

# HOLDOUT-LOSS-BASED DATA SELECTION FOR LLM FINETUNING VIA IN-CONTEXT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Fine-tuning large pretrained language models is a common approach for aligning them with human preferences, but noisy or off-target examples can dilute supervision. While small, well-chosen datasets often match the performance of much larger ones, systematic and efficient ways to identify high-value training data remain underexplored. Many current methods rely on heuristics or expensive retraining. We present a theoretically grounded, resource-efficient framework for data selection and reweighting. At its core is an In-Context Approximation (ICA) that estimates the holdout loss a model would incur after training on a candidate example by conditioning on a small, curated holdout set in context. ICA requires no reference model and no additional finetuning. Under a local linearization, ICA is equivalent to a first-order update toward the holdout optimum, motivating its use as a proxy for data value. We derive per-example weights from ICA scores, dynamically reweighting gradient updates as model parameters evolve. Across SFT, DPO, and SimPO, and over diverse backbones and datasets, ICA-based reweighting consistently improves model alignment with minimal overhead. We analyze sensitivity to score update frequency and the choice of  $k$  holdout examples for in-context demonstrations, and note limitations for rapidly drifting on-policy updates, highlighting directions for future work. Code and prompts will be released.

## 1 INTRODUCTION

Fine-tuning has become the standard approach for adapting large pretrained language models to downstream applications and aligning them with human intent (Wei et al., 2021; Ouyang et al., 2022). This process typically starts with supervised fine-tuning (SFT) on instruction-response pairs (Wei et al., 2021), followed by preference-based alignment using pairwise preference data, e.g., RLHF (Christiano et al., 2017; Ouyang et al., 2022), DPO (Rafailov et al., 2023), and simPO (Meng et al., 2024). Since fine-tuning effectively steers the pretrained model toward desired behaviors, the quality of training data plays a central role: high-quality examples provide clear alignment signals, while noisy or inconsistent data can severely degrade performance.

However, training data for fine-tuning, typically collected from human annotators or model-generated outputs, often contain errors, inconsistencies, or redundancies. For example, Gao et al. (2024) report that 20–40% of preference pairs in LM alignment are noisy, and that alignment performance is highly sensitive to such noise. In addition, recent studies have shown that smaller but carefully curated datasets can yield alignment performance comparable to much larger ones (Zhou et al., 2023; Chen et al., 2023). Despite the recognized importance and promise of effective data selection for fine-tuning, systematic yet principled approaches remain underexplored.

A central challenge is the absence of consensus on what constitutes “valuable” data. While some approaches adopt hand-crafted heuristics or leverage large models as judges, a more principled perspective is to define data value by its impact on downstream performance, which is the ultimate goal of fine-tuning. Nevertheless, directly measuring each example’s contribution to downstream performance would require retraining (and evaluating) the model across all possible candidates, a process that is computationally infeasible.

To overcome this challenge, prior work has explored different strategies, including influence-function-based methods that approximate model performance change using Taylor expansions (Pruthi et al., 2020; Xia et al., 2024; Wang et al., 2024), surrogate models that fit a linear

approximation of the loss (Engstrom et al., 2024), and bilevel optimization or meta-learning approaches (Shen et al., 2024; Grangier et al., 2023; Calian et al., 2025). Beyond these analytical and learning-based approaches, empirical studies examine correlations between hand-picked notions of data quality and downstream performance to inform heuristic data selection criteria (Liu et al., 2023; Bukharin & Zhao, 2023; Zhao et al., 2024; Lu et al., 2023; Li et al., 2023a; Cao et al., 2023; Li et al., 2023b). These methods, while promising, are often computationally expensive or lack theoretical grounding.

In this work, we propose a computationally efficient and theoretically grounded framework for data selection in model fine-tuning. Following prior work (Mindermann et al., 2022; Xia et al., 2024; Calian et al., 2025), we aim to select a subset of training data such that the model trained on it minimizes the holdout loss, i.e., the loss on this holdout set. Each example’s value is quantified by the holdout loss the model would incur if trained with that example. Computing this naively is intractable, so we build on a tractable approximation introduced by RHO-Loss (Mindermann et al., 2022), which estimates the holdout loss from a Bayesian perspective. Nonetheless, efficiently computing this approximation remains challenging, as it would normally require retraining on updated training and holdout sets at each step.

RHO-Loss mitigates per-step retraining but requires a reference model trained on the holdout set. We introduce an in-context approximation (ICA) that eliminates the need for both additional fine-tuning and a fixed reference model. Building on the insight that in-context learning performs implicit fine-tuning (Dai et al., 2023), we provide the holdout set as in-context demonstrations at each training step, simulating one step of fine-tuning. The resulting data values derived from ICA, termed ICA scores, enable dynamic evaluation of each example’s utility as the model evolves. These scores are then used to reweight gradient updates during fine-tuning, prioritizing examples that most reduce holdout loss. Experiments show that training with ICA-based reweighting consistently improves model alignment for SFT, DPO, and SimPO across diverse datasets and backbone models.

## 2 RELATED WORK

A central question in data selection is how to determine what makes a training example valuable. This is often done by quantifying each example’s impact on a downstream proxy, typically measured as the loss on a small, high-quality holdout set. Methods in this category include influence-function formulations (Pruthi et al., 2020; Xia et al., 2024; Wang et al., 2024), Data Shapley (Ghorbani & Zou, 2019; Wang et al., 2025), and learned scorers such as Datamodels (Engstrom et al., 2024), meta-learning frameworks (Calian et al., 2025), and optimization-based approaches (Grangier et al., 2023; Shen et al., 2024; Gu et al., 2025; Pan et al., 2025). Particularly relevant to our work is RHO-Loss (Mindermann et al., 2022), which estimates the holdout loss a model would incur if trained on a particular example, but requires a separate reference model. Another related approach is One-Shot Learning (Li et al., 2023b), which also leverages in-context learning, though in a different manner from our method.

In addition to analytic or learned approaches for quantifying data value, hand-crafted notions of data quality have been studied, with empirical analyses assessing how these proxies correlate with fine-tuned model performance and inform data selection (Liu et al., 2023; Bukharin & Zhao, 2023; Zhao et al., 2024; Lu et al., 2023; Li et al., 2023a; Cao et al., 2023; Huang & Goyal, 2025; Yu et al., 2025; Deng et al., 2025; Morimura et al., 2024). Beyond quantifying data value with respect to downstream performance, prior work has also explored selecting data to directly match the target distribution, using techniques such as gradient alignment (Killamsetty et al., 2021), importance resampling (Xie et al., 2023; Katharopoulos & Fleuret, 2018), or optimal transport (Kang et al., 2024). In contrast to these data-centric methods, a separate line adopts an algorithmic perspective, directly modifying learning objectives and proposing extensions to SFT or DPO to account for instance-specific differences or improve generalization (Wu et al., 2024; D’Oosterlinck et al., 2025; Wu et al., 2025).

## 3 PRELIMINARY

**In-context learning** LLMs have demonstrated strong in-context learning (ICL) capabilities (Brown et al., 2020; Cao et al., 2023). In ICL, a few demonstration examples are concatenated

with a query to form the model input. The model then identifies patterns from these examples and makes predictions without any parameter updates.

Formally, given a query  $\mathbf{x}$ , the model predicts an output  $\mathbf{y}$  conditioned on a demonstration set  $C$ . In this work, we adopt a demonstration set of the following form, which includes an optional task instruction  $I$  followed by  $k$  demonstration examples:

$$C = \{I, s(\mathbf{x}_1, \mathbf{y}_1), \dots, s(\mathbf{x}_k, \mathbf{y}_k)\}.$$

**Supervised fine-tuning** SFT adapts pretrained LLMs to produce responses with desired characteristics using instruction–response pairs. Let  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i^*)\}_{i=1}^{|\mathcal{D}|}$  denote the instruction dataset, where  $\mathbf{y}_i^*$  is the reference response for query  $\mathbf{x}_i$ . SFT updates the model parameters by minimizing the sentence-level cross-entropy:

$$\mathcal{L}_{\text{SFT}}(\theta) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}^*) \sim \mathcal{D}} [-\log \pi_{\theta}(\mathbf{y}^* | \mathbf{x})]$$

where  $\pi_{\theta}(\mathbf{y} | \mathbf{x})$  is the probability of generating  $\mathbf{y}$  given  $\mathbf{x}$  under the model parameterized by  $\theta$ .

**Preference-based alignment methods** RLHF (Christiano et al., 2017; Ouyang et al., 2022) is a widely used method to align LLMs with human preferences. It uses data of the form  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_{w,i}, \mathbf{y}_{l,i})\}_{i=1}^{|\mathcal{D}|}$ , where  $\mathbf{y}_{w,i}$  and  $\mathbf{y}_{l,i}$  are the preferred and dispreferred responses for a prompt  $\mathbf{x}_i$ . The standard RLHF pipeline first learns a reward model and then optimizes a policy using RL algorithms such as PPO.

DPO (Rafailov et al., 2023) provides an off-policy alternative that bypasses the RL step, learning directly from preference data without training a reward model. Specifically, DPO solves:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \frac{\pi_{\theta}(\mathbf{y}_w | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_w | \mathbf{x})} - \beta \frac{\pi_{\theta}(\mathbf{y}_l | \mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_l | \mathbf{x})} \right) \right]$$

where  $\pi_{\text{ref}}$  is a reference policy,  $\sigma$  is the sigmoid function, and  $\beta$  controls the trade-off between adhering to the reference and incorporating new preference data.

A recent variant, SimPO (Meng et al., 2024), extends DPO by eliminating the need for a reference model and introducing a length-normalized implicit reward based on the average log probability of a sequence. Its objective is:

$$\mathcal{L}_{\text{SimPO}}(\theta) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim \mathcal{D}} \left[ \log \sigma \left( \frac{\beta}{|\mathbf{y}_w|} \pi_{\theta}(\mathbf{y}_w | \mathbf{x}) - \frac{\beta}{|\mathbf{y}_l|} \pi_{\theta}(\mathbf{y}_l | \mathbf{x}) - \gamma \right) \right]$$

where  $\gamma > 0$  is a target reward margin ensuring that the reward difference between winning and losing responses exceeds this threshold.

## 4 HOLDOUT-LOSS-BASED DATA SELECTION VIA IN-CONTEXT LEARNING

In this section, we formally introduce the problem formulation and present a holdout-loss-based data selection framework that leverages in-context learning for efficient computation.

### 4.1 PROBLEM FORMULATION

We consider the problem of fine-tuning a pretrained model on a large, but potentially noisy, training dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{|\mathcal{D}|}$ , where  $(\mathbf{x}, \mathbf{y})$  may represent an instruction–response pair or a preference pair. The goal is to optimize model performance on a smaller, higher-quality holdout set  $\mathcal{D}_{\text{ho}} = \{(\mathbf{x}_i^{\text{ho}}, \mathbf{y}_i^{\text{ho}})\}_{i=1}^{|\mathcal{D}_{\text{ho}}|}$ . For example,  $\mathcal{D}_{\text{ho}}$  could be a corrected subset of  $\mathcal{D}$  generated by a stronger model or a manually curated set from a target domain.

However, training on all of  $\mathcal{D}$  may be inefficient and suboptimal due to the presence of noise. While one could train directly on  $\mathcal{D}_{\text{ho}}$ , its small size often leads to overfitting and fails to leverage the information in the larger set. Therefore, we aim to select a subset  $\bar{\mathcal{D}}^* \subset \mathcal{D}$  such that a model trained on  $\bar{\mathcal{D}}^*$  minimizes the loss on  $\mathcal{D}_{\text{ho}}$  (i.e., the holdout loss). Formally, the problem can be framed as

$$\bar{\mathcal{D}}^* = \arg \min_{\bar{\mathcal{D}} \subset \mathcal{D}} \mathcal{L}(\mathcal{D}_{\text{ho}}; \theta^*(\bar{\mathcal{D}})), \quad \text{where} \quad \theta^*(\bar{\mathcal{D}}) = \arg \min_{\theta} \mathcal{L}(\bar{\mathcal{D}}; \theta). \quad (1)$$

We denote by  $\theta^*(\mathcal{S})$  the model parameters obtained by training on any dataset  $\mathcal{S}$  and  $\mathcal{L}(\mathcal{S}; \theta)$  the loss on  $\mathcal{S}$  under model  $\theta$ .

Equation 1 would be prohibitively expensive to solve naively, as it requires training on every candidate subset  $\bar{\mathcal{D}} \subset \mathcal{D}$  and evaluating the holdout loss. Instead of solving it directly, we recast the problem in terms of quantifying the contribution of each training example to reducing holdout loss. Concretely, we assign a score to each example that reflects this contribution. These scores can then be used either to select the example that most reduces the holdout loss, or to reweight gradient updates in the optimizer according to each example’s contribution.

In the following section, we detail how to approximate these scores without performing actual training and how to compute them efficiently using in-context learning.

## 4.2 COMPUTING APPROXIMATED HOLDOUT LOSS VIA IN-CONTEXT LEARNING

**Approximating holdout loss via Bayesian framework** Consider sequential (greedy) data selection, where points are added one at a time. Let  $\mathcal{D}_t$  be the dataset selected up to step  $t$ , and  $\theta_t := \theta^*(\mathcal{D}_t)$  the model trained on it. Then, the subset selection problem in Equation 1 reduces to selecting, at each step, the example  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$  that, when added to  $\mathcal{D}_t$ , minimizes the holdout loss:

$$(\mathbf{x}^*, \mathbf{y}^*) = \arg \min_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \mathcal{L}(\mathcal{D}_{\text{ho}}; \theta^*(\mathcal{D}_t \cup \{(\mathbf{x}, \mathbf{y})\})). \quad (2)$$

Accordingly, each training example’s contribution can be measured by the holdout loss in Equation 2. To compute this holdout loss without performing actual training, we adopt an approximation derived by Mindermann et al. (2022) in the general supervised learning setting, based on probabilistic modeling. Specifically, considering the negative log-likelihood as the loss function ( $\ell(\mathbf{y} | \mathbf{x}; \theta) = -\log p(\mathbf{y} | \mathbf{x}; \theta)$ ), and, applying Bayes’ rule under the conditional independence assumption, the holdout loss of the model trained with a particular example can be approximated as (see Appendix A for a reproduction of the derivation in Mindermann et al. (2022))<sup>1</sup>:

$$\mathcal{L}(\mathcal{D}_{\text{ho}}; \theta^*(\mathcal{D}_t \cup \{(\mathbf{x}, \mathbf{y})\})) \approx \ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}})) - \ell(\mathbf{y} | \mathbf{x}; \theta_t) - \mathcal{L}(\mathcal{D}_{\text{ho}}; \theta_t) \quad (3)$$

where we use  $\mathcal{L}$  for losses over a set, and  $\ell$  for the per-example loss. Omitting the term independent of  $(\mathbf{x}, \mathbf{y})$  and reversing the sign, we define the remaining expression as the holdout-loss score of each example at step  $t$ :

$$s_{\text{ho}}(\mathbf{x}, \mathbf{y}; \theta_t) = \ell(\mathbf{y} | \mathbf{x}; \theta_t) - \ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}})). \quad (4)$$

The optimal example can then be approximately found by maximizing this score:

$$(\mathbf{x}^*, \mathbf{y}^*) = \arg \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} s_{\text{ho}}(\mathbf{x}, \mathbf{y}; \theta_t). \quad (5)$$

The holdout-loss score provides a tractable tool for data selection. However, since  $\mathcal{D}_t$  is updated with each newly added example, the model trained on  $\mathcal{D}_t \cup \mathcal{D}_{\text{ho}}$  must also be updated at every selection step, which still incurs substantial computational overhead. To mitigate this cost, prior work (Mindermann et al., 2022) approximates the retraining by training a model only on  $\mathcal{D}_{\text{ho}}$  and reusing it across all selection steps, i.e.,  $\ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}})) \approx \ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_{\text{ho}}))$ . This results in the reducible holdout loss (RHO-Loss) criterion employed for data selection in their work.

However, RHO-Loss uses a fixed reference model rather than re-evaluating each example’s impact on the holdout loss as the model is updated. This simplification can introduce bias in estimating each example’s contribution (Wang et al., 2024). To enable an efficient and adaptive data selection criterion, we introduce an in-context approximation that removes auxiliary training entirely by approximating  $\ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}}))$  via in-context learning. This technique is described in detail below.

<sup>1</sup>This Bayesian framework also extends to pairwise preference data (e.g., for DPO and simPO), where  $\mathbf{y}$  represents a preference pair  $(\mathbf{y}_w, \mathbf{y}_l)$  with  $\mathbf{y}_w \succ \mathbf{y}_l$ , and the loss is defined correspondingly for the chosen preference learning model (e.g., Bradley-Terry as used in DPO).

**Efficient computation via in-context learning** Motivated by the finding that in-context learning can implicitly perform gradient-based model updates on the provided demonstrations (Dai et al., 2023), we avoid retraining on  $\mathcal{D}_t \cup \mathcal{D}_{\text{ho}}$  by using the holdout set as in-context demonstrations. Specifically, we introduce the following in-context approximation (ICA) to compute the second term in equation 4 without retraining:

$$\ell(\mathbf{y} \mid \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}})) \approx \ell(\mathbf{y} \mid \mathbf{x}, \mathcal{D}_{\text{ho}}; \theta_t). \quad (6)$$

By conditioning on the holdout set  $\mathcal{D}_{\text{ho}}$ , this approximation effectively simulates one step of fine-tuning on the set, providing a computationally tractable estimate of the model’s state after training on  $\mathcal{D}_t \cup \mathcal{D}_{\text{ho}}$ .

Applying ICA to equation 4 yields a computationally efficient selection criterion, which we refer to as the in-context approximation score (ICA score):

$$s_{\text{ICA}}(\mathbf{x}, \mathbf{y}; \theta_t) := \ell(\mathbf{y} \mid \mathbf{x}; \theta_t) - \ell(\mathbf{y} \mid \mathbf{x}, \mathcal{D}_{\text{ho}}; \theta_t). \quad (7)$$

We can now use the ICA score in place of the holdout-loss score in equation 5 as the data selection criterion. This approach is not only more computationally efficient, but also enables dynamic re-evaluation of each example’s impact on the holdout loss as the model evolves.

In practice, examples are often selected in batches rather than sequentially. We next describe how to leverage the ICA score for batch selection using a reweighting strategy.

#### 4.3 BATCH SELECTION VIA A REWEIGHTING STRATEGY

**ICA score-based reweighting** For batch selection, we use a reweighting strategy that up-weights high-scoring examples and downweights lower-scoring ones. Unlike hard selection, which only chooses top examples and may reduce batch diversity or ignore interactions among data points (Wang et al., 2024), reweighting leverages the gradient signal from the entire batch, improving training stability. We apply ICA score-based reweighting in our main experiments and ablate this choice in Section 5.3, comparing it to percentile-based filtering.

Concretely, consider gradient-based training (e.g., Adam or SGD) with mini-batch updates. At each iteration  $t$ , a batch  $B_t \subset \mathcal{D}$  of size  $n_B$  is sampled. For each example in the batch, we compute its ICA score as defined in equation 7 and convert these scores into continuous weights in the range  $[0, 1]$  via max-min normalization. We adopt max-min instead of the softmax used by Wang et al. (2024); Calian et al. (2025) because it preserves the relative differences between scores and avoids the exponential amplification that can distort the contribution of low- and high-scoring examples. The weights are given by

$$w(\mathbf{x}_i, \mathbf{y}_i; \theta_t) = \frac{s(\mathbf{x}_i, \mathbf{y}_i; \theta_t) - \min_{j \in B_t} s(\mathbf{x}_j, \mathbf{y}_j; \theta_t)}{\max_{j \in B_t} s(\mathbf{x}_j, \mathbf{y}_j; \theta_t) - \min_{j \in B_t} s(\mathbf{x}_j, \mathbf{y}_j; \theta_t)}. \quad (8)$$

These weights, reflecting the relative utility of each example for minimizing the holdout loss, are used to scale their contributions to the batch gradient:

$$\mathbf{g}_t = \sum_{i=1}^{|B_t|} w(\mathbf{x}_i, \mathbf{y}_i; \theta_t) \nabla_{\theta} \ell(\mathbf{x}_i, \mathbf{y}_i; \theta_t). \quad (9)$$

The resulting weighted gradient  $\mathbf{g}_t$  is used to update the model parameters via a standard optimizer.

Because the ICA score, and hence the weights, evolve throughout training, this reweighting strategy adaptively adjusts each example’s contribution to the gradient updates according to the holdout loss the model would achieve if trained on that example. Algorithm 1 outlines how the ICA score is computed and used to reweight gradient updates during fine-tuning.

**Practical implementation** In practice, we apply the following two techniques to further improve the efficiency of Algorithm 1.

When computing the in-context approximation (line 6 or equation 6), including the entire holdout set as demonstrations is often infeasible due to prompt length constraints. One solution is to divide the full holdout set into multiple subsets, compute the score for each, and average the results; however,

**Algorithm 1** Reweighting training examples using ICA scores for LLM fine-tuning

---

```

1: Input: Training set  $\mathcal{D}$ ; Holdout set  $\mathcal{D}_{\text{ho}}$ ; Pre-trained model parameters  $\theta$ ; Number of training
   steps  $T$ ; Batch size  $n_B$ ; Optimizer OPTIMIZER
2: Initialize  $\theta_0 \leftarrow \theta$ 
3: for  $t = 0, \dots, T - 1$  do
4:   Sample candidate set  $B_t \subset \mathcal{D}$  of size  $n_B$ 
5:   for  $i = 1, \dots, n_B$  do
6:      $\text{ConditionalLoss}_i \leftarrow \ell(\mathbf{y}_i \mid \mathbf{x}_i, \mathcal{D}_{\text{ho}}; \theta_t)$ 
7:      $\text{Loss}_i \leftarrow \ell(\mathbf{y}_i \mid \mathbf{x}_i; \theta_t)$ 
8:      $\text{Score}_i \leftarrow \text{Loss}_i - \text{ConditionalLoss}_i$ 
9:   end for
10:  Compute per-example weights within the batch using equation 8
11:  Compute the reweighted batch gradient  $\mathbf{g}_t$  on  $B_t$  using equation 9
12:   $\theta_{t+1} \leftarrow \text{OPTIMIZER}(\theta_t, \mathbf{g}_t)$ 
13: end for
Return Finetuned model parameters  $\theta_T$ 

```

---

this can be computationally intensive. Instead, we select the top- $k$  holdout examples most similar to each candidate via  $k$ -nearest neighbor (kNN) search in an embedding space. To further improve efficiency, we periodically update the scores  $R$  times during training, instead of recomputing them at every iteration. Experimental results show that even with these practical approximations, our method consistently improves alignment performance.

The complete procedure incorporating these techniques is presented as Algorithm 2 in Appendix B.4. We ablate the effects of different choices of  $k$  and  $R$  in Section 5.3 and analyze the computational overhead of our implementation in Section 5.4.

## 5 EXPERIMENTAL RESULTS

We apply our method to both SFT and preference-based alignment (DPO and SimPO), comparing ICA score-based reweighting to standard training (without reweighting) and to reweighting using scores computed by baseline methods, across multiple models and datasets.

### 5.1 EXPERIMENTAL SETUP

**Evaluation protocol and metrics** Since our goal is to select a subset of training data that minimizes holdout loss, which serves as a proxy for downstream performance, we evaluate our method by measuring how closely the model’s outputs align with the test set targets. To this end, we define the win rate as the percentage of pairwise comparisons in which a response from a model trained with our method is judged closer to the target output than that of a model trained with standard (non-reweighted) or baseline methods. Comparisons are performed by GPT-4o (2024-08-06). Due to computational constraints, we train and evaluate each model only once. The standard deviation of evaluation results is small (see Appendix B.2), so reporting a single evaluation run is sufficient.

**Datasets** We consider two data selection scenarios. **High-quality selection** prioritizes expert-level data using a smaller curated holdout set: for SFT, we use Alpaca as training data and sample high-quality holdout examples from its curated version, Alpaca-cleaned (Taori et al., 2023); for preference optimization (DPO and SimPO), we use UltraFeedback-binarized (Cui et al., 2023), which provides preference pairs with scalar quality scores, enabling the construction of high-quality holdout and test sets. **Domain-relevant selection** prioritizes examples relevant to a target domain using a domain-specific holdout set: for SFT, we use Yahoo\_Answers\_Topic<sup>2</sup>, and for preference optimization, we use SHP-2 (Ethayarajh et al., 2022); both datasets contain domain labels, allowing selection of a

<sup>2</sup>[https://huggingface.co/datasets/community-datasets/yahoo\\_answers\\_topics/viewer?views%5B%5D=train](https://huggingface.co/datasets/community-datasets/yahoo_answers_topics/viewer?views%5B%5D=train)

Table 1: Win rates of our method for SFT across models and datasets (higher values indicate better performance), comparing outputs from our method to each baseline.

Ours against	Alpaca			Yahoo_Answers_Topic		
	w/o	RHO-Loss	One-Shot	w/o	RHO-Loss	One-Shot
LLaMA 3B	77.81	48.96	57.03	78.90	46.73	62.33
LLaMA 8B	80.55	49.92	62.11	85.10	54.03	66.93
Qwen 4B	71.21	50.56	56.82	80.30	49.93	58.73
Qwen 8B	82.92	51.21	58.33	82.93	54.43	63.13

target domain for evaluating out-of-domain alignment. Details of the dataset splits (training, test, and holdout) are provided in Appendix B.1.

**Models and training** We evaluate our approach on multiple model families and scales, including LLaMA-3-8B-Instruct, LLaMA-3-3B-Instruct, Qwen-3-8B, and Qwen-3-4B. Models are fine-tuned using two parameter-updating paradigms: full-parameter fine-tuning and parameter-efficient LoRA (Hu et al., 2021). Due to space constraints, detailed training configurations are provided in Appendix B.3, and LoRA fine-tuning results with our method are reported in Appendix B.8.

**Default setting** By default, we perform full-parameter training and use the ICA score to reweight training examples. We compute ICA scores using the two practical techniques described in Section 4.3 (see Appendix B.4 for additional implementation details). Specifically, for each candidate, we use the top  $k = 3$  holdout examples as in-context demonstrations and update scores for all training examples  $R = 1$  time (computing scores for all training examples only at the initialization step  $t = 0$ ). When selecting the top  $k$  holdout examples in embedding space, we adopt all-mpnet-base-v2 (Reimers & Gurevych, 2019) to compute embeddings.

**Baselines** We compare our method with standard training (without reweighting) and with reweighting using scores from the following methods in place of the ICA score: (1) **RHO-Loss** (Mindermann et al., 2022), which approximates the holdout loss score using a model trained on the holdout set; (2) **One-shot learning** (Li et al., 2023b), which scores each candidate as the difference between the one-shot loss with the candidate included as context and the zero-shot loss without it. Detailed formulas for computing scores with these two baseline methods are provided in Appendix B.5.

## 5.2 MAIN RESULTS

We report the performance of our method across SFT, DPO, and SimPO, comparing against standard training and baseline approaches, in Tables 1, 2, and 3.

**Comparison to standard training** Across all datasets and model families, incorporating our reweighting method leads to consistently better alignment than standard training without reweighting. The improvements hold for both SFT and preference-based alignment, demonstrating that our method provides a robust advantage across different learning paradigms.

**Comparison to baselines** When compared to reweighting using baseline scoring methods, our approach achieves consistent gains over one-shot learning across all settings, with win rates consistently above 50%, and often exceeding 60%. Against RHO-Loss, which approximates the same selection criterion in equation 4 but requires training an auxiliary model, our method achieves comparable performance, in some cases surpassing it, while avoiding the cost of training a reference model. These patterns are consistent across both LLaMA and Qwen families and hold for models of different sizes (3B and 8B).

Table 2: Win rates of our method for DPO across models and datasets (higher values indicate better performance), comparing outputs from our method to each baseline.

Ours against	UltraFeedback-binarized			StanfordNLP/SHP-2		
	w/o	RHO-Loss	One-Shot	w/o	RHO-Loss	One-Shot
LLaMA 3B	61.05	46.85	57.60	79.90	56.70	60.20
LLaMA 8B	64.00	48.25	58.40	77.20	48.70	55.10
Qwen 4B	64.30	51.90	54.60	79.20	49.90	57.00
Qwen 8B	64.85	49.90	60.45	70.40	47.60	54.40

Table 3: Win rates of our method for simPO across models and datasets (higher values indicate better performance, comparing outputs from our method to each baseline.

Ours against	UltraFeedback-binarized			StanfordNLP/SHP-2		
	w/o	RHO-Loss	One-Shot	w/o	RHO-Loss	One-Shot
LLaMA 3B	62.20	47.20	55.45	93.90	50.70	55.90
LLaMA 8B	64.70	49.90	57.70	63.30	46.50	51.30
Qwen 3B	64.95	47.65	51.20	66.00	46.80	49.60
Qwen 8B	65.55	50.60	54.35	77.10	45.50	53.60

### 5.3 ABLATION STUDIES

We perform ablation studies on LLaMA-3B-Instruct trained on Yahoo\_Answers\_Topics to examine the effect of key design choices in our method on model performance, with results detailed in Appendix B.6. Key findings are presented below.

**Large  $k$  is not required for good performance** To reduce computational cost, we use the top- $k$  holdout examples selected via kNN in embedding space as in-context demonstrations, rather than the full holdout set. Using  $k = 3$  as the default, smaller or larger values yield no improvement, with win rates below 50% relative to the default: 43.5% for  $k = 1$ , 48.0% for  $k = 5$ , and 46.0% for  $k = 10$ . This indicates that a small number of holdout examples is sufficient to maintain strong alignment while preserving computational efficiency.

**More frequent score updates improve alignment** We investigate the effect of the total number of score computations  $R$  on model performance. In the default setting, scores are computed once at initialization ( $R = 1$ ). Increasing  $R$  to 3, 5, and 9 improves alignment, with win rates of 50.73%, 52.77%, and 51.6%, respectively. This suggests that additional score updates can further enhance the alignment performance of our method.

**Filtering is less effective** We compare percentile-based filtering, which retains examples above a specified percentile, at different thresholds to the default reweighting method. Simple filtering is generally less effective, with win rates below 50% relative to reweighting. A threshold of the 75th percentile yields a higher win rate (48.67%) than the 50th (40.80%) and 90th (40.07%). These results indicate that retaining some lower-scoring examples can be beneficial, but too many degrade performance, so the threshold must be chosen carefully. Adaptive reweighting eliminates this need for manual selection of filtering threshold by automatically adjusting example importance.

**More advanced embeddings improve performance** To evaluate the effect of different embedding models for selecting the top- $k$  holdout examples, we test BAAI/bge-m3 (Chen et al., 2024), a stronger embedding model. BAAI/bge-m3 yields higher win rates (52%) compared to the default all-mpnet-base-v2, suggesting that more capable embeddings can further improve the alignment achieved by our method.



## 5.4 ANALYSIS

**Computational complexity** Our method introduces two sources of overhead beyond standard fine-tuning: (i) a one-time precomputation to obtain embeddings, and (ii) periodic score updates. The latter is dominated by the in-context approximation term, which involves additional forward passes and accounts for most of the computational cost. We measure runtime on four NVIDIA A6000 GPUs (48 GB each) for LLaMA-3B-Instruct fine-tuned on Yahoo\_Answers\_Topics. Precomputation is negligible, on the order of seconds, and the additional time for score computation and reweighting is reported as a percentage relative to the runtime of standard fine-tuning. Our method adds only  $\sim 1.5\%$  overhead, compared to roughly 10% for RHO-Loss and 4% for one-shot learning. Full results are provided in Appendix B.7.

**Score distribution across training examples** We analyze SFT on Yahoo\_Answers\_Topic to examine the scores assigned to training examples. Using a holdout set sampled from the Sports topic, we compute the average score for examples from each domain (Figure 1). Examples from Sports receive higher scores under both our method and RHO-Loss compared to One-Shot Learning and other topics, demonstrating that our method achieves strong alignment while requiring less computational overhead. Additional analyses in Appendix B.10 examine response patterns produced by models trained with our method, and Appendix B.11 presents instruction-response pairs with high and low ICA scores, offering further insight into our method’s effectiveness.

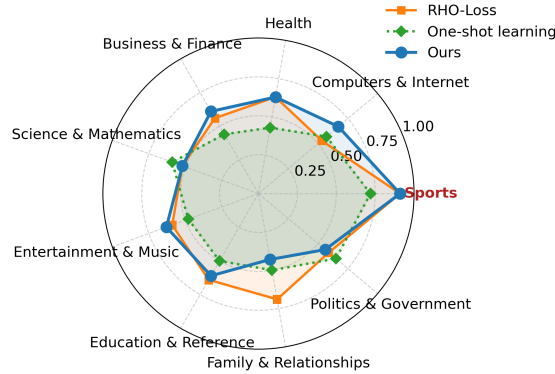


Figure 1: Average scores assigned to training examples from each domain (normalized to  $[0, 1]$ ). The target domain is highlighted in red, and different colors indicate different scoring methods, with our method in blue. Higher scores indicate stronger alignment with the target domain.

## 6 CONCLUSION AND DISCUSSION

We propose a theoretically grounded, resource-efficient framework for data selection in LLM fine-tuning. Our approach leverages an in-context approximation (ICA) to estimate the holdout loss of the model after including each candidate example in training, without requiring additional finetuning or a reference model. The resulting ICA scores are used to dynamically reweight gradient updates during fine-tuning. Empirical results show that ICA score-based reweighting consistently improves model alignment across SFT, DPO, and SimPO over diverse datasets and backbone architectures, with only marginal computational overhead ( $\sim 1.5\%$ ).

We also note limitations and directions for future work. Our method relies on a high-quality holdout set as a proxy for the test distribution; noisy or unrepresentative holdouts may reduce generalization to unseen data, even if alignment appears strong. Additionally, ICA is inherently off-policy, so applying it to on-policy methods such as PPO would require frequent recomputation of scores for all examples as new data is generated and model parameters continuously evolve, creating a computational bottleneck. Addressing these challenges is a promising direction for future work.

## ETHICS STATEMENT

All authors have read, adhere to, and explicitly acknowledge the ICLR Code of Ethics. This work does not involve human subjects or private user data, and all datasets used in our experiments are publicly available and appropriately cited, with full details of data processing and holdout/test splits provided in the appendix to ensure transparency and reproducibility.

All experimental procedures and analyses adhere to accepted research integrity standards, and we disclose any potential conflicts of interest: none exist. This work does not involve applications that directly pose safety or legal risks beyond standard LLM use.

The assistive use of LLMs in preparing this paper, to refine grammar and improve the clarity of the text, is documented in Appendix C, and the authors take full responsibility for all content.

## REPRODUCIBILITY STATEMENT

All code and materials used in our experiments are available at <https://anonymous.4open.science/r/datawork>. To facilitate reproducibility, we provide detailed descriptions of all datasets used, including holdout and test set construction and dataset splits, in Appendix B.1. The algorithm used in our experiments is detailed in Appendix B.4, with full training parameters, LoRA configurations, and evaluation settings provided in Appendix B.3. Prompts used for querying LLMs when computing ICA scores, as well as prompts used for GPT to judge win rates, are provided in Appendix B.9.

For theoretical results, our method builds on the holdout loss approximation introduced in Mindermann et al. (2022), and we reproduce the complete derivation in Appendix A. Together, these materials allow others to reproduce our experiments and verify our results.

## REFERENCES

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Alexander Bukharin and Tuo Zhao. Data diversity matters for robust instruction tuning. *arXiv preprint arXiv:2311.14736*, 2023.
- Dan A Calian, Gregory Farquhar, Iurii Kemaev, Luisa M Zintgraf, Matteo Hessel, Jeremy Shar, Junhyuk Oh, András György, Tom Schaul, Jeffrey Dean, et al. Datarater: Meta-learned dataset curation. *arXiv preprint arXiv:2505.17895*, 2025.
- Yihan Cao, Yanbin Kang, Chi Wang, and Lichao Sun. Instruction mining: Instruction data selection for tuning large language models. *arXiv preprint arXiv:2307.06290*, 2023.
- Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2024. URL <https://arxiv.org/abs/2402.03216>.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpapasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. Ultrafeedback: Boosting language models with high-quality feedback. 2023.
- Damai Dai, Yutao Sun, Li Dong, Yaru Hao, Shuming Ma, Zhifang Sui, and Furu Wei. Why can gpt learn in-context? language models implicitly perform gradient descent as meta-optimizers, 2023. URL <https://arxiv.org/abs/2212.10559>.

- Xun Deng, Han Zhong, Rui Ai, Fuli Feng, Zheng Wang, and Xiangnan He. Less is more: Improving llm alignment via preference data selection. arXiv preprint arXiv:2502.14560, 2025.
- Karel D’Oosterlinck, Winnie Xu, Chris Develder, Thomas Demeester, Amanpreet Singh, Christopher Potts, Douwe Kiela, and Shikib Mehri. Anchored preference optimization and contrastive revisions: Addressing underspecification in alignment. Transactions of the Association for Computational Linguistics, 13:442–460, 2025.
- Logan Engstrom, Axel Feldmann, and Aleksander Madry. Dsdm: Model-aware dataset selection with datamodels, 2024. URL <https://arxiv.org/abs/2401.12926>.
- Kawin Ethayarajh, Yejin Choi, and Swabha Swayamdipta. Understanding dataset difficulty with  $\mathcal{V}$ -usable information. In International Conference on Machine Learning, pp. 5988–6008. PMLR, 2022.
- Yang Gao, Dana Alon, and Donald Metzler. Impact of preference noise on the alignment performance of generative language models. arXiv preprint arXiv:2404.09824, 2024.
- Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. In International conference on machine learning, pp. 2242–2251. PMLR, 2019.
- David Grangier, Pierre Ablin, and Awni Hannun. Bilevel optimization to learn training distributions for language modeling under domain shift. In NeurIPS 2023 Workshop on Distribution Shifts: New Frontiers with Foundation Models, 2023.
- Yuxian Gu, Li Dong, Hongning Wang, Yaru Hao, Qingxiu Dong, Furu Wei, and Minlie Huang. Data selection via optimal control for language models, 2025. URL <https://arxiv.org/abs/2410.07064>.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- Chengyu Huang and Tanya Goyal. Dcrm: A heuristic to measure response pair quality in preference optimization, 2025. URL <https://arxiv.org/abs/2506.14157>.
- Feiyang Kang, Hoang Anh Just, Yifan Sun, Himanshu Jahagirdar, Yuanzhi Zhang, Rongxing Du, Anit Kumar Sahu, and Ruoxi Jia. Get more for less: Principled data selection for warming up fine-tuning in llms. arXiv preprint arXiv:2405.02774, 2024.
- Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning with importance sampling. In International conference on machine learning, pp. 2525–2534. PMLR, 2018.
- Krishnateja Killamsetty, Sivasubramanian Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match: Gradient matching based data subset selection for efficient deep model training. In International Conference on Machine Learning, pp. 5464–5474. PMLR, 2021.
- Ming Li, Yong Zhang, Zhitao Li, Jiuhai Chen, Lichang Chen, Ning Cheng, Jianzong Wang, Tianyi Zhou, and Jing Xiao. From quantity to quality: Boosting llm performance with self-guided data selection for instruction tuning. arXiv preprint arXiv:2308.12032, 2023a.
- Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiayi Yang, Min Yang, Lei Zhang, Shuzheng Si, Ling-Hao Chen, Junhao Liu, Tongliang Liu, et al. One-shot learning as instruction data prospector for large language models. arXiv preprint arXiv:2312.10302, 2023b.
- Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. arXiv preprint arXiv:2312.15685, 2023.
- Keming Lu, Hongyi Yuan, Zheng Yuan, Runji Lin, Junyang Lin, Chuanqi Tan, Chang Zhou, and Jingren Zhou. # instag: Instruction tagging for analyzing supervised fine-tuning of large language models. arXiv preprint arXiv:2308.07074, 2023.

- Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235, 2024.
- Sören Mindermann, Jan M Brauner, Muhammed T Razzak, Mrinank Sharma, Andreas Kirsch, Winnie Xu, Benedikt Höltingen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, et al. Prioritized training on points that are learnable, worth learning, and not yet learnt. In *International Conference on Machine Learning*, pp. 15630–15649. PMLR, 2022.
- Tetsuro Morimura, Mitsuki Sakamoto, Yuu Jinnai, Kenshi Abe, and Kaito Ariu. Filtered direct preference optimization. *arXiv preprint arXiv:2404.13846*, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.
- Rui Pan, Dylan Zhang, Hanning Zhang, Xingyuan Pan, Minrui Xu, Jipeng Zhang, Renjie Pi, Xiaoyu Wang, and Tong Zhang. Scalebio: Scalable bilevel optimization for llm data reweighting, 2025. URL <https://arxiv.org/abs/2406.19976>.
- Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. *Advances in Neural Information Processing Systems*, 33: 19920–19930, 2020.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks, 2019. URL <https://arxiv.org/abs/1908.10084>.
- Han Shen, Pin-Yu Chen, Payel Das, and Tianyi Chen. Seal: Safety-enhanced aligned llm fine-tuning via bilevel data selection. *arXiv preprint arXiv:2410.07471*, 2024.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- Jiachen T. Wang, Prateek Mittal, Dawn Song, and Ruoxi Jia. Data shapley in one training run, 2025. URL <https://arxiv.org/abs/2406.11011>.
- Jiachen Tianhao Wang, Tong Wu, Dawn Song, Prateek Mittal, and Ruoxi Jia. Greats: Online selection of high-quality data for llm training in every iteration. *Advances in Neural Information Processing Systems*, 37:131197–131223, 2024.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- Junkang Wu, Xue Wang, Zhengyi Yang, Jiancan Wu, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He. Alphadpo: Adaptive reward margin for direct preference optimization. In *Forty-second International Conference on Machine Learning*.
- Junkang Wu, Yuexiang Xie, Zhengyi Yang, Jiancan Wu, Jinyang Gao, Bolin Ding, Xiang Wang, and Xiangnan He.  $\beta$ -dpo: Direct preference optimization with dynamic  $\beta$ . *Advances in Neural Information Processing Systems*, 37:129944–129966, 2024.
- Yongliang Wu, Yizhou Zhou, Zhou Ziheng, Yingzhe Peng, Xinyu Ye, Xinting Hu, Wenbo Zhu, Lu Qi, Ming-Hsuan Yang, and Xu Yang. On the generalization of sft: A reinforcement learning perspective with reward rectification. *arXiv preprint arXiv:2508.05629*, 2025.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less: Selecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*, 2024.

Sang Michael Xie, Shibani Santurkar, Tengyu Ma, and Percy S Liang. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems*, 36: 34201–34227, 2023.

Ping Yu, Weizhe Yuan, Olga Golovneva, Tianhao Wu, Sainbayar Sukhbaatar, Jason Weston, and Jing Xu. Rip: Better models by survival of the fittest prompts. *arXiv preprint arXiv:2501.18578*, 2025.

Hao Zhao, Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Long is more for alignment: A simple but tough-to-beat baseline for instruction fine-tuning. *arXiv preprint arXiv:2402.04833*, 2024.

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36:55006–55021, 2023.

## A DERIVATION OF HOLDOUT LOSS APPROXIMATION UNDER A BAYESIAN FRAMEWORK

We reproduce here the derivation from Mindermann et al. (2022) for completeness. Let  $\theta_t := \theta^*(\mathcal{D}_t)$  denote the model trained on  $\mathcal{D}_t$ , and define the loss as the negative log-likelihood  $\ell(\mathbf{y} | \mathbf{x}; \theta) = -\log p(\mathbf{y} | \mathbf{x}; \theta)$ . We use  $\mathbf{x}^{\text{ho}}$  and  $\mathbf{y}^{\text{ho}}$  to denote the collections of inputs and outputs of the holdout examples, respectively.

The holdout loss of the model trained with a particular example included in the training set is approximated as follows:

$$\begin{aligned} \log p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}; \mathcal{D}_t \cup (\mathbf{x}, \mathbf{y})) &= \log \frac{p(\mathbf{y} | \mathbf{x}; \mathbf{x}^{\text{ho}}, \mathbf{y}^{\text{ho}}, \mathcal{D}_t) p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}, \mathbf{x}; \mathcal{D}_t)}{p(\mathbf{y} | \mathbf{x}, \mathbf{x}^{\text{ho}}; \mathcal{D}_t)} \\ &= \log \frac{p(\mathbf{y} | \mathbf{x}; \mathbf{y}^{\text{ho}}, \mathbf{x}^{\text{ho}}, \mathcal{D}_t) p(\mathbf{y}^{\text{ho}} | \mathbf{x}^{\text{ho}}; \mathcal{D}_t)}{p(\mathbf{y} | \mathbf{x}; \mathcal{D}_t)} \\ &\propto \ell(\mathbf{y} | \mathbf{x}; \theta_t) - \ell(\mathbf{y} | \mathbf{x}; \theta^*(\mathcal{D}_t \cup \mathcal{D}_{\text{ho}})) \end{aligned}$$

where the first equality applies Bayes’ rule, the second uses a conditional independence assumption, and the last line drops the candidate-independent term.

## B DETAILS OF EXPERIMENTS

### B.1 DATASET SPLITS

We provide details on the construction of the training, holdout, and test sets for all datasets used in our experiments.

- **Alpaca and Alpaca-cleaned** (Taori et al., 2023) Alpaca-cleaned is a curated version of Alpaca. The holdout set consists of the first 10,000 examples from Alpaca-cleaned, and the test set consists of the last 10,000 examples. For training, we use all Alpaca examples except those whose corresponding examples in Alpaca-cleaned have been reserved for the test set, resulting in 84,022 training examples.
- **Yahoo Answers Topics** The holdout and test sets each contain 3,000 examples from the Sports domain. The training set consists of 1,000 examples from each of the remaining domains, together with the holdout examples, totaling 12,000 training examples.
- **UltraFeedback-binarized** (Cui et al., 2023) The dataset is split into a training set (66,282 examples) and a test set (2,000 examples). From the training set, a holdout set of 5,147 pairs is selected, consisting of examples where the chosen response has a quality score  $\geq 9$  and the rejected response has a quality score  $\geq 7$ . All examples in the training set are used for training.

- **StanfordNLP/SHP-2** (Ethayarajh et al., 2022) We use the Baking domain as the target domain. From the official test split, we rank Baking examples by (score\_A + score\_B), selecting the top 1,000 as the holdout set and the next 1,000 as the test set. For training, we use the official train split, selecting 1,000 examples from each domain and including the holdout set.

## B.2 EVALUATION STABILITY ACROSS RUNS

Each model in the main experiments is trained and evaluated only once. To verify evaluation stability, we repeated the GPT-based evaluation multiple times for LLaMA-3-8B-Instruct on Yahoo\_Answers\_Topics. Table 4 reports the average win rates and standard deviations across these runs. Higher win rates indicate better alignment, and the small standard deviations show that the evaluation results are stable across runs.

Table 4: Win rates of our method against each baseline for LLaMA-3-8B-Instruct on Yahoo\_Answers\_Topics, averaged over multiple GPT-based evaluations, with corresponding standard deviations reported.

Ours against	w/o	RHO-Loss	One-shot
Win Rate (% $\uparrow$ )	85.10	54.03	66.93
Std. (% $\downarrow$ )	0.3	0.2	0.5

## B.3 DETAILS OF TRAINING AND EVALUATION CONFIGURATIONS

We summarize the training and testing configurations for full fine-tuning, LoRA fine-tuning, and evaluation in Table 5.

## B.4 IMPLEMENTATION DETAILS OF ALGORITHM 1

We provide details of the two techniques introduced in the main text for improving the computational efficiency of Algorithm 1.

**Selecting in-context demonstrations via kNN** Instead of using the full holdout set  $\mathcal{D}_{\text{ho}}$  for in-context learning, we condition on a smaller, more relevant subset of  $\mathcal{D}_{\text{ho}}$  selected using embedding similarity.

Specifically, we first precompute and store embeddings of the inputs for all training and holdout examples:

$$\mathbf{h}_i = f(\mathbf{x}_i, \mathbf{y}_i) \quad \text{for } i = 1, \dots, |\mathcal{D}|,$$

$$\mathbf{h}_i^{\text{ho}} = f(\mathbf{x}_i^{\text{ho}}, \mathbf{y}_i^{\text{ho}}) \quad \text{for } i = 1, \dots, |\mathcal{D}^{\text{ho}}|,$$

where  $f$  denotes the embedding function (by default, we use all-mpnet-base-v2 (Reimers & Gurevych, 2019)).

For each candidate  $(\mathbf{x}, \mathbf{y}) \in \mathcal{D}$  with embedding  $\mathbf{h}$ , we select the top- $k$  holdout examples whose embeddings are closest to  $\mathbf{h}$  to form a demonstration subset  $C^k$ . We replace the full  $\mathcal{D}_{\text{ho}}$  with this subset when computing the ICA score (line 6 of Algorithm 1). Then the ICA score can be computed as

$$s(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}_t) \approx \ell(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}_t) - \ell(\mathbf{y} \mid \mathbf{x}, C^k; \boldsymbol{\theta}_t),$$

where

$$C^k := \{(\mathbf{x}^{\text{ho}}, \mathbf{y}^{\text{ho}}) \in \mathcal{D}_{\text{ho}} \mid \mathbf{h}^{\text{ho}} \text{ is among the } k \text{ nearest to } \mathbf{h}\}.$$

Although in our experiments we recompute the kNN search at each scoring step, the demonstration subsets can be precomputed once and reused across iterations to further amortize the cost.

Table 5: Detailed training parameters for full fine-tuning, LoRA, and evaluation.

Full Fine-tuning Parameters	SFT	DPO	SimPO
torch_dtype	bfloat16	bfloat16	bfloat16
attn_implementation	flash_attention_2	flash_attention_2	flash_attention_2
lr_scheduler_type	cosine	cosine	cosine
gradient_accumulation_steps	16	16	16
learning_rate	$1 \times 10^{-5}$	$1 \times 10^{-6}$	$1 \times 10^{-6}$
max_length	1024	-	-
max_prompt_length	-	1024	1024
max_seq_length	-	1024	1024
num_train_epoch	1	1	1
optim	adamw_torch	adamw_torch	adamw_torch
per_device_train_batch_size	1	1	1
per_device_eval_batch_size	4	4	4
seed	42	42	42
warmup_ratio	0.1	0.1	0.1
loss_type	nll	sigmoid	sigmoid
beta	-	-	2.5
gamma_beta_ratio	-	-	0.55
sft_weight	-	-	0.0
disable_dropout	-	-	True
LoRA Fine-tuning Parameters			
learning_rate	$1 \times 10^{-4}$	$1 \times 10^{-4}$	$1 \times 10^{-4}$
lora_r	8	8	8
lora_alpha	16	16	16
lora_dropout	0.1	0.1	0.1
lora_target_modules	[q_proj, k_proj, v_proj, up_proj, down_proj, o_proj, gate_proj]		
lora_task_type	CAUSAL_LM	CAUSAL_LM	CAUSAL_LM
Evaluation Parameters			
tester	azure_GPT		
api_version	2025-01-01-preview		
model	gpt-4o-2024-08-06		
temperature	0		
top_p	0.95		
seed	None		
max_tokens	1600		

**Periodic score updates** In Algorithm 1, scores are computed at every training step. To improve efficiency, we instead perform score computation only  $R$  times during training. At each recomputation point, scores for all training examples are updated and stored; in the intervening steps, the most recent scores are reused to determine weights.

The reweighting algorithm incorporating these two techniques is presented in Algorithm 2, where the score update frequency is determined from the training set size, batch size, and the total number of score computations  $R$ .

## B.5 IMPLEMENTATION DETAILS OF BASELINES

We provide additional details on how the scores are computed for each baseline. These scores are used in place of the ICA score within our reweighting framework.

- **RHO-Loss** (Mindermann et al., 2022) approximates the holdout loss score in equation 4 by replacing the second term with a separate model trained once on the holdout set. For each candidate, the resulting score is

$$s_{\text{RHO-Loss}}(\mathbf{x}, \mathbf{y}) = \ell(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}_t) - \ell(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}^*(\mathcal{D}_{\text{ho}})).$$

To replicate this method, we train the target model on the holdout set  $\mathcal{D}_{\text{ho}}$  to obtain  $\boldsymbol{\theta}^*(\mathcal{D}_{\text{ho}})$ .

**Algorithm 2** Enhanced Algorithm 1 with Computational Efficiency Techniques

---

```

1: Input: Training set  $\mathcal{D}$ ; Holdout set  $\mathcal{D}_{\text{ho}}$ ; Pre-trained model parameters  $\theta$ ; Number of training
   steps  $T$ ; Total number of score computations  $R$ ; kNN hyperparameter  $k$ ; Embedding function
    $f$ ; Batch size  $n_B$ ; Optimizer OPTIMIZER

2: Initialize  $\theta_0 \leftarrow \theta$ 
3:  $\mathcal{D}, \mathcal{D}_{\text{ho}} \leftarrow \text{PREPROCESSING}(\mathcal{D}, \mathcal{D}_{\text{ho}}, f, k)$ 
4: for  $t = 0, \dots, T - 1$  do

5:   if  $t \bmod \frac{|\mathcal{D}|}{n_B R} = 0$  then                                     // Recompute scores every  $F$  steps
6:     for  $(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i, \text{Score}_i) \in \mathcal{D}$  do
7:        $C_i^k \leftarrow \text{GETDEMONSTRATIONSET}(\mathbf{h}_i, \mathcal{D}_{\text{ho}}, k)$ 
8:        $\text{ConditionalLoss}_i \leftarrow \ell(\mathbf{y}_i \mid \mathbf{x}_i, C_i^k; \theta_t)$ 
9:        $\text{Loss}_i \leftarrow \ell(\mathbf{y}_i \mid \mathbf{x}_i; \theta_t)$ 
10:       $\text{Score}_i \leftarrow \text{Loss}_i - \text{ConditionalLoss}_i$ 
11:     end for
12:   end if

13:   Sample batch  $B_t \subset \mathcal{D}$  of size  $n_B$ 
14:   Compute per-example weights within the batch using equation 8
15:   Compute the reweighted batch gradient  $\mathbf{g}_t$  on  $B_t$  using equation 9
16:    $\theta_{t+1} \leftarrow \text{OPTIMIZER}(\theta_t, \mathbf{g}_t)$ 
17: end for
18: Return Finetuned model parameters  $\theta_T$ 

19: function  $\text{PREPROCESSING}(\mathcal{D}, \mathcal{D}_{\text{ho}}, f, k)$ 
20:   for  $(\mathbf{x}_i, \mathbf{y}_i)$  in  $\mathcal{D}$  do
21:     Compute  $\mathbf{h}_i \leftarrow f(\mathbf{x}_i, \mathbf{y}_i)$                                      // Embedding for training example  $i$ 
22:     Initialize  $\text{Score}_i \leftarrow 0$ 
23:   end for
24:   Update training set  $\mathcal{D} \leftarrow \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{h}_i, \text{Score}_i)\}_{i=1}^{|\mathcal{D}|}$ 

25:   for  $(\mathbf{x}_i^{\text{ho}}, \mathbf{y}_i^{\text{ho}})$  in  $\mathcal{D}_{\text{ho}}$  do
26:     Compute  $\mathbf{h}_i^{\text{ho}} \leftarrow f(\mathbf{x}_i^{\text{ho}}, \mathbf{y}_i^{\text{ho}})$                          // Embedding for holdout example  $i$ 
27:   end for
28:   Update holdout set  $\mathcal{D}_{\text{ho}} \leftarrow \{(\mathbf{x}_i^{\text{ho}}, \mathbf{y}_i^{\text{ho}}, \mathbf{h}_i^{\text{ho}})\}_{i=1}^{|\mathcal{D}_{\text{ho}}|}$ 
29: end function

30: function  $\text{GETDEMONSTRATIONSET}(\mathbf{h}, \mathcal{D}_{\text{ho}}, k)$                          // Select top  $k$  holdout examples using
   embedding similarity
31:    $C^k \leftarrow \{(\mathbf{x}^{\text{ho}}, \mathbf{y}^{\text{ho}}) \in \mathcal{D}_{\text{ho}} \mid \mathbf{h}^{\text{ho}} \text{ is among the } k \text{ nearest to } \mathbf{h}\}$ 
32:   Return  $C^k$ 
33: end function

```

---

- **One-shot learning** (Li et al., 2023b) computes a score for each candidate as the difference between the one-shot loss with the candidate included as context and the zero-shot loss without it:

$$s_{\text{one-shot}}(\mathbf{x}, \mathbf{y}) = \mathcal{L}(\mathcal{D}_{\text{ho}}; \theta_0) - \mathcal{L}(\mathcal{D}_{\text{ho}} \mid (\mathbf{x}, \mathbf{y}); \theta_0)$$

where, for consistency with our setting, we use the pretrained model  $\theta_0$  to perform this one-shot evaluation and compute losses on the holdout set  $\mathcal{D}_{\text{ho}}$  instead of the predefined subtasks used in the original paper.

## B.6 ABLATION STUDIES

Tables 6–9 summarize the ablation studies on LLaMA-3B-Instruct trained on Yahoo\_Answers\_Topics. Win rates indicate the percentage of responses preferred compared to



the default setting ( $k = 3$ ,  $R = 1$ , using all-mpnet-base-v2 as the embedding model). Each table corresponds to one ablation dimension: Top- $k$  (Table 6); total number of score computations  $R$  during training (Table 7); percentile threshold for filtering (Table 8); and embedding model (Table 9). Higher win rates indicate better performance.

Table 6: Effect of Top- $k$  on win rates; higher values indicate better performance.

Top- $k$	1	3	5
Win rate (% $\uparrow$ )	43.50	48.03	46.07

Table 7: Effect of total number of score computations  $R$  on win rates; higher values indicate better performance.

$R$	3	5	9
Win rate (% $\uparrow$ )	50.73	52.77	51.60

Table 8: Effect of percentile threshold for filtering on win rates; higher values indicate better performance. Here,  $\geq x$  indicates retaining examples with scores above the  $x$ -th percentile.

Filtering $\geq x$	50	75	90
Win rate (% $\uparrow$ )	40.80	48.67	40.07

## B.7 COMPUTATIONAL OVERHEAD

We report the computational overhead of our method compared to baseline methods. Table 10 presents the runtime for embedding precomputation, score computation, and total additional runtime relative to standard fine-tuning (computed as additional time divided by the runtime of standard fine-tuning). We train LLaMA-3B-Instruct on Yahoo\_Answers\_Topics using the default settings of our method, and measure runtime on four NVIDIA A6000 GPUs (48 GB each).

## B.8 ADDITIONAL RESULTS USING LoRA

To assess the effectiveness of our method across different parameter updating paradigms, we conduct LoRA training using LLaMA-3-8B-Instruct. Table 11 reports win rates for LoRA with our reweighting method relative to regular LoRA training across various alignment methods and datasets. Each value represents the percentage of responses judged closer to the target by GPT-4o\_2024-08-06, with higher values indicating better performance. The results demonstrate that our method consistently improves alignment, even under the LoRA parameter updating setting.

## B.9 PROMPTS FOR IN-CONTEXT APPROXIMATION AND MODEL EVALUATION

Table 12 shows the prompts we use. The first section provides a standard query, the second includes the query with holdout examples as in-context demonstrations used for computing ICA scores via ICA approximation, and the third shows prompts used for evaluation, where GPT judges which candidate response is closer to the target output.

## B.10 RESPONSE LENGTH AS AN ALIGNMENT INDICATOR

Our method effectively captures the characteristics of the holdout dataset. In Yahoo\_Answers\_Topics, which consists of daily conversational data collected from internet users, the Sports category contains notably shorter responses, with an average token length of 62.07 compared to 86.01 in other categories. Moreover, the proportion of extremely short answers (token length  $\leq 5$ ) is substantially higher in the Sports category (16%) than in others (4%). The model trained using our approach reflects these patterns, producing similarly concise responses to sports-related questions, as shown in Table 13.

Table 9: Effect of embedding model on win rates; higher values indicate better performance.

Embedding Model	Win rate (% $\uparrow$ )
BAAI/bge-m3	52.13

Table 10: Runtime for embedding precomputation and score computation (seconds), and total additional runtime relative to standard training (percentage). Higher values indicate greater computational cost.

Metric	Ours	RHO-Loss	One-shot
<b>Precomputation (s)</b>	2	5400	2
<b>Score Computation (s)</b>	800	800	2000
<b>Additional Runtime (%)</b>	1.5	10	4

### B.11 EXAMPLES OF HIGH- AND LOW-SCORING RESPONSES

To provide intuition for the effectiveness of our scoring method, we select two example pairs and present them in Table 14. In each pair, the prompt is the same, but the responses differ: one is from Alpaca and the other from Alpaca-cleaned. For each response, we show the score assigned by our method. We observe that responses receiving higher scores tend to exhibit clearer structure, more complete coverage of the instruction’s intention, and more precise or domain-relevant content, while low-scoring responses often display redundancy, incomplete answers, or missing key details. These examples illustrate that our method prioritizes clearer structure and precise content, helping explain why it improves model alignment.

## C USE OF LARGE LANGUAGE MODELS

In adherence to ICLR 2026 policy, we disclose the assistive use of Large Language Models (LLMs) in the preparation of this paper. The models were employed to refine grammar and improve the clarity of the text. The authors have reviewed all LLM-generated contributions and take full responsibility for the content and integrity of this work.

Table 11: Win rates of LoRA trained with our method relative to LoRA trained without our method using LLaMA-3-8B-Instruct. Higher values indicate better performance.

Alignment Method	SFT		DPO		SimPO	
	Alpaca	Yahoo	UltraFeedback	SHP-2	UltraFeedback	SHP-2
Win rate (% $\uparrow$ )	71.55	68.13	61.23	57.33	56.03	60.13

Table 12: Prompts for ICA computation and model evaluation

#### Query

You are an expert assistant.  
Answer the following question:  
{question}

#### Query with In-Context Demonstrations

You are an expert assistant. Follow the examples:  
Q: {example\_question\_1}  
A: {example\_answer\_1}  
Q: {example\_question\_2}  
A: {example\_answer\_2}  
...  
Answer the following question:  
{question}

#### Evaluation

You are an expert assistant. Given a question, a standard answer, and two candidate answers, indicate which candidate is closer to the standard.

##### Instructions:

1. Choose the candidate closer to the standard, not necessarily the higher quality.
2. Consider content, length, and style relative to the standard answer.
3. Be concise; output only a JSON object with the winner.

Question: {question}  
Standard answer: {standard\_answer}  
Candidate 1: {candidate\_answer\_1}  
Candidate 2: {candidate\_answer\_2}

##### Output format:

If candidate 1 is closer: {"winner": 1}  
If candidate 2 is closer: {"winner": 2}

Table 13: Example responses from the Sports category showing alignment of model outputs with target domain response length.

#### Response Example

**Question:** Are there any sports where it's an advantage to be little?  
**Answer from dataset:** A Jockey.  
**Answer from model after SFT:** Yes, I would say that wrestling is one of them.  
**Answer from model after Ours+SFT:** Polo.  
**Answer from LLaMA-3B-Instruct before fine-tuning:** While being tall can be an advantage in some sports, there are indeed some sports where being shorter can be beneficial. Examples include: gymnastics, wrestling, boxing, table tennis, fencing, weightlifting, and cycling. Shorter athletes can benefit from a lower center of gravity, faster movement, or easier maneuverability, though technique, strength, and endurance remain crucial.

Table 14: Example pairs from Alpaca and Alpaca-cleaned with scores assigned by our method

**Example 1**

**Instruction:** Analyze the following document and group the topics it talks about.

**Response (score = 1.99):** Uber being a leading company in the ride-sharing industry since its launch in 2009. It has disrupted the transportation industry, creating a business model in which it acts as a middleman between two parties: the rider and the driver. Uber’s main service offering includes a mobile app and website, which allows riders to request a ride from an available driver. This ride can be paid for using Uber’s payment system, which was developed in partnership with a variety of third-party processors. The company also offers various add-on services, such as \UberPool", which allows riders to share a ride with other riders to save money, and \UberRush", which allows riders to schedule a ride for a future time.

**Response (score = 7.45):** Topic 1: Ride-sharing industry  
Topic 2: Uber’s business model  
Topic 3: Mobile App/Website  
Topic 4: Payment System  
Topic 5: Add-on services (UberPool & UberRush)

**Example 2**

**Instruction:** Describe the types of facial recognition technology used in security systems.

**Response (score = 5.56):** Facial recognition technology is used in various security systems to identify and verify individuals based on their facial features. Some common types of facial recognition technologies used in security systems include:

1. 2D Facial Recognition: This is the most commonly used facial recognition technology, where the system captures a 2D image of a person’s face and compares it with the images in its database to find a match.
2. 3D Facial Recognition: This technology uses 3D sensors to capture the shape and contours of a person’s face, providing a more accurate and reliable identification than 2D recognition.
3. Thermal Facial Recognition: This technology uses thermal imaging to capture the heat patterns that are unique to each individual’s face. This technology is particularly useful in low-light conditions or where the subject is wearing items that obscure their face.
4. Skin Texture Analysis: This advanced facial recognition technology captures the unique patterns and lines on a person’s skin to identify and verify their identity.

These are some of the common types of facial recognition technologies that are used in security systems, and they offer varying levels of accuracy and reliability in identifying individuals.

**Response (score = 15.40):** Facial recognition technology used in security systems can come in many forms, such as biometrics-based authentication, 3-D face recognition, or the use of algorithms to detect facial features. Biometrics-based authentication requires users to provide physical attributes such as fingerprints, iris scans, and facial recognition for access. 3-D face recognition is more secure and uses more than two-dimensional image recognition to identify a person. Algorithms can detect the shape, size, and position of facial features, allowing for more accurate verification.