard facility location framework [40]. While the greedy optimization strategy is well-established, our approach uniquely incorporates the composite distance metric d_{ij}^{SFT} to balance KN and IF budgets. This formulation allows us to leverage the submodular property of the objective, theoretically guaranteeing a $(1 - \frac{1}{e})$ -approximation of the optimal solution (Algorithm 2).

The algorithm begins by computing the composite distance metric $d_{ij}^{\text{SFT}} = \frac{1}{\alpha} d_{ij}^{\text{KN}} + \frac{1}{1 - \alpha} d_{ij}^{\text{IF}}$ for all pairs of data points (Line 3). Following Eq. (26), the optimization objective $\mathcal{B}(S)$ is defined (Line 4), incorporating an auxiliary element s_0 to implicitly transform

Algorithm 2 GreedyCoresetSelection

```
1: function GREEDYCORESETSELECTION( $\mathcal{D}, \gamma, \alpha, d^{\text{KN}}, d^{\text{IF}}$ )
2:    $k \leftarrow \lceil \gamma |\mathcal{D}| \rceil, C \leftarrow \emptyset$ 
3:    $d_{ij}^{\text{SFT}} \leftarrow \frac{1}{\alpha} d_{ij}^{\text{KN}} + \frac{1}{1-\alpha} d_{ij}^{\text{IF}}$ 
4:   Let  $\mathcal{B}(S) = \sum_{i \in \mathcal{D}} \min_{c \in S \cup \{s_0\}} d_{ic}^{\text{SFT}}$  ▷ with auxiliary  $s_0$ 
5:   while  $|C| < k$  do
6:      $u^* \leftarrow \arg \max_{u \in \mathcal{D} \setminus C} (\mathcal{B}(C) - \mathcal{B}(C \cup \{u\}))$ 
7:      $C \leftarrow C \cup \{u^*\}$ 
8:   end while
9:    $w_j \leftarrow |\{i \in \mathcal{D} : j = \arg \min_{c \in C} d_{ic}^{\text{SFT}}\}|$  for all  $j \in C$ 
10:  return  $C, W = \{w_j\}_{j \in C}$ 
11: end function
```

the minimization problem into a maximization one. This objective represents the total gradient approximation error, measuring how well the coreset covers the full dataset \mathcal{D} .

The core of Algorithm 2 is the iterative greedy selection process (Lines 5-8). At each iteration, the algorithm directly identifies the candidate data point u^* that maximizes the marginal gain in the objective function (Line 6). This step effectively selects the point that maximally reduces the gradient approximation error. The selected point is then added to C (Line 7). This greedy strategy ensures a compact yet representative coreset.

After selecting K data points, the algorithm calculates their weights (Line 9) by assigning each data point in \mathcal{D} to its nearest representative in C under the composite distance metric. The weight w_j for each coreset instance is equal to the number of mapped data points. In practice, each greedy step picks the instruction–response pair that most reduces the remaining composite distance for other examples, so the algorithm naturally focuses on regions of the SFT corpus that are still poorly covered.

5.4 Upper-level Optimization Problem

The upper-level optimization in Eq. (16) seeks the optimal partition parameter $\alpha \in (0, 1)$ that minimizes the coreset size. Different values of α decide which aspect of the selected training examples is emphasized, that is, whether their contribution to KN or to IF matters more. We establish the existence of at least one local minimum through upper and lower bounds, which enables efficient optimization via ternary search.

Notation and Setup. For any subset $C \subseteq D$, recall the error bounds defined in Eq. (17). For a given error budget $\tau \geq 0$, we define the single-constraint minimum size functions:

$$K_{KN}(\tau) := \min\{|C| : C \subseteq D, \mathcal{B}_{KN}(C) \leq \tau\}, \quad (27)$$

$$K_{IF}(\tau) := \min\{|C| : C \subseteq D, \mathcal{B}_{IF}(C) \leq \tau\}. \quad (28)$$

These functions are monotonically non-increasing in τ since relaxing the error constraint cannot increase the minimum coreset size.

The dual-constraint minimum size function, corresponding to the lower-level optimal value in Eq. (16), is:

$$K^*(\alpha) := \min\{|C| : C \subseteq D, \mathcal{B}_{KN}(C) \leq \alpha\epsilon, \mathcal{B}_{IF}(C) \leq (1-\alpha)\epsilon\}. \quad (29)$$

Lower Bound: Unimodal Structure. We first establish a lower bound for $K^*(\alpha)$ by observing that any feasible coreset must satisfy both constraints independently:

THEOREM 5.4 (LOWER BOUND). For any $\alpha \in (0, 1)$:

$$K^*(\alpha) \geq \max\{K_{KN}(\alpha\epsilon), K_{IF}((1-\alpha)\epsilon)\}. \quad (30)$$

PROOF. Let C^* be an optimal coreset for $K^*(\alpha)$. Since C^* must satisfy $\mathcal{B}_{KN}(C^*) \leq \alpha\epsilon$, we have $|C^*| \geq K_{KN}(\alpha\epsilon)$ by the definition of K_{KN} . Similarly, from $\mathcal{B}_{IF}(C^*) \leq (1-\alpha)\epsilon$, we get $|C^*| \geq K_{IF}((1-\alpha)\epsilon)$. Thus, $K^*(\alpha) = |C^*| \geq \max\{K_{KN}(\alpha\epsilon), K_{IF}((1-\alpha)\epsilon)\}$. \square

This lower bound exhibits a unimodal structure: $K_{KN}(\alpha\epsilon)$ is non-increasing in α (as $\alpha\epsilon$ increases), while $K_{IF}((1-\alpha)\epsilon)$ is non-decreasing in α (as $(1-\alpha)\epsilon$ decreases). The maximum of these two monotonic functions with opposite trends creates a unimodal lower bound with at least one minimum point.

Upper Bound: U-shaped Structure. We construct an upper bound through a union-based approach:

THEOREM 5.5 (UPPER BOUND). For any $\alpha \in (0, 1)$:

$$K^*(\alpha) \leq K_{KN}(\alpha\epsilon) + K_{IF}((1-\alpha)\epsilon). \quad (31)$$

PROOF. Let C_{KN} be an optimal coreset of size $K_{KN}(\alpha\epsilon)$ satisfying $\mathcal{B}_{KN}(C_{KN}) \leq \alpha\epsilon$, and C_{IF} be an optimal coreset of size $K_{IF}((1-\alpha)\epsilon)$ satisfying $\mathcal{B}_{IF}(C_{IF}) \leq (1-\alpha)\epsilon$. The union $C = C_{KN} \cup C_{IF}$ satisfies both constraints:

$$\mathcal{B}_{KN}(C) \leq \mathcal{B}_{KN}(C_{KN}) \leq \alpha\epsilon, \quad (32)$$

$$\mathcal{B}_{IF}(C) \leq \mathcal{B}_{IF}(C_{IF}) \leq (1-\alpha)\epsilon, \quad (33)$$

where the first inequalities follow from the monotonicity of error bounds with respect to set inclusion. Thus, C is feasible for $K^*(\alpha)$, yielding $K^*(\alpha) \leq |C| \leq |C_{KN}| + |C_{IF}| = K_{KN}(\alpha\epsilon) + K_{IF}((1-\alpha)\epsilon)$. \square

This upper bound exhibits U-shaped behavior: as $\alpha \rightarrow 0$, we have $K_{KN}(\alpha\epsilon) \rightarrow |D|$ (requiring nearly all data to satisfy a stringent KN constraint), while $K_{IF}((1-\alpha)\epsilon)$ remains bounded. Similarly, as $\alpha \rightarrow 1$, $K_{IF}((1-\alpha)\epsilon) \rightarrow |D|$ while $K_{KN}(\alpha\epsilon)$ remains bounded. Thus, the upper bound is large at both extremes.

Efficient Optimization via Ternary Search. The proven bounds and their monotonic properties enable efficient optimization of the upper-level problem via ternary search. The algorithm iteratively narrows the search interval by evaluating two internal points and discarding the subinterval with worse performance. This divide-and-conquer approach converges to a local minimum in $O(\log(1/\delta))$ iterations, where δ is the convergence threshold.

Equivalence Between Subset Size and GA Error Minimization. The dual formulation in Eq. (15) (minimizing coreset size subject to error constraint) and the original objective in Eq. (14) (minimizing total GA error subject to size constraint) are equivalent through duality in constrained optimization [40]. This duality relationship ensures that solving either formulation yields the same optimal trade-off between coreset size and gradient approximation error, allowing us to choose the more convenient formulation for our bi-level optimization framework.

Time Complexity. Our bi-level optimization framework achieves efficient computation through its hierarchical structure. The inner loop (greedy coreset selection) has the same time complexity as the traditional coreset selection algorithm [40], requiring $O(|D|^2 \cdot K)$ operations to compute pairwise distances and perform greedy selection. The outer loop employs ternary search over α , which

converges to a locally optimal α^* in $O(\log(1/\delta))$ iterations where δ is the convergence threshold. Therefore, the overall time complexity of our method is $O(\log(1/\delta) \cdot |D|^2 \cdot K)$.

6 EXPERIMENT

In this section, we use 4 massive datasets to fine-tune 3 popular base models to demonstrate the efficiency and effectiveness of BRIEF. We organize the results around five research questions:

RQ1 (Overall performance). How does BRIEF compare with other baselines in accuracy under various selection ratios?

RQ2 (Generalization). Does BRIEF generalize across base models?

RQ3 (Efficiency). How efficient is BRIEF compared with baselines?

RQ4 (Mechanism). Why does BRIEF work—what do the validation of the automatically selected auxiliary variable α and the loss decomposition analysis reveal?

RQ5 (Hyperparameter). What is the impact of the hyperparameter on BRIEF?

Due to space constraints, additional mechanism analyses, such as the investigation of error-vector angles and gradient conflict visualizations, are detailed in our technical report [54].

6.1 Experiment Setup

Training Settings. We evaluate our method on three widely used foundation models (i.e., Llama-3.1-8B [11], Qwen3-4B [64], and Mistral-Nemo-12B [42]). Mistral-Nemo-12B is a 12-billion parameter model jointly developed by Mistral AI and NVIDIA. Specifically, we train each model for 4 epochs with a maximum token length of 2048 on 8 A800 GPUs. Following the setting of Tulu-3-SFT-Mixture [27], we employ AdamW optimizer and a maximum learning rate of $5e-6$. A linear learning rate scheduler with 0.03 warm-up ratio is used and the batch size is set to 128.

Data Selection Settings. Following LESS [62], we first conduct a warm-up phase by training the model for one epoch on a randomly sampled 5% subset from the candidate data pool D , allowing the model to capture the data distribution of the candidate dataset. Then, this warmed-up model is used as the basis for gradient computation during the data selection period. For methods requiring a reference set, following previous work [61], we sample 5 examples from the training set of MMLU dataset [20] to serve as the reference set.

Datasets. We use four popular datasets as candidate data pools D for our experiments.

(1) Tulu-3-SFT-Mixture [27] comprises a diverse mixture of approximately 900k instruction following examples drawn from various public sources (e.g., CoCoNot [12], FLAN v2 [36], etc.), which emphasizes high-quality responses in reasoning, general knowledge, and multilingual tasks.

(2) Magpie-Pro-10K-GPT4o-mini [45] contains 10k professionally curated instruction–response pairs generated with GPT-4o-mini. This dataset emphasizes high-quality, diverse data covering reasoning, commonsense and practical application scenarios.

(3) Code-74k-ShareGPT [3] is a code-centric dataset containing approximately 74k examples distilled from ShareGPT interactions, which focuses on programming-related queries.

(4) MegaScience [13] is a comprehensive scientific SFT candidate pool covering diverse fields such as math, medicine, biology, physics, and chemistry.

Baselines. We compare BRIEF with several baselines that are classified into two categories according to whether they rely on a reference set.

Reference-based methods: These methods select data points with a distribution similar to a reference set (e.g., MMLU).

(1) DSIR [63] assigns weights to the candidate training data points based on the similarity of n -gram features between the training data points and the reference set, and then selects top $k\%$ data points according to these estimated weights.

(2) BM25 [52] selects top $k\%$ data points by TF-IDF score, which measures the word similarity between the training data points and the reference set.

(3) LESS [62] selects $k\%$ training data points that have the highest influence scores between the reference set and the data points in the candidate pool.

Reference-free methods: These methods aim to enhance the general capabilities of LLMs (e.g., selecting a set of data points without relying on a reference set and evaluating it on multiple different downstream tasks).

(1) Random randomly samples $k\%$ of the data points from D , which are then used for fine-tuning.

(2) Total fine-tunes the target model using all training data points in the candidate data pool D , i.e., full-dataset training.

(3) IFD [30] selects $k\%$ of the data points that pose greater instruction following difficulty for the target model based on the IFD metric.

(4) TAGCOS [69] selects $k\%$ data points from a candidate data pool D by clustering gradients to identify the coreset.

(5) DELIFT [2] is a recent data-efficient instruction fine-tuning method that computes a pairwise utility score between training examples and then applies submodular optimization to select a diverse and informative $k\%$ subset. We follow its Instruction Tuning setting, which requires no reference set, and compute the utilities using the training model.

(6) CRAIG [40] reformulates the problem of minimizing gradient approximation (GA) error as a submodular optimization problem and uses a greedy algorithm to select the $k\%$ data points as the coreset.

(7) GRAD-MATCH [23] selects a coreset of $k\%$ data points whose weighted gradients match those of the full training set, and optimizes this gradient-matching objective with an orthogonal matching pursuit (OMP) algorithm.

(8) BRIEF is our solution.

Evaluation Datasets. To comprehensively evaluate the capabilities of fine-tuned models, we evaluate them on various downstream tasks covering the following categories:

(1) **General Tasks:** MMLU-PRO [57] evaluates broad general knowledge across 14 domains; MMLU [20] evaluates multidisciplinary knowledge with multiple-choice questions across 57 subjects; DROP [10] tests reading comprehension with discrete reasoning; AGIEval [73] benchmarks standardized exams and math-cloze tasks; KorBench [37] is about diverse reasoning challenges across symbolic and logical problem solving.

(2) **Mathematical Tasks:** GSM8K [8] focuses on grade school math problems, while MATH [21] covers competition-level mathematics.

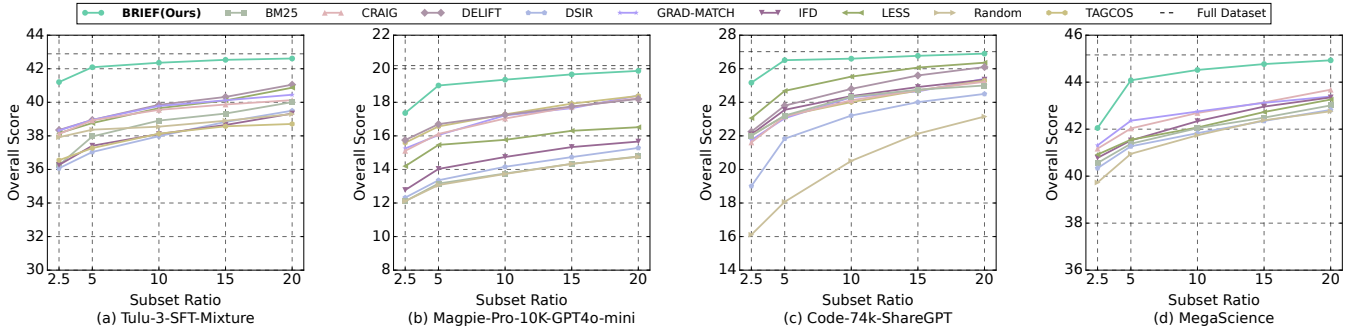


Figure 3: Overall Performance Comparison with Baselines.

(3) *Coding Tasks*: HumanEval [7] evaluates code synthesis, whereas LCB [48] spans code generation and test output prediction.

Evaluation Metrics. Evaluations are conducted using the OpenCompass [9] framework and the overall score is reported for comparison. For coding tasks (i.e., HumanEval and LCB), we use pass@1, which represents the fraction of problems solved correctly on the first attempt. For other tasks, we report accuracy as the primary metric. We also report ExaFLOPs

6.2 RQ1: Overall Performance

Figure 3 compares BRIEF with other baselines using Llama-3.1-8B on average accuracy (Y-axes), with varying data selection ratios: 2.5%, 5%, 10%, 15% and 20% (X-axes). In summary, BRIEF outperforms other baselines on all candidate data pools at different data selection ratios.

For reference-based methods (i.e., BM25, DSIR and LESS), they primarily select data associated with a reference set (MMLU in our experiment setting). This causes the trained model to excel only on specific downstream tasks. However, such data selection undermines the model’s generalization ability to other downstream tasks, leading to performance degradation under diverse evaluation scenarios. For example, under the same 10% data selection ratio, BRIEF outperforms LESS by 3.59% on the Magpie-Pro-10K-GPT4o-mini dataset, highlighting its superior effectiveness in data selection.

For reference-free methods, BRIEF achieves a 3.722% accuracy improvement over random selection on the Tulu-3-SFT-Mixture dataset at a 5% selection ratio. Additionally, unlike IFD which overlooks knowledge capabilities, BRIEF optimizes them jointly. This leads to a 3.08% improvement over IFD on Code-74k-ShareGPT at a 2.5% selection ratio. Similarly, BRIEF outperforms DELIFT by approximately 3.14% on Tulu-3-SFT-Mixture at a 5% selection ratio, as DELIFT only considers pairwise relationships based on loss, rather than the relationship between the subset and the whole dataset. Finally, while coreset methods (i.e., CRAIG, GRAD-MATCH, and TAGCOS) are strong baselines, they rely on raw SFT gradients that cause interference between capabilities. BRIEF addresses this by performing fine-grained minimization of KN and IF error vectors instead of just the overall SFT error. Consequently, BRIEF outperforms CRAIG by approximately 3.165% and GRAD-MATCH by approximately 2.99% on Tulu-3-SFT-Mixture at a 5% selection ratio. The same trend is observed on MegaScience, demonstrating

that our method is robust to different domains and maintains good performance.

6.3 RQ2: Generalization Across Models

In this section, we report the performance of BRIEF across 9 downstream tasks under different base models. As shown in Table 2, we fine-tune Llama-3.1-8B, Qwen3-4B and Mistral-Nemo-12B using 5% data points selected from Tulu-3-SFT-Mixture. Specifically, BRIEF exhibits superior accuracy compared to all baselines like GRAD-MATCH, IFD, LESS, CRAIG, DELIFT and TAGCOS. Notably, BRIEF shows a 2.18% accuracy improvement over IFD, a leading baseline at Qwen3-4B. We observe that reference-based methods like LESS or BM25 sometimes perform better on MMLU and MMLU-Pro (e.g., LESS outperforms our method by 3.03% on MMLU with Llama-3.1-8B), which is attributed to their use of reference sets drawn from these benchmarks. Nevertheless, our method achieves the best Overall Score across different model families and sizes, demonstrating that it is robust to model architecture and scale.

6.4 RQ3: Efficiency

To assess the practical feasibility of our approach, we compare the total computational costs required to achieve the same accuracy across three phases: the warm-up phase, the data selection phase, and the training phase. As shown in Figure 4, BRIEF achieves the lowest EFLOPs among all baseline methods while attaining performance comparable to that of using the full candidate data pool, reducing the total computational cost by 3 \times . Meanwhile, Table 3 presents an analysis of the time complexity of BRIEF compared with other coreset methods (e.g., CRAIG and TAGCOS), along with a comparison of the actual time consumption. We note that BRIEF and other coreset methods exhibit similar time complexity. However, when selecting coresets for comparable performance, BRIEF significantly improves the model’s efficiency by 2–3 \times .

6.5 RQ4: Mechanism Analyses

Effectiveness of α Selection To further validate the effectiveness of our method for choosing α , we enumerate various values of α and evaluate the performance of the coresets obtained from the inner optimization corresponding to each α (shown in Figure 5). Notably, $\alpha = 0.612$ is the value automatically determined by the

Table 2: Comprehensive evaluation of Llama-3.1-8B, Qwen3-4B and Mistral-Nemo-12B on 9 downstream tasks. Results are reported for various data selection methods at a 5% sampling ratio. The best scores are highlighted in bold and the second-best scores are underlined. We run each experiment three times and present the averaged results.

Method (5%)	DROP	AGIEval	General			Mathematical		Coding		Overall
			MMLU-PRO	MMLU	KorBench	GSM8K	MATH	HumanEval	LCB	
<i>Llama-3.1-8B</i>										
DSIR	48.00	39.20	33.79	61.67	22.80	54.66	23.36	35.37	14.52	37.041
BM25	46.81	36.13	35.29	61.19	21.76	64.44	27.22	32.93	15.90	37.964
LESS	47.13	35.56	<u>34.05</u>	62.06	19.68	<u>61.90</u>	23.26	49.39	15.84	38.763
IFD	49.02	35.52	32.44	51.90	21.36	61.25	25.88	44.15	15.20	37.413
DELIFT	50.89	38.66	32.31	51.51	23.12	62.42	26.38	49.49	15.77	38.950
TAGCOS	47.41	37.64	30.39	42.12	22.60	63.84	25.18	50.00	16.26	37.272
GRAD-MATCH	52.01	38.53	31.57	53.11	23.03	62.56	25.33	<u>49.25</u>	<u>15.18</u>	<u>38.952</u>
CRAIG	<u>57.66</u>	37.95	31.34	51.96	21.12	61.64	24.72	49.39	14.55	38.925
Random	51.03	38.84	33.71	51.09	<u>23.36</u>	60.27	26.54	44.51	15.97	38.368
BRIEF (Ours)	60.09	39.83	33.92	59.03	24.40	66.41	<u>27.05</u>	51.61	16.48	42.090
<i>Qwen3-4B</i>										
DSIR	62.88	47.39	39.66	60.05	46.56	79.87	57.16	72.93	36.58	55.898
BM25	68.92	48.06	<u>56.27</u>	76.37	51.52	79.49	<u>57.26</u>	73.05	37.97	60.990
LESS	75.46	48.17	<u>53.72</u>	<u>72.33</u>	49.28	79.27	<u>56.48</u>	69.39	41.18	60.586
IFD	76.13	47.97	55.49	<u>67.98</u>	<u>51.97</u>	80.87	57.51	71.33	41.51	<u>61.201</u>
DELIFT	76.13	<u>48.22</u>	56.02	68.17	<u>50.21</u>	<u>79.98</u>	55.81	71.29	41.23	<u>60.784</u>
TAGCOS	76.68	47.79	55.71	67.75	50.60	77.16	55.06	70.05	41.03	60.203
GRAD-MATCH	<u>77.05</u>	47.64	56.14	67.52	50.94	78.16	55.76	71.33	<u>42.01</u>	60.728
CRAIG	76.38	47.52	55.19	68.50	51.00	78.37	56.27	72.51	<u>40.33</u>	60.674
Random	75.21	47.09	51.99	67.46	51.20	75.88	56.62	<u>73.22</u>	41.85	60.058
BRIEF (Ours)	78.59	51.30	59.11	69.33	55.51	83.02	57.10	74.01	42.46	63.381
<i>Mistral-Nemo-12B</i>										
DSIR	52.69	32.30	33.58	56.44	20.09	58.82	26.47	45.45	15.95	37.977
BM25	52.12	33.65	<u>35.19</u>	<u>62.11</u>	20.78	<u>60.81</u>	24.53	45.84	16.37	39.044
LESS	53.57	32.91	<u>33.77</u>	63.86	20.54	59.14	26.95	46.03	16.19	39.218
IFD	54.28	<u>34.49</u>	33.89	57.34	21.62	59.54	26.19	46.62	16.81	38.976
DELIFT	54.67	<u>34.10</u>	33.94	57.45	21.90	60.18	27.01	46.94	17.10	39.254
TAGCOS	54.42	33.40	33.90	56.90	21.29	59.22	26.67	<u>47.20</u>	<u>17.46</u>	38.940
GRAD-MATCH	55.08	34.18	34.63	57.65	21.49	59.91	<u>27.38</u>	<u>46.54</u>	<u>16.72</u>	<u>39.287</u>
CRAIG	<u>55.47</u>	33.90	34.10	57.06	<u>21.98</u>	59.46	26.90	46.29	16.54	39.073
Random	53.99	32.79	33.79	57.00	21.90	59.63	26.28	45.92	16.31	38.623
BRIEF (Ours)	57.29	36.14	35.37	58.46	23.78	62.98	28.09	48.58	18.36	41.006

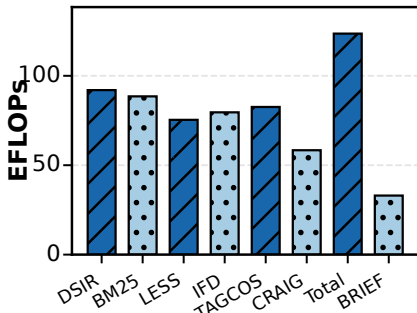


Figure 4: EFLOPs Analysis.

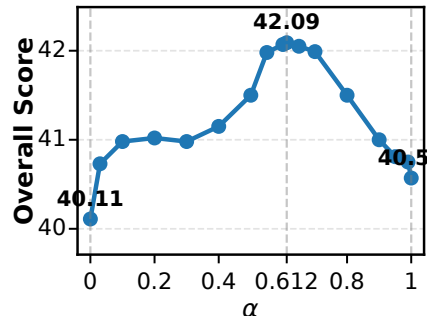


Figure 5: Impact of α .

Dataset	0%	5%(default)	25%	100%
<i>Llama-3.1-8B</i>				
Tulu-3-SFT-Mixture	40.15	42.09	42.31	42.55
Magpie-Pro	17.55	19.11	19.23	19.37
Code-74k-ShareGPT	23.38	26.51	26.87	26.96
Average Score	27.03	29.24	29.47	29.63
<i>Qwen3-4B</i>				
Tulu-3-SFT-Mixture	60.11	63.38	63.51	63.73
Magpie-Pro	28.06	31.56	31.61	31.72
Code-74k-ShareGPT	41.15	43.73	43.90	44.01
Average Score	43.11	46.22	46.34	46.49

Table 4: Warm-up ratio ablation. Overall score (higher is better).

BRIEF algorithm, which achieves the highest performance across all tested α values, thus validating the effectiveness of our method. **Effectiveness of Loss Decomposition.** To verify our optimization objective, we present the GA errors in both the KN and IF spaces under different coreset optimization objectives, along with their corresponding accuracy, as shown in Figure 6. Specifically, the x-axis represents the GA error in the IF gradient space, and the y-axis

represents the GA error in the KN gradient space; points closer to the origin indicate smaller GA errors. We compare BRIEF with the following methods: CRAIG uses gradients from the standard SFT loss: it selects a coreset without distinguishing KN and IF; in BRIEF we extend the α from $(0, 1)$ to $[0, 1]$, where $\alpha = 1$ only uses the KN gradient and $\alpha = 0$ only uses the IF gradient, in which case our method reduces to a Craig-style single-gradient variant. This

Table 3: The time complexity and the actual time of each coreset method when achieving comparable accuracy. Here, $|D|$ denotes the candidate data pool size, K the subset size, δ the convergence threshold, and c the number of clusters.

Method	Data Selection		Total	Accuracy
	Complexity	Actual (h)	Actual (h)	
TAGCOS	$O(D \cdot K + D \cdot c)$	0.6	28.9	42.083
CRAIG	$O(D ^2 \cdot K)$	0.7	23.8	42.071
BRIEF (Ours)	$O(\log(1/\delta) D ^2 K)$	2.4	13.6	42.090

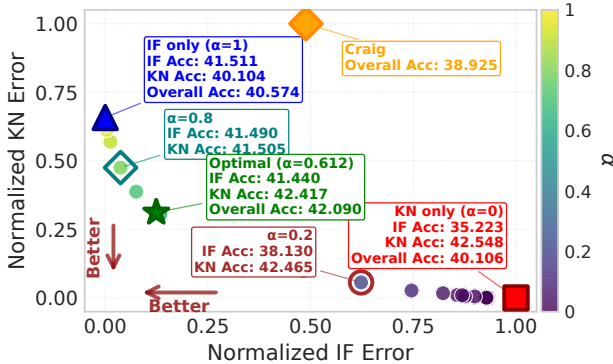


Figure 6: Error of KN and IF under varying α

result highlights a critical limitation in traditional coreset selection approaches: when optimizing with the full SFT gradient, the KN and IF components always interfere destructively with each other. **Analysis of α on instruction following and knowledge.** We categorize benchmarks into instruction following (e.g., DROP, AGIEval, and KorBench) and knowledge-related tasks (e.g., MMLU, MMLU-PRO, GSM8K, MATH, HumanEval, and LCB) to analyze the impact of α . Figure 6 illustrates a clear trade-off: increasing α improves instruction following but compromises knowledge accuracy, while decreasing α does the reverse. Crucially, BRIEF automatically selects an α located in the balanced region, ensuring high performance on both task types.

6.6 RQ5: Impact of Hyperparameters

We now study two hyperparameters in BRIEF: the warm-up data ratio and the training budget (number of fine-tuning epochs).

Effect of warm-up ratio. We first study how the warm-up ratio used to compute gradients affects BRIEF. Using Llama-3.1-8B and Qwen3-4B on different candidate data pools, we compare BRIEF when gradients are computed on the base model (no warm-up) and on models trained with warm-up ratios of 5%, 25%, and 100% of the candidate data. Table 4 shows that using the base model without warm-up greatly hurts BRIEF’s performance, likely due to the shift between the pre-training and fine-tuning input distributions [62]. As the warm-up ratio increases, the model performance improves, but the total training cost (FLOPs) also grows. Notably, a 5% warm-up already achieves performance comparable to using 25% data, at

Table 5: Effect of training budget on overall score and total EFLOPs. Each cell reports Score / EFLOPs.

Method	1 epoch	2 epochs	4 epochs	6 epochs
Full Data (100%)	40.90 / 28.95	41.96 / 60.88	42.88 / 123.60	42.02 / 190.74
Random (5%)	36.50 / 1.55	37.22 / 3.10	38.37 / 6.08	37.72 / 10.01
CRAIG (5%)	36.91 / 15.54	37.89 / 17.22	38.93 / 20.27	38.21 / 23.28
BRIEF (5%)	39.99 / 27.67	41.10 / 29.03	42.09 / 32.18	41.28 / 35.33

a much lower cost. Therefore, we use a 5% warm-up ratio as the default, which is consistent with prior work [62, 69].

Effect of training budget. We also add an explicit ablation on training budget to study how longer training interacts with data selection. The result is reported in Table 5, where we report overall score and total EFLOPs as we increase the number of fine-tuning epochs. For this experiment, we fix Llama-3.1-8B and Tulu-3-SFT-Mixture as the candidate data pool, keep a 5% subset size for all methods, and vary the number of fine-tuning epochs from 1 to 2, 4, and 6. For each epoch setting, we fine-tune the model with BRIEF, CRAIG, random selection, and a full-dataset baseline.

We make three observations. First, training longer does not always help: all methods peak at the 4th epoch and then drop at the 6th epoch due to overfitting, a phenomenon also observed in prior work on LLM [55]. Second, full-data training obtains the best accuracy at each epoch, while BRIEF closely matches this performance and consistently outperforms all other subset baselines. Third, although BRIEF has an initial cost for data selection, it becomes efficient quickly: by the 4th epoch, this cost is offset, and BRIEF saves about 70% EFLOPs compared to full-data training.

7 CONCLUSION

This work takes a simple but informative algebraic view of the SFT loss, showing that it naturally decomposes into knowledge-related and instruction-following components and clarifying why conventional coreset methods can fail to preserve both capabilities in instruction tuning. Building on this perspective, we propose an objective that explicitly minimizes the aggregate GA error over these two components and use a bi-level optimization framework. Experiments on state-of-the-art LLMs and real-world datasets demonstrate that our approach, BRIEF, reduces fine-tuning costs by up to 3 \times while improving accuracy by 5%.

ACKNOWLEDGMENTS

Chengliang Chai is supported by the NSFC (U25B2019, 62472031), the National Key Research and Development Program of China (2024YFC3308200), Beijing Nova Program, and Huawei. Yuping Wang is supported by the NSFC (U23A20297, U23A20317). Ye Yuan is supported by the NSFC (Grant Nos. 62225203, 62532001), Fundamental and Interdisciplinary Disciplines Breakthrough Plan of the Ministry of Education of China (Grant No. JYB2025XDXM108), Beijing Natural Science Foundation (Grant No. L241010). Guoren Wang is supported by the NSFC (62427808, U2001211), and the Liaoning Revitalization Talents Program (XLYC2204005). Lei Cao is supported by the NSF (DBI-2327954) and Amazon Research Awards.

REFERENCES

- [1] Amro Abbas, Kushal Tirumala, Dániel Simig, Surya Ganguli, and Ari S Morcos. 2023. Semdedup: Data-efficient learning at web-scale through semantic deduplication. *arXiv preprint arXiv:2303.09540* (2023).
- [2] Ishika Agarwal, Krishnateja Killamsetty, Lucian Popa, and Marina Danilevsky. [n.d.]. DELIFT: Data Efficient Language model Instruction Fine-Tuning. In *The Thirteenth International Conference on Learning Representations*.
- [3] Ajibawa. 2025. Code-74k-ShareGPT. <https://huggingface.co/datasets/ajibawa-2023/Code-74k-ShareGPT>. <https://huggingface.co/datasets/ajibawa-2023/Code-74k-ShareGPT> HuggingFace dataset.
- [4] Chengliang Chai, Jiabin Liu, Nan Tang, Ju Fan, Dongjing Miao, Jiayi Wang, Yuyu Luo, and Guoliang Li. 2023. Goodcore: Data-effective and data-efficient machine learning through coresets selection over incomplete data. *Proceedings of the ACM on Management of Data* 1, 2 (2023), 1–27.
- [5] Chengliang Chai, Jiayi Wang, Nan Tang, Ye Yuan, Jiabin Liu, Yuhao Deng, and Guoren Wang. 2023. Efficient coresets selection with cluster-based methods. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 167–178.
- [6] Jie Chen, Zhipeng Chen, Jiapeng Wang, Kun Zhou, Yutao Zhu, Jinhao Jiang, Yingqian Min, Wayne Xin Zhao, Zhicheng Dou, Jiabin Mao, et al. 2024. Towards Effective and Efficient Continual Pre-training of Large Language Models. *arXiv preprint arXiv:2407.18743* (2024).
- [7] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374* (2021).
- [8] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168* (2021).
- [9] OpenCompass Contributors. 2023. OpenCompass: A Universal Evaluation Platform for Foundation Models. <https://github.com/open-compass/opencompass>.
- [10] Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning Over Paragraphs. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 2368–2378.
- [11] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [12] Vidhisha Balachandran Pradeep Dasigi Valentina Pyatkin Abhilasha Ravichander Sarah Wiegrefe Nouha Dziri Khyathi Chandu Jack Hessel Yulia Tsvetkov Noah A. Smith Yejin Choi Hannaneh Hajishirzi Faeze Brahman, Sachin Kumar. 2024. The Art of Saying No: Contextual Noncompliance in Language Models. (2024).
- [13] Run-Ze Fan, Zengzhi Wang, and Pengfei Liu. 2025. MegaScience: Pushing the Frontiers of Post-Training Datasets for Science Reasoning. *arXiv preprint arXiv:2507.16812* (2025). <https://arxiv.org/abs/2507.16812>
- [14] Dan Feldman. 2020. Core-sets: Updated survey. *Sampling techniques for supervised or unsupervised tasks* (2020), 23–44.
- [15] Satoru Fujishige. 2005. *Submodular functions and optimization*. Vol. 58. Elsevier.
- [16] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027* (2020).
- [17] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. 2023. Textbooks are all you need. *arXiv preprint arXiv:2306.11644* (2023).
- [18] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirog Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948* (2025).
- [19] Muhammad Usman Hadi, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, Seyedali Mirjalili, et al. 2023. A survey on large language models: Applications, challenges, limitations, and practical usage. *Authorea Preprints* (2023).
- [20] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. *Proceedings of the International Conference on Learning Representations (ICLR)* (2021).
- [21] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. [n.d.]. Measuring Mathematical Problem Solving With the MATH Dataset. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- [22] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR* 1, 2 (2022), 3.
- [23] Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. 2021. Grad-match: Gradient matching based data subset selection for efficient deep model training. In *International Conference on Machine Learning*. PMLR, 5464–5474.
- [24] Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. 2021. Glist: Generalization based data subset selection for efficient and robust learning. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 8110–8118.
- [25] Krishnateja Killamsetty, Xujiang Zhao, Feng Chen, and Rishabh Iyer. 2021. Retrieve: Coresets selection for efficient and robust semi-supervised learning. *Advances in neural information processing systems* 34 (2021), 14488–14501.
- [26] Andreas Krause and Daniel Golovin. 2014. Submodular function maximization. *Tractability* 3, 71–104 (2014), 3.
- [27] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. 2024. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124* (2024).
- [28] Changho Lee, Janghoon Han, Seonghyeon Ye, Stanley Jungkyu Choi, Honglak Lee, and Kyunghoon Bae. 2024. Instruction Matters: A Simple yet Effective Task Selection for Optimized Instruction Tuning of Specific Tasks. *arXiv preprint arXiv:2404.16418* (2024).
- [29] Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. 2024. Selective Reflection-Tuning: Student-Selected Data Recycling for LLM Instruction-Tuning. In *Findings of the Association for Computational Linguistics ACL 2024*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand and virtual meeting, 16189–16211. <https://aclanthology.org/2024.findings-acl.958>
- [30] Ming Li, Yong Zhang, Shwai He, Zhitao Li, Hongyu Zhao, Jianzong Wang, Ning Cheng, and Tianyi Zhou. 2024. Superfiltering: Weak-to-strong data filtering for fast instruction-tuning. *arXiv preprint arXiv:2402.00530* (2024).
- [31] Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. 2024. From Quantity to Quality: Boosting LLM Performance with Self-Guided Data Selection for Instruction Tuning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, Kevin Duh, Helena Gomez, and Steven Bethard (Eds.). Association for Computational Linguistics, Mexico City, Mexico, 7595–7628. <https://aclanthology.org/2024.naacl-long.421>
- [32] Yunshui Li, Binnyun Hui, Xiaobo Xia, Jiayi Wang, Min Yang, Lei Zhang, Shuzheng Si, Ling-Hao Chen, Junhao Liu, Tongliang Liu, et al. 2023. One-shot learning as instruction data prospector for large language models. *arXiv preprint arXiv:2312.10302* (2023).
- [33] Yiming Li, Yanyan Shen, and Lei Chen. 2022. Camel: Managing data for efficient stream learning. In *Proceedings of the 2022 International Conference on Management of Data*. 1271–1285.
- [34] Zhenghao Lin, Zhibin Gou, Yeyun Gong, Xiao Liu, Yelong Shen, Ruochen Xu, Chen Lin, Yujiu Yang, Jian Jiao, Nan Duan, et al. 2024. Rho-1: Not all tokens are what you need. *arXiv preprint arXiv:2404.07965* (2024).
- [35] 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. *arXiv preprint arXiv:2312.15685* (2023).
- [36] Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The Flan Collection: Designing Data and Methods for Effective Instruction Tuning. *arXiv preprint arXiv:2301.13688* (2023).
- [37] Kaijing Ma, Xeron Du, Yunran Wang, Haoran Zhang, ZhoufutuWen, Xingwei Qu, Jian Yang, Jiaheng Liu, minghao liu, Xiang Yue, Wenhao Huang, and Ge Zhang. 2025. KOR-Bench: Benchmarking Language Models on Knowledge-Orthogonal Reasoning Tasks. In *The Thirteenth International Conference on Learning Representations*. <https://openreview.net/forum?id=SVRRQ8goQo>
- [38] Max Marion, Ahmet Üstün, Luiza Pozzobon, Alex Wang, Marzieh Fadaee, and Sara Hooker. 2023. When less is more: Investigating data pruning for pretraining llms at scale. *arXiv preprint arXiv:2309.04564* (2023).
- [39] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196* (2024).
- [40] Baharan Mirzasoileman, Jeff Bilmes, and Jure Leskovec. 2020. Coresets for data-efficient training of machine learning models. In *International Conference on Machine Learning*. PMLR, 6950–6960.
- [41] Baharan Mirzasoileman, Kaidi Cao, and Jure Leskovec. 2020. Coresets for robust training of deep neural networks against noisy labels. *Advances in Neural Information Processing Systems* 33 (2020), 11465–11477.
- [42] Mistral AI and NVIDIA. 2024. Mistral-Nemo-Base-2407. <https://huggingface.co/mistralai/Mistral-Nemo-Base-2407>.
- [43] Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. 2024. Scaling data-constrained language models. *Advances in Neural Information Processing Systems* 36 (2024).

- [44] Alexander Munteanu and Chris Schwiegelshohn. 2018. Coresets-methods and history: A theoreticians design pattern for approximation and streaming algorithms. *KI-Künstliche Intelligenz* 32 (2018), 37–53.
- [45] Mxode. 2025. Magpie-Pro-10K-GPT4o-mini. <https://huggingface.co/datasets/Mxode/Magpie-Pro-10K-GPT4o-mini>. <https://huggingface.co/datasets/Mxode/Magpie-Pro-10K-GPT4o-mini>.
- [46] 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. *arXiv preprint arXiv:2307.06435* (2023).
- [47] George L Nemhauser, Laurence A Wolsey, and Marshall L Fisher. 1978. An analysis of approximations for maximizing submodular set functions—I. *Mathematical programming* 14 (1978), 265–294.
- [48] OpenCompass. [n.d.]. *CodeBench: LCBench - LeetCode Contest Benchmark Dataset*. <https://github.com/open-compass/CodeBench/>
- [49] 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. *Advances in neural information processing systems* 35 (2022), 27730–27744.
- [50] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116* (2023).
- [51] Omead Pooladzandi, David Davini, and Baharan Mirzasoileiman. 2022. Adaptive second order coreset for data-efficient machine learning. In *International Conference on Machine Learning*. PMLR, 17848–17869.
- [52] Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends® in Information Retrieval* 3, 4 (2009), 333–389.
- [53] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024).
- [54] Chaoyuan Shen, Chi Zhang, Chengliang Chai, Jiacheng Wang, Jia Yuan, Yuping Wang, Ye Yuan, Guoren Wang, and Lei Cao. 2025. *BRIEF: Bi-level Coreset Selection for Efficient Instruction Tuning in LLMs*. Technical Report. https://github.com/HR10108/BRIEF/blob/main/BRIEF_Technical_Report.pdf Technical Report.
- [55] Zhengyan Shi, Adam X Yang, Bin Wu, Laurence Aitchison, Emine Yilmaz, and Aldo Lipani. 2024. Instruction tuning with loss over instructions. *Advances in Neural Information Processing Systems* 37 (2024), 69176–69205.
- [56] Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. 2024. Dolma: An open corpus of three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159* (2024).
- [57] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhrranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. 2024. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems* 37 (2024), 95266–95290.
- [58] Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2019. CCNet: Extracting high quality monolingual datasets from web crawl data. *arXiv preprint arXiv:1911.00359* (2019).
- [59] Alexander Wettig, Aatmik Gupta, Saumya Malik, and Danqi Chen. 2024. QuRating: Selecting High-Quality Data for Training Language Models. In *Forty-first International Conference on Machine Learning*. <https://openreview.net/forum?id=GLGYqPwjy>
- [60] Xuansheng Wu, Wenlin Yao, Jianshu Chen, Xiaoman Pan, Xiaoyang Wang, Ninghao Liu, and Dong Yu. 2024. From Language Modeling to Instruction Following: Understanding the Behavior Shift in LLMs after Instruction Tuning. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 2341–2369.
- [61] Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. Less: Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333* (2024).
- [62] Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. 2024. LESS: Selecting Influential Data for Targeted Instruction Tuning. In *Forty-first International Conference on Machine Learning*.
- [63] Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. 2023. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems* 36 (2023), 34201–34227.
- [64] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. *arXiv preprint arXiv:2505.09388* (2025).
- [65] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115* (2024).
- [66] Zichun Yu, Spandan Das, and Chenyan Xiong. 2024. MATES: Model-Aware Data Selection for Efficient Pretraining with Data Influence Models. *arXiv preprint arXiv:2406.06046* (2024).
- [67] Chi Zhang, Huaping Zhong, Hongtao Li, Chengliang Chai, Jiawei Hong, Yuhao Deng, Jiacheng Wang, Tian Tan, Yizhou Yan, Jiantao Qiu, et al. 2025. Not All Documents Are What You Need for Extracting Instruction Tuning Data. *arXiv preprint arXiv:2505.12250* (2025).
- [68] Chi Zhang, Huaping Zhong, Kuan Zhang, Chengliang Chai, Rui Wang, Xinlin Zhuang, Tianyi Bai, Jiantao Qiu, Lei Cao, Ju Fan, et al. 2024. Harnessing diversity for important data selection in pretraining large language models. *arXiv preprint arXiv:2409.16986* (2024).
- [69] Jipeng Zhang, Yaxuan Qin, Renjie Pi, Weizhong Zhang, Rui Pan, and Tong Zhang. 2024. TAGCOS: Task-agnostic Gradient Clustered Coreset Selection for Instruction Tuning Data. *arXiv preprint arXiv:2407.15235* (2024).
- [70] Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792* (2023).
- [71] Yifan Zhang, Yifan Luo, Yang Yuan, and Andrew C Yao. 2024. Autonomous data selection with language models for mathematical texts. In *ICLR 2024 Workshop on Navigating and Addressing Data Problems for Foundation Models*.
- [72] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).
- [73] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364* (2023).