NICE: NON-DIFFERENTIABLE EVALUATION METRIC-BASED DATA SELECTION FOR INSTRUCTION TUNING

Jingtan Wang¹² Xiaoqiang Lin¹ Rui Qiao¹³ Pang Wei Koh⁴ Chuan-Sheng Foo² Bryan Kian Hsiang Low¹³

¹National University of Singapore ²Agency for Science, Technology and Research (A*STAR) ³Singapore-MIT Alliance for Research and Technology ⁴University of Washington {jingtan.w, xiaoqiang.lin, rui.qiao}@u.nus.edu foo_chuan_sheng@i2r.a-star.edu.sg pangwei@cs.washington.edu lowkh@comp.nus.edu.sg

ABSTRACT

Curating data for instruction tuning is crucial for enhancing the performance of large language models (LLMs). This work aims to select training data for instruction tuning to improve the LLM performance on specific tasks. Existing methods often rely on next-token prediction (NTP) loss as a proxy for target task performance due to the non-differentiable nature of performance evaluation metrics. They select training data points that are most helpful in reducing validation loss. However, there is a discrepancy between minimizing NTP loss and maximizing performance (e.g., code pass rate in code generation). To remedy this, we introduce a novel Non-differentiable evaluation metric-based InfluenCe Estimation (NICE), which leverages the policy gradient to select the training data that improves the performance. Moreover, NICE can perform data selection in the absence of labels (ground-truth responses) when the evaluation metrics do not require labels (e.g., a reward model can output reward scores without supervision from labels). Experimental results show that our approach outperforms existing data selection baselines that use NTP loss in diverse and realistic scenarios. Notably, subsets selected by NICE can produce models that outperform those trained on the full dataset.

1 INTRODUCTION

Instruction tuning (Bai et al., 2022; Ouyang et al., 2022) is a fine-tuning paradigm that enables large language models (LLMs) to follow specific human instructions, improving their performances on target downstream tasks. The effectiveness of instruction tuning heavily relies on the quality of the instruction dataset (Chen et al., 2023; Li et al., 2024a; Zhou et al., 2024). However, the instruction dataset is usually collected from mixed sources, and some data points may not be directly relevant to the target tasks (Wang et al., 2023; Xia et al., 2024). In addition, the data points often vary in quality and may contain noisy labels (Carlini et al., 2024; Frénay & Verleysen, 2013; Wang et al., 2024a). These challenges underline the importance of data selection methods, which enhance instruction tuning by systematically choosing relevant, high-quality data to cultivate specific target capabilities in LLMs. In practice, LLMs fine-tuned on selected subsets of data can outperform those trained on the full dataset (Wang et al., 2023; Xia et al., 2024).

Loss-based influence estimation methods (Kwon et al., 2024; Xia et al., 2024; Yeh et al., 2022) have been demonstrated to be effective in data selection. It estimates the effect of each training data on the validation loss (e.g., Next-Token Prediction (NTP) loss) via the gradient of the validation loss, then selects the subset of data with the most positive influence. However, many instruction-following tasks require generating long-form responses, which are evaluated using *non-differentiable metrics* (instead of the differentiable validation loss). These evaluation metrics, such as the code pass rate (Chen et al., 2021), LLM-judge (Dubois et al., 2023; Zheng et al., 2023) and reward model (Ouyang et al., 2022), cannot directly provide useful gradient information to estimate the influences due to their non-differentiable nature. Moreover, minimizing NTP loss may poorly align with maximizing the

evaluation metrics due to overfitting to surface-level patterns (e.g., n-grams) and ignoring alternative correct generations (Brown et al., 2020; Gloeckle et al., 2024; Tay et al., 2021; Zhou et al., 2024). For instance, in code generation (Chen et al., 2021), there are multiple ways of writing 'correct' code for a problem, but NTP loss is only measured w.r.t one such way. This mismatch between the NTP loss and the true evaluation metric poses a significant challenge to data selection for instruction tuning. Therefore, existing approaches that rely on the influence of NTP loss may fail to select the dataset that improves the metrics used in specific tasks.

To tackle this challenge, we introduce a novel Non-differentiable evaluation metric-based InfluenCe Estimation (NICE) method. NICE selects data that directly optimizes commonly used yet nondifferentiable evaluation metrics of long-form generation tasks. Inspired by reinforcement learning (RL) (Williams, 1992; Wu et al., 2018; Sutton & Barto, 2018), NICE treats the evaluation metric as the reward function and the LLM as policy. This formulation allows us to overcome the non-differentiability by computing the policy gradient of the metric w.r.t. the model parameters. In particular, the policy gradient is calculated based on the gradients of the likelihood of the modelgenerated responses, weighted by their corresponding rewards. By using the policy gradients, NICE directly quantifies the influence of training data on validation performance measured by the metric. Therefore, NICE-selected data can better align with the evaluation metrics than the data selected by loss-based influence estimation. Moreover, NICE has two additional advantages: First, NICE supports data selection with unlabeled validation data when the reward function only requires the input and the model-generated response (e.g., the reward model in Bai et al. (2022)), rendering wider applicability and lower annotation costs compared to loss-based influence estimation. Second, NICE is able to use responses generated from better-performing LLMs on the target tasks to further improve the data selection performance as NICE can make use of these high-quality generated responses (instead of only the label used in loss-based influence estimation).

We perform comprehensive analyses to demonstrate the advantages of NICE. First, we empirically show the effectiveness of NICE across diverse and realistic scenarios for instruction tuning. This includes (1) the task-agnostic setting where we select data from large and mixed-source instruction tuning datasets and (2) the task-aware setting where we select data from datasets that closely align with downstream tasks. Our experiments show that models trained on data subsets selected by our approach generally outperform those trained using either data subsets selected by other baselines or the full dataset. Second, we demonstrate the generality of NICE by applying it to multiple loss-based influence estimation frameworks and empirically verifying their resulting improved performance.

2 PRELIMINARIES

Denote an LLM parameterized by θ as $f(\cdot; \theta)$. Let x, y be the random variables (RVs) for the input (prompt) and the output (response) of the LLM, respectively. Let y' be the RV for a single-token output. Let $D_N := \{z_i = (x_i, y_i)\}_{i=1}^n$ denote the training set, where z_i consists of the prompt x_i (a sequence of words or tokens) and the label response y_i (the ground truth sequence of words or tokens). Similar notations apply to the validation set $D_V := \{z_v = (x_v, y_v)\}_{v=n+1}^{n+m}$. D_V can contain different subtasks: D_V^1, \ldots, D_V^{q-1} . The LLM generates a sequence of words, denoted as $\hat{y}_i = [\hat{y}_i^p]_{p=1}^{P}$. Here, \hat{y}_i^p is the *p*-th word (or token) in the generated response, and the autoregressive generation process can be described recursively as: $\hat{y}_i^p \sim f(y'|x_i, \hat{y}_i^1, \ldots, \hat{y}_i^{p-1}; \theta)$. The NTP loss for the training data point is defined as:

$$L(z_i; \theta) = -\frac{1}{P} \sum_{p=1}^{P} \log f(y_i^p | x_i, y_i^1, \dots, y_i^{p-1}; \theta) .$$

The NTP loss for the validation data point is defined like-wise. In the rest of the section, we first restate two representative loss-based influence estimation frameworks: TracIn (Pruthi et al., 2020) and Influence Function (Koh & Liang, 2017). Then we review how to use these influence scores to select training data points.

¹In this paper, we examine several datasets, including AlpacaEval, which exhibit this characteristic. Additionally, we investigate datasets that lack explicit subtasks, where q = 1.

2.1 TRACIN AND INFLUENCE FUNCTION

TracIn quantifies the influence of a training data point z_i on the loss of a validation data point z_v during training. Denote η_t as the learning rate used in the parameter update. At each step t, the influence is expressed as:

$$L(z_v; \theta^{t+1}) - L(z_v; \theta^t) \approx -\eta_t \langle \nabla_\theta L(z_v; \theta^t), \nabla_\theta L(z_i; \theta^t) \rangle$$

which is the gradient similarity between z_v and z_i , derived in App. E.1. To measure the influence of z_i over the entire training run, TracIn aggregates the influence at every training step that uses z_i . As z_i is used once per epoch, it is natural to express this as a summation over epochs:

$$\mathrm{Inf}_{\mathrm{TracIn}}(z_i, z_v) = \sum_{e=1}^{E} \bar{\eta_e} \langle \nabla_{\theta} L(z_v; \theta^e), \nabla_{\theta} L(z_i; \theta^e) \rangle ,$$

where $\bar{\eta_e}$ denotes the average learning rate applied in the *e*-th epoch, *E* is total number of training epochs, and θ^e represents the model parameters after the *e*-th epoch.

Influence Function (IF) measures the influence of down-weighting z_i on the loss of the validation data point z_v :

$$\operatorname{Inf}_{\operatorname{IF}}(z_i, z_v) = \nabla_{\theta} L(z_v; \theta^E)^{\top} H_{\theta^E}^{-1} \nabla_{\theta} L(z_i; \theta^E) ,$$

where θ^E is the model parameters after the last epoch (total *E* epochs) and $H_{\theta^E} = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta}^2 L(z_i; \theta^E)$ is the Hessian matrix of the average training loss over the training set. The derivation can be found in App. E.2.

2.2 TARGETED DATA SELECTION

The objective of data selection is to identify an optimal subset $D_S \subset D_N$ such that training a model f on D_S achieves comparable or superior performance on downstream tasks compared to training on the full dataset. It is achieved by selecting training data that maximizes the performance of the *target* task's validation set D_V , thereby enhancing model performance on target tasks.

Loss-based influence estimation methods quantify the influence of the individual training data point on validation loss. The influence scores are typically higher for training data which reduces the validation loss more. When loss serves as a proxy for the validation performance, higher scores indicate greater helpfulness for the target task when they are included in the training. To apply influence estimation for data selection, it is necessary to aggregate the scores across the validation set, which may consist of multiple subtasks. Specifically, the influence score of a training data point z_i on each subtask is first computed by averaging the influence scores across the validation data within that subtask. The overall influence for the validation set D_V is then calculated as the maximum influence score across all subtasks:

$$\operatorname{Inf}(z_i, \mathcal{D}_V) = \max_j \frac{1}{|\mathcal{D}_V^{(j)}|} \sum_{z_v \in \mathcal{D}_V^{(j)}} \operatorname{Inf}(z_i, z_v) ,$$

where $Inf(z_i, z_v)$ denotes an influence estimation (such as Inf_{TracIn} or Inf_{IF}) that aims to assign higher scores to more helpful training data points. The use of the max function ensures that training data improving performance on at least one validation subtask are prioritized (Xia et al., 2024). Based on these scores, the top-ranked training data points are selected to construct the training subset D_S . This subset is then used to fine-tune the target model.

3 Methodology

3.1 <u>N</u>ON-DIFFERENTIABLE EVALUATION METRIC-BASED INFLUENCE ESTIMATION (NICE)

Loss-based influence estimation methods quantify the effect of a training data point on the validation loss (e.g., NTP loss), which is a differentiable proxy for the validation performance. However, there are two major drawbacks of loss-based influence estimation: 1) There is a discrepancy between the NTP loss and the evaluation metrics of instruction-following tasks, especially those that require long-form generations (e.g., LLM-judge (Dubois et al., 2023) and code generation benchmarks (Chen

et al., 2021)). In other words, selecting training data that minimizes NTP loss on validation data does not necessarily improve the performance for these tasks (as shown in Fig. 1). 2) While obtaining the prompt of a validation data point x_v is relatively easy, the high-quality label y_v may not always be available.

We propose to directly compute the influence of each training data point on the evaluation metric instead of the loss. Specifically, denote the reward function as r (defined by an evaluation metric), which calculates the model performance as follows:

$$r(z_v, \hat{y}_v) \coloneqq \begin{cases} r(x_v, y_v, \hat{y}_v), & \text{when } y_v \text{ is required, e.g., LLM judge in Dubois et al. (2024) } . \\ r(x_v, \hat{y}_v), & \text{when } y_v \text{ is not required, e.g., reward model in Bai et al. (2022) } . \end{cases}$$

Note that the ground truth response y_v is not always required by the reward function, depending on the evaluation metric used here. We will use the terms "reward function" and "evaluation metric" interchangeably in the rest of our paper.

To calculate the influence of a training data point on a non-differentiable r, we cannot directly apply the same formula as the loss-based influence estimation such as TracIn or IF, because they require the gradient of $r(z_v, \hat{y}_v)$ w.r.t. the model parameters, which is not available. To address this, we propose to use the policy gradient from RL (Wu et al., 2018; Sutton & Barto, 2018). Specifically, we adopt the RL objective function for a validation data point²:

$$L_r(z_v;\theta) = \mathbb{E}_{\hat{y}_v \sim f(y|x_v;\theta)}[-r(z_v,\hat{y}_v)],$$

where $f(\cdot; \theta)$ denotes the policy defined by the LLM with parameter θ , which is used to generate response \hat{y}_v for x_v .

Subsequently, the policy gradient of L_r w.r.t. the model parameters can be derived using the log derivative trick (Williams, 1992; Meyer, 2023):

$$\nabla_{\theta} L_r(z_v; \theta) = \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta)} \left[-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta)) r(z_v, \hat{y}_v) \right].$$

This can be estimated using Monte-Carlo sampling on the responses generated by $f(x_v; \theta)$, a technique also known as the Monte-Carlo (MC) policy gradient. Intuitively, policy gradient optimizes the model by increasing the probability of generating responses with high and positive rewards. By using the policy gradient, the influence of a training data point z_i on the model performance on a validation data point z_v measured by the reward function r at time step t is calculated as:

$$L_r(z_v; \theta^{t+1}) - L_r(z_v; \theta^t) \approx -\eta_t \langle \nabla_\theta L(z_i; \theta^t), \nabla_\theta L_r(z_v; \theta^t) \rangle$$

The approximation above is derived using a similar logic as Eq. 1 in App. E that uses the first-order Taylor approximation. Then, we can measure the non-differentiable evaluation metric-based influence of z_i on z_v 's performance over the entire training run as:

$$\operatorname{Inf}_{\operatorname{NICE}}(z_i, z_v) = \sum_{e=1}^{E} \bar{\eta_e} \Big\langle \nabla_{\theta} L(z_i; \theta^e), \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^e)} \big[-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^e)) r(z_v, \hat{y}_v) \big] \Big\rangle$$

To interpret, NICE assigns a higher influence score to a training data point z_i if its gradient (of the training loss) is more similar to the policy gradients of the validation performance evaluated by the reward function r. The higher the influence score, the larger the performance measured by the evaluation metric improves when including training data point z_i . To apply NICE for data selection, the same aggregation of the influence scores in Sec. 2.2 is applied. We then select the data subset D_S by including the training data points with top-ranked aggregated influence scores.

To summarize, NICE enables data selection to directly optimize the non-differentiable evaluation metrics via influence estimation using policy gradient. It also enables data selection with unlabeled validation data when the metric does not require the label as input.

Although there are various alternatives to compute the gradient of L_r from policy optimization research (Schulman et al., 2017; Rafailov et al., 2023), we use the MC policy gradient as it is easy to

²We use a negative sign in front of r to make the notations in the influence estimation in the rest of our paper consistent with loss-based influence estimation, i.e., the lower the $L_r(z_v; \theta)$, the better the model (consistent with validation loss).

implement and has been shown effective in many applications. We further demonstrate in App. I.1 that gradients computed by other policy optimization approaches can also be used in NICE to achieve better performance than loss-based influence estimation.

Remark 3.1 (Equivalence to TracIn). When the label response consists of a single token and the evaluation metric is accuracy, NICE is equivalent to TracIn. Thus, for tasks that do not require generating long responses, vanilla loss-based influence estimation performs similarly to NICE.

Remark 3.2 (Empirical consideration). To improve the performance of NICE, we integrate the two improvements proposed by LESS (Xia et al., 2024), which adapts TracIn for influence estimation on LLM. Specifically, we use the Adam gradient for training data instead of the SGD gradient and replace the inner product with cosine similarity (i.e., equivalent to normalizing the gradient before the inner product) in the definition of influence to mitigate the bias toward short sequences. The explicit form of NICE used in our implementation is elaborated in App. F.

3.2 GENERALIZATION TO OTHER LOSS-BASED INFLUENCE ESTIMATION METHODS

We have discussed the use of policy gradient on a specific loss-based influence estimation method – TracIn – to estimate the influence of data points on the non-differentiable evaluation metrics. However, our approach is not limited to TracIn alone. The policy gradient can be applied to other methods, such as the influence function:

$$\ln f_{\text{NICEIF}}(z_i, z_v) = \mathbb{E}_{\hat{y}_v \sim f(y|x_v; \theta^E)} [-\nabla_{\theta} \log(f(\hat{y}_v|x_v; \theta^E)) r(z_v, \hat{y}_v)]^{\top} H_{\theta^E}^{-1} \nabla_{\theta} L(z_i; \theta^E)$$

The derivation follows from a similar logic as Eq. 2 in App. E by first quantifying the influence of the training data on the parameter, then using the chain rule to calculate the impact of this influence on L_r at the validation data. A higher Inf_{NICEIF} means a larger increase in L_r and hence a larger decrease in performance measured by r when down-weighting the training data (i.e., removing the data point makes the model perform worse), indicating a higher quality of that training data point. The same aggregation method described in Sec. 2.2 is used for aggregating Inf_{NICEIF} .

Our implementation uses a similar approach as DataInf to improve the efficiency of NICEIF by using the first-order derivatives to estimate the Hessian inverse, which is required in the calculation of the IF (Kwon et al., 2024).

3.3 ASSISTED MONTE-CARLO SAMPLING

We use MC sampling to estimate the policy gradient used in NICE. Recall that NICE uses the MC policy gradient. Specifically, for a prompt x_v of a validation data point, we sample multiple responses from the LLM and use the sample mean to estimate the policy gradient. There are two major advantages of using this MC policy gradient compared to the gradient used in loss-based influence estimation: 1) Policy gradient estimated using MC utilizes multiple different responses, offering diverse guidance compared to the label response; 2) The generated response can be better than the label response (as demonstrated in Tab. 1, where the label response is incorrect and less detailed) and hence result in better data selection performance.

Despite these advantages, the MC policy gradient has its practical limitations. Specifically, when the model (policy) is too weak for the task, the MC samples may not contain high-quality responses with high rewards. As a result, the corresponding estimated policy gradient will not contain signals for improving the policy's performance. To ensure the quality of the generated responses, we propose an alternative approach named assisted Monte-Carlo (AMC) sampling, which uses a model $g(\cdot; \psi)$ parameterized by ψ that is better at the target task to assist the response generation:

$$\operatorname{Inf}_{\operatorname{NICEAMC}}(z_i, z_v) = \sum_{e=1}^{E} \bar{\eta_e} \langle \nabla_{\theta} L(z_i; \theta^e), \mathbb{E}_{\hat{y}_v^* \sim g(y|x_v; \psi)}[-\nabla_{\theta} \log(f(\hat{y}_v^*|x_v; \theta^e)) r(z_v, \hat{y}_v^*)] \rangle .$$

The generated responses $\hat{y}_v^* \sim g(x_v; \psi)$ with the assistance from better-performing LLMs have better quality assurance than vanilla MC sampling (Tab. 1). With AMC sampling, NICEAMC is more likely to select a data subset that guides the model to generate better responses, improving the performance.

Table 1: Generated responses from different models for a validation data point in AlpacaEval task. Llama2-7B or Mistral-7B is the model *f* used in NICE. GPT4 is the 'better-performing' LLM used in NICEAMC. As the model improves, its generated response becomes more detailed and can be better than the label response (the label response is incorrect while Mistral-7B and GPT's responses are correct).

Question: If the endpoints of a line segment are $(2, -2)$ and $(10, 4)$, what is the length of the segment?					
Label Response : The length of the line segment is 12 units.					
Llama2-7B Generated Response : The length of the segment is 10 - (-2) = 12 units.					
Mistral-7B Generated Response : The length is the distance between the endpoints. To find this distance, use the following					
formula: $(final_x - init_x)^2 + (final_y - init_y)^2 = distance^2$					
In this case, the values are $(10, 4)$ and $(2, -2)$. The length is $sqrt(8^2 + 6^2) = sqrt(100) = 10$.					
GPT4 Generated Response: To find the length of a line segment given its endpoints, you can use the distance formula. The					
distance formula between two points (x_1, y_1) and (x_2, y_2) in the coordinate plane is: $d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$. For the					
points given, $(x_1, y_1) = (2, -2)$ and $(x_2, y_2) = (10, 4)$. Plugging these into the formula: $d = \sqrt{(10-2)^2 + (4-(-2))^2}$ $d = 10, 4, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5,$					
$\sqrt{8^2+6^2}$ $d = \sqrt{64+36}$ $d = \sqrt{100}$ $d = 10$. Therefore, the length of the line segment is 10 units.					

4 EXPERIMENTS

We experimentally analyze model performance under two different settings: task-agnostic setting and task-aware setting. We evaluate our method on tasks that require long-form generations and use evaluation metrics that are commonly used to evaluate the corresponding tasks. In this section, we describe the evaluation tasks (i.e., the task that we wish to improve on via instruction tuning) in Sec. 4.1, experimental settings in Sec. 4.2, experimental details in Sec. 4.3, main results in Sec. 4.4, and additional analysis in Sec. 4.5.

4.1 EVALUATION TASKS

We use four evaluation tasks, including AlpacaEval (Dubois et al., 2024), TLDR (Stiennon et al., 2020), RLHF (Bai et al., 2022), and HumanEval (Chen et al., 2021).

Dataset	$ D_V $	$ D_{\text{test}} $	Evaluation metrics
AlpacaEval	10	795	length-controlled win rate
TLDR	322	6553	reward model
RLHF	2192	2354	reward model
HumanEval	10	154	pass@k

Table 2: Detailed information about the evaluation datasets.

AlpacaEval is a compilation of prompt-response pairs aimed at assessing language models' instruction following capability. We use the length-controlled win rate to ensure a fair evaluation (Dubois et al., 2024). HumanEval evaluates code generation from natural language instructions using the pass@k metric, which is the probability of having at least one correct solution (pass a specified unit test) when sampling k responses from the model. Pass@k with k > 1 is practical in real-world scenarios when expected behaviors are known and test cases are readily available. It evaluates the test-time scaling capability of the model by allowing multiple candidate solutions (k) to be generated and validated to effectively identify the correct code. We adopt pass@100 because it achieves substantially higher accuracy than smaller k (Chen et al., 2021), making it more practical for real-world scenarios. For completeness, we also provide detailed results for different k in App. H.1. **TLDR** contains polished text summaries. The evaluation metric is the reward model (OpenAssistant, 2023) trained on human feedback, measuring the the quality of summaries and alignment with human preference. **RLHF** consists of prompt-response pairs where each includes a "chosen" response that aligns better with human preferences (we use only the "chosen" columns). We use a trained reward model (Ray2333, 2024) as the evaluation metric to evaluate the helpfulness of the responses. For the dataset splits, we randomly select 10 examples from both AlpacaEval and HumanEval as the validation set, with the remainder as the test set. For RLHF, we sample 5% from the training dataset as the validation set since the original dataset only contains train and test splits. For TLDR, we use 5% of the original validation set, given that the original dataset size is too large. We study the robustness of our method to different validation splits in App. H.2. We provide a summary of these evaluation tasks in Tab. 2 and additional details on these tasks in App. B.1. Unless specified, the results reported are evaluated on the test set.

4.2 EXPERIMENTAL SETTINGS

Our problem setup focuses on targeted data selection, assuming access to validation data during data selection. We further consider two distinct settings, namely "task-agnostic" and "task-aware" settings, categorized by whether the knowledge of the downstream task is available when forming the initial pool of training data (for selection) in the data preparation stage. Specifically:

Task-agnostic Setting. In this setting, a large, diverse, mixed-source pool of instruction tuning training set is collected without the knowledge of the downstream task, before data selection. We use Tulu (Wang et al., 2023) as the training dataset, which consists of Chain of Thought (COT) (Wei et al., 2022), Databricks Dolly (DOLLY) (Conover et al., 2023), Open Assistant 1 (OASST) (Köpf et al., 2023), and FLAN V2 (Longpre et al., 2023). Intuitively, this mixed-source pool of data may contain irrelevant data (e.g., assistant-style conversations) w.r.t. the targeted task (e.g., coding task). Additional details for the training set are in App. A.1.

Task-aware Setting. In this setting, the training set is collected specifically for the downstream evaluation task (hence task-aware). We consider two evaluation tasks here: **RLHF** and **HumanEval**, with the evaluation datasets the same as in Sec. 4.1. For RLHF, the training data is 95% of the original training set used for the helpfulness assistant, as provided in Bai et al. (2022). For HumanEval, we adopt the CodeAlpaca 20k (Chaudhary, 2023) dataset as the training set, which is a crowd-sourced collection of code-related instruction-response pairs, designed to fine-tune language models for better performance in code generation and understanding. Additional training set details are in App. A.1.

4.3 EXPERIMENTAL DETAILS

Efficient Data Selection for LLM. To improve the efficiency of data selection, we train models with LoRA (Low-Rank Adaptation). We adopt the warmup training which trains the LLM on a randomly selected subset of training data for influence estimation and the number of warmup epochs is the *E*. Additionally, random projections are applied to the LoRA gradients, preserving the essential inner products while reducing the dimensionality of the gradient to reduce the memory requirement (Johnson, 1984). Detailed time-complexity analysis is elaborated in App. H.3.

Models and Hyperparameters. Our primary evaluation of NICE focuses on two LLMs: Llama2-7B (Touvron et al., 2023), and Mistral-7B (Jiang et al., 2023), with performance averaged over three seeds. Larger and state-of-the-art models, including Llama2-13B and Llama3-8B, were also tested on the RLHF dataset, presented in App. H.4. We perform warmup training for 4 epochs on 5% of the training set for the task-agnostic setting, and 20% for the task-aware setting due to its smaller training set size. We project the LoRA gradient into an 8192-dimensional vector. The influence estimates for data points are obtained by the respective data selection approaches, with the top 5% data points (ranked by influence) to be the selected data subset D_S for the task-agnostic setting (20% for the task-aware setting). Multinomial sampling (Chatterjee & Cancedda, 2010) is used to generate the MC samples for NICE. We generate 20 MC samples for all evaluation tasks, except for HumanEval, where we generate 500 samples due to the difficulty of the task (i.e., responses having low code pass rates). For NICEAMC, we use GPT-4 as $g(\cdot, \psi)$. Note that when the reward function does not require labels, the ground truth label is not used by our approach, detailed in App. C.

Baselines. We evaluate NICE against a variety of baselines. The most straightforward baseline is Random, where data points are randomly sampled from the training set for instruction tuning. We also employ BM25 (Robertson et al., 2009), which ranks training data based on relevance to the validation data, and then selects the top-ranked data points to form D_S . Another baseline, DSIR (Xie et al., 2023), selects D_S based on n-gram lexical feature matching between training and validation distributions. Representation-based Data Selection (RDS) (Zhang et al., 2018) ranks the training instances using the cosine similarity of the features between the training and the validation data, and we adopt SentenceBERT (Reimers, 2019) embedding as the features. LESS (Xia et al., 2024) uses loss-based influence estimation to select the training data with top influence scores. TSDS (Liu et al., 2024b) also leverages loss-based gradient features, and further optimizes for distribution alignment and diversity via optimal transport and kernel density estimation, respectively. Note that all baselines, except random, are calculated using complete data points (i.e., the concatenation of prompt and response). More implementation details about the baselines are in App. D.

Table 3: Comparison of NICE and NICEAMC in both task-agnostic and task-aware settings for
Llama2-7B and Mistral-7B. Bold numbers indicate the top-performing selected subset. A purple
cell suggests that NICE outperforms LESS which uses loss-based influence estimation. Underlined
numbers show that the subset selected by our approach exceeds the performance of the full dataset.
Subscript numbers represent standard deviations.

Model & Dataset		Full	Random	RDS	BM25	DSIR	TSDS	LESS	NICE	NICEAMC
Task-agnosti	с									
Llama2-7B	AlpacaEval TLDR RLHF HumanEval	22.59 2.40 2.31 47.44	$\begin{array}{c} 16.13_{\pm 1.18} \\ 1.80_{\pm 0.08} \\ 2.05_{\pm 0.11} \\ 44.30_{\pm 2.36} \end{array}$	14.70 2.08 1.87 45.29	19.60 2.15 2.83 46.19	20.27 1.53 2.57 42.22	$\begin{array}{c} 17.40_{\pm 2.44} \\ 2.19_{\pm 0.29} \\ 1.01_{\pm 0.12} \\ 43.68_{\pm 1.82} \end{array}$	$\begin{array}{c} 26.94_{\pm 2.37} \\ 3.37_{\pm 0.29} \\ 1.44_{\pm 0.07} \\ 47.50_{\pm 1.57} \end{array}$	$\frac{\underline{27.61}_{\pm 2.13}}{\underline{3.61}_{\pm 0.78}}$ $\frac{\underline{2.82}_{\pm 0.10}}{\underline{48.59}_{\pm 2.08}}$	$\frac{\underline{30.45}_{\pm 2.40}}{\underline{3.55}_{\pm 0.40}}\\ \underline{\underline{3.03}}_{\pm 0.02}\\ \underline{45.10}_{\pm 2.84}$
Mistral-7B	AlpacaEval TLDR RLHF HumanEval	33.77 2.79 2.56 83.63	$\begin{array}{c} 24.99_{\pm 4.28} \\ 3.06_{\pm 0.24} \\ 2.13_{\pm 0.04} \\ 85.56_{\pm 1.27} \end{array}$	21.70 2.90 1.78 84.15	28.47 2.41 2.88 84.09	29.31 3.48 2.94 79.17	$\begin{array}{c} 35.84_{\pm 0.53} \\ 3.28_{\pm 0.41} \\ 1.83_{\pm 0.15} \\ 82.78_{\pm 1.25} \end{array}$	$\begin{array}{c} 41.09_{\pm 1.56} \\ 4.40_{\pm 0.12} \\ 1.70_{\pm 0.09} \\ 85.24_{\pm 0.45} \end{array}$	$\begin{array}{c} \underline{41.43}_{\pm 3.00}\\ \underline{4.80}_{\pm 0.12}\\ \underline{3.10}_{\pm 0.06}\\ \underline{85.59}_{\pm 1.41}\end{array}$	$\begin{array}{c} \underline{47.40}_{\pm 2.94}\\ \underline{4.59}_{\pm 0.20}\\ \underline{3.42}_{\pm 0.05}\\ \underline{85.67}_{\pm 0.34} \end{array}$
Task-aware										
Llama2-7B	RLHF HumanEval	1.01 51.27	${\begin{array}{c} 1.04_{\pm 0.04}\\ 51.91_{\pm 1.61}\end{array}}$	0.66 54.74	1.29 52.23	1.43 53.10	$\begin{array}{c} 0.97_{\pm 0.02} \\ 49.85_{\pm 3.17} \end{array}$	${\begin{array}{c} 1.62 _{\pm 0.05} \\ 52.67 _{\pm 0.71} \end{array}}$	$\tfrac{1.69}{55.09}_{\pm 1.66}$	$\tfrac{1.32_{\pm 0.05}}{50.67_{\pm 1.24}}$
Mistral-7B	RLHF HumanEval	0.99 84.27	${\begin{array}{c} 1.05_{\pm 0.04}\\ 83.34_{\pm 2.54}\end{array}}$	0.56 86.75	1.31 84.81	1.31 79.91	${\begin{array}{c} 1.15_{\pm 0.06}\\ 85.51_{\pm 1.28}\end{array}}$	$\frac{1.29_{\pm 0.13}}{85.26_{\pm 1.13}}$	$\frac{\underline{1.71}}{\underline{87.35}_{\pm 1.03}}$	$\tfrac{1.35}{84.18_{\pm 1.63}}$

4.4 RESULTS

The results of NICE and NICEAMC are presented in Tab. 3. We summarize the findings below.

NICE Outperforms the Loss-based Influence Estimation. Our results show that NICE consistently outperforms LESS, which uses loss-based influence estimation for data selection, across various models, settings, and datasets. This result verifies that using the estimated influence on the evaluation metric is more helpful for data selection than that on the loss.

No Labels? No Problem! NICE Outperform Baselines that Use the Label Response. For tasks like TLDR, RLHF, and HumanEval, NICE or NICEAMC uses only unlabeled validation data (i.e., only prompts). Surprisingly, they outperform baselines that use labeled data (both prompts and label responses).

Less Is More: Subset Outperforms the Full Dataset. We find that the subset selected by NICE or NICEAMC can outperform the full dataset, demonstrating the value of carefully curated data over larger, less refined datasets.

Assisted Monte-Carlo Sampling Can Boost Data Selection. Adopting AMC sampling in policy gradient has the potential to further improve the performance of NICE, especially when the initial pool of training data is large. Responses generated from better-performing models can effectively guide data selection, enabling models trained on the selected subset to achieve better performance.

However, NICEAMC does not always have improved performance compared to NICE, particularly in the task-aware setting with a smaller selection pool (of training data). When the training set is small, there may not be enough training data points with gradients close to the policy gradient of NICEAMC. We perform a simple experiment to verify this intuition: restricting NICEAMC to compute score from the RLHF training set yields a performance of 1.26, but expanding the selection pool to the combination of RLHF set and a large instruction tuning set (COT, DOLLY, OASST, Flan V2) increases the performance to 3.35. Note that the additional controlled experiment only expanded the selection pool without altering the initial warmup process. In contrast, the performance of NICE improves from 1.68 to 2.44 with a larger pool. This comparison shows a clear advantage of NICEAMC when the size of the training data to select from is large. Consequently, for the task-agnostic setting when the size of training data is large, we can prioritize using NICEAMC.

4.5 Additional Analysis

Unless specified, the experiments in the section below are conducted on the Llama2-7B model.

Generalizing NICE to the Influence Function (IF). We demonstrate the effectiveness of the main idea of NICE beyond the TracIn framework by extending it to IF. We adopt DataInf as an efficient

implementation for IF. We compute the DataInf on the last checkpoint of the warmup model and adopt random projection to reduce the dimensionality of gradients and store the gradients, eliminating the need for computing the gradients again after the computation of the Hessian. More details on the implementation of DataInf are provided in the App. E.3. We compared the performance of selecting data with the vanilla (loss-based) DataInf and the DataInf enhanced by the policy gradient: NICEIF, in the task-agnostic setting (Sec. 4.2). As shown in Tab. 4, NICEIF generally outperforms DataInf, demonstrating the effectiveness of our approach and showing that the concept of NICE of applying policy gradient for influence estimation can be readily applied to other loss-based influence estimation methods to improve the performance.

Table 4: Comparison between Influence Function (DataInf) and NICEIF on Llama2-7B in *task-agnostic* setting. NICEIF consistently outperforms the DataInf.

Method	AlpacaEval	TLDR	RLHF	HumanEval
IF (DataInf)	11.11	2.01	0.83	37.40
NICEIF	20 .44	3.97	1.89	39.68

The Discrepancy Between NTP Loss and Evaluation Metrics. Previous works have discussed the discrepancy between validation loss and downstream performance in instruction tuning for LLMs (Tay et al., 2021; Xia et al., 2024). In these scenarios, minimizing validation loss does not necessarily correspond to improving validation performance, especially when the task requires long-form generations. Empirical observations described in Fig. 1 and further results in App. H.5 verify this discrepancy: The minimized validation loss is achieved at step around 250. However, that checkpoint is the worst-performing checkpoint (lowest validation reward). The reward can be further increased in later steps, even if the loss increases. We additionally provide a comparison between NICE and LESS in App. H.6, which demonstrates that NICE-selected data optimizes towards the direction of increasing validation performance, while LESS-selected data may not consistently improve validation performance.



Figure 1: Discrepancy between NTP loss and performance (i.e., measured by the reward model here) of the validation set for RLHF task in the last few training steps. The checkpoint with the minimized loss (highest negative loss) corresponds to the checkpoint with a relatively worse performance (having the lowest reward among the five checkpoints). The performance may continue to increase even if the loss increases (negative loss decreases).

The Effect of the Number of Monte-Carlo Samples. We empirically study the effect of the number of Monte-Carlo (MC) samples used in approximating policy gradient on the data selection performance (measured by reward for the RLHF dataset). The results in Fig. 2 indicate a positive correlation between performance and the number of MC samples, which shows the potential to further improve NICE by using more MC samples. However, generating additional MC samples is computationally expensive. We use 20 MC samples for the majority of tasks since it is relatively less computationally expensive, while sufficient to achieve better performance than other baselines. The number of MC samples can also affect the stability of our approach, elaborated in App. H.7.

Data Addition. In Fig. 3, we plot the performance (measured by reward) against the percent of data points selected by NICE and Random for RLHF dataset. The performance of the model trained on a randomly selected subset increases gradually as more data is used. In contrast, for the NICE-selected subset in the task-aware setting (right of Fig. 3), performance rises slightly from 5% to 25% but drops sharply beyond 25%. The task-aware training data is more relevant to the downstream task, so a small percentage can miss useful data, while a large percentage can include irrelevant or harmful data, harming performance. In the task-agnostic setting, performance declines as the selection percentage



Figure 2: Performance of LLMs trained on data selected by NICE for RHLF dataset when different numbers of MC samples are used. There is a positive correlation between performance and generated MC samples. Using the sampling size of 20 provides good performance while increasing the sampling size has the potential to improve the performance.

increases, likely because only a small fraction is relevant. This experiment also demonstrates the importance of data selection to exclude data points that are not useful to model performance.



Figure 3: Performance versus percent of data points selected by NICE or Random for RLHF dataset. As the percentage selected by NICE increases, performance may drop due to the inclusion of low-score, irrelevant, or harmful points. When using data selection, the trained model consistently outperforms those trained on randomly selected data, even outperforming the full dataset, underscoring the importance of data selection.

5 RELATED WORKS

Various approaches are proposed to estimate the influence of training data in fine-tuning LLMs. LESS (Xia et al., 2024) adapts the TracIn framework to estimate the influence of data points in instruction tuning. Kwon et al. (2024) and Choe et al. (2024) scale up the IF by speeding up the computation of the Hessian inverse. Lin et al. (2024) studies the token-level influence function for LLMs. These approaches above are all loss-based influence estimation methods that aim to approximate the influence of data on the validation loss. Consequently, they fall short for generations tasks whose evaluation metrics align poorly with NTP loss, as elaborated in Sec. 4.5. On the other hand, the works of Kwon et al. (2024); Choe et al. (2024); Lin et al. (2024) focus on data attribution and hence do not optimize for data selection. Other works, including Park et al. (2023) which approximates the data model to estimate the influence, and Wang et al. (2024a) which scales up the Shapley value to estimate the influence only consider classification tasks, and hence can not be used in instruction tuning. Our work focuses on the influence estimation for tasks that require extensive generations (i.e., generating long responses) instead of only classification. More related works on data curation for instruction tuning are discussed in App. G.

6 CONCLUSION

We propose NICE, a novel influence estimation approach that selects training data to directly optimize non-differentiable evaluation metrics via policy gradient, rather than relying on NTP loss. When using a reward function that does not require label response, NICE can perform data selection without relying on costly annotated labels. Experimental results show that our approach outperforms existing data selection methods across diverse scenarios. Of note, despite the superior performance achieved by NICE and NICEAMC, the computational cost of these approaches is not negligible even with the acceleration (e.g., the use of LoRA and random projection). Further explorations can be done to study other computationally efficient ways of computing the gradients.

ACKNOWLEDGEMENT

Jingtan Wang is supported by the Institute for Infocomm Research of Agency for Science, Technology and Research (A*STAR). This research is supported by the National Research Foundation (NRF), Prime Minister's Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) programme. The Mens, Manus, and Machina (M3S) is an interdisciplinary research group (IRG) of the Singapore MIT Alliance for Research and Technology (SMART) centre. The computational work for this article was partially performed on resources of the National Supercomputing Centre, Singapore (https://www.nscc.sg).

REFERENCES

- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, Colin Raffel, Shiyu Chang, Tatsunori Hashimoto, and William Yang Wang. A survey on data selection for language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

MS Bartlett. Approximate confidence intervals. *Biometrika*, 40(1/2):12–19, 1953.

- Gantavya Bhatt, Yifang Chen, Arnav Das, Jifan Zhang, Sang Truong, Stephen Mussmann, Yinglun Zhu, Jeff Bilmes, Simon Du, Kevin Jamieson, et al. An experimental design framework for label-efficient supervised finetuning of large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pp. 6549–6560, 2024.
- 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. In *Advances in Neural Information Processing Systems*, 2020.
- Alexander Bukharin and Tuo Zhao. Data diversity matters for robust instruction tuning. *arXiv* preprint arXiv:2311.14736, 2023.
- 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.
- Nicholas Carlini, Matthew Jagielski, Christopher A Choquette-Choo, Daniel Paleka, Will Pearce, Hyrum Anderson, Andreas Terzis, Kurt Thomas, and Florian Tramèr. Poisoning web-scale training datasets is practical. In 2024 IEEE Symposium on Security and Privacy (SP), pp. 407–425. IEEE, 2024.
- Samidh Chatterjee and Nicola Cancedda. Minimum error rate training by sampling the translation lattice. In Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pp. 606–615, Cambridge, MA, October 2010. Association for Computational Linguistics.
- Sahil Chaudhary. Code alpaca: An instruction-following llama model for code generation. *GitHub* repository, 2023.
- Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpagasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*, 2023.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Sang Keun Choe, Hwijeen Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya Pratapa, Willie Neiswanger, Emma Strubell, Teruko Mitamura, et al. What is your data worth to GPT? LLM-scale data valuation with influence functions. arXiv preprint arXiv:2405.13954, 2024.

- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world's first truly open instruction-tuned llm. *Company Blog of Databricks*, 2023.
- Qianlong Du, Chengqing Zong, and Jiajun Zhang. Mods: Model-oriented data selection for instruction tuning. arXiv preprint arXiv:2311.15653, 2023.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. Alpacafarm: A simulation framework for methods that learn from human feedback. In *Advances in Neural Information Processing Systems*, pp. 30039–30069, 2023.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024.
- Benoît Frénay and Michel Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2013.
- Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Roziere, David Lopez-Paz, and Gabriel Synnaeve. Better and faster large language models via multi-token prediction. In *International Conference on Machine Learning*, 2024.
- Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, et al. Studying large language model generalization with influence functions. *arXiv preprint arXiv:2308.03296*, 2023.
- 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. *arXiv preprint arXiv:2106.09685*, 2021.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. arXiv preprint arXiv:2310.06825, 2023.
- William B Johnson. Extensions of lipshitz mapping into hilbert space. In *Conference modern analysis* and probability, 1984, pp. 189–206, 1984.
- Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *International Conference on Machine Learning*, pp. 1885–1894. PMLR, 2017.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations - democratizing large language model alignment. In Advances in Neural Information Processing Systems, pp. 47669–47681, 2023.
- Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou. DataInf: Efficiently estimating data influence in lora-tuned LLMs and diffusion models. In *Proc. ICLR*, 2024.
- 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. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, June 2024a.
- Yunshui Li, Binyuan Hui, Xiaobo Xia, Jiaxi Yang, Min Yang, Lei Zhang, Shuzheng Si, Ling-Hao Chen, Junhao Liu, Tongliang Liu, Fei Huang, and Yongbin Li. One-shot learning as instruction data prospector for large language models. In *Proceedings of the 62nd Annual Meeting of the Association* for Computational Linguistics (Volume 1: Long Papers), pp. 4586–4601, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.252.

- Huawei Lin, Jikai Long, Zhaozhuo Xu, and Weijie Zhao. Token-wise influential training data retrieval for large language models. *arXiv preprint arXiv:2405.11724*, 2024.
- Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. In *The Twelfth International Conference on Learning Representations*, 2024a.
- Zifan Liu, Amin Karbasi, and Theodoros Rekatsinas. TSDS: Data selection for task-specific model finetuning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b.
- Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pp. 22631–22648. PMLR, 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. In *The Twelfth International Conference on Learning Representations*, 2024.
- James Martens et al. Deep learning via hessian-free optimization. In *International Conference on Machine Learning*, volume 27, pp. 735–742, 2010.
- David Meyer. Notes on policy gradients and the log derivative trick for reinforcement learning, 2023. Lecture Notes.
- OpenAssistant. Reward model deberta v3 large v2, 2023. URL https://huggingface.co/ OpenAssistant/reward-model-deberta-v3-large-v2. Accessed: 2025-01-06.
- 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. In *Proc. NeurIPS*, pp. 27730–27744, 2022.
- Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. TRAK: Attributing model behavior at scale. In *International Conference on Machine Learning*, 2023.
- Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating training data influence by tracing gradient descent. In *Advances in Neural Information Processing Systems*, volume 33, pp. 19920–19930. Curran Associates, Inc., 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. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Ray2333. Gpt-2 large harmless reward model. https://huggingface.co/Ray2333/gpt2-large-harmless-reward_model, 2024. Accessed: 2024-06-18.
- N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019.
- Stephen Robertson, Hugo Zaragoza, et al. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389, 2009.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, pp. 3008–3021, 2020.

Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.

Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from pre-training and fine-tuning transformers. *arXiv preprint arXiv:2109.10686*, 2021.

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Jingtan Wang, Xiaoqiang Lin, Qiao Rui, Foo Chuan-Sheng, and Low Bryan Kian Hsiang. Helpful or harmful data? fine-tuning-free shapley attribution for explaining language model predictions. In *International Conference on Machine Learning*, 2024a.
- Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. Diversity measurement and subset selection for instruction tuning datasets. *arXiv preprint arXiv:2402.02318*, 2024b.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. How far can camels go? exploring the state of instruction tuning on open resources. *Advances in Neural Information Processing Systems*, 36:74764–74786, 2023.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pp. 24824–24837, 2022.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992.
- Lijun Wu, Fei Tian, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. A study of reinforcement learning for neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium, October-November 2018. Association for Computational Linguistics.
- Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less: Selecting influential data for targeted instruction tuning. In *International Conference on Machine Learning*, 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.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- Chih-Kuan Yeh, Ankur Taly, Mukund Sundararajan, Frederick Liu, and Pradeep Ravikumar. First is better than last for language data influence. *Advances in Neural Information Processing Systems*, 35:32285–32298, 2022.
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595, 2018.
- Yingxiu Zhao, Bowen Yu, Binyuan Hui, Haiyang Yu, Minghao Li, Fei Huang, Nevin L. Zhang, and Yongbin Li. Tree-instruct: A preliminary study of the intrinsic relationship between complexity and alignment. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 16776–16789, Torino, Italia, May 2024. ELRA and ICCL.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proc. NeurIPS*, volume 36, 2023.
- 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, 2024.

A TRAINING

A.1 TRAINING DATASET

For the task-agnostic setting, four processed training datasets that are utilized are described in Wang et al. (2023). These datasets, annotated or authored by humans, are detailed in Tab. 5. FLAN V2 and COT are based on existing NLP datasets, while DOLLY and OASST feature open-ended responses from humans, demonstrating diverse formats, lengths, and tasks.

For the task-aware setting, we use 95% of the *RLHF* task's training data as our training set, while the remaining 5% serves as the validation set. The original RLHF data point is a pair of responses, with one marked as 'chosen' by human annotators for being more helpful than the other. We only use these 'chosen' responses for training, validation, and test. During human labeling, only the last response from the assistant is compared; hence the last response is considered more helpful. Consequently, we use only the last-turn response as labels. For the *HumanEval* task's training data, we adopt the CodeAlpaca 20k (Chaudhary, 2023) dataset with the addition of the original instruction format for each data point.

For fine-tuning Llama2, we adopt the 'Tulu' format following from the study by Wang et al. (2023).

Tulu Instruction Format

```
<user>
What can you help me with?
<assistant>
I'd like to show off how chat templating works!
```

For fine-tuning the Mistral and Llama3 models, we utilize the respective instruction formats predefined by each model, as detailed below:

Mistral Instruction Format

[INST] What can you help me with? [/INST] I'd like to show off how chat templating works!</s>

Llama3 Instruction Format

```
<|start_header_id|>user<|end_header_id|>
```

```
What can you help me with?<|eot_id|><|start_header_id|>assistant <|end_header_id|>
```

I'd like to show off how chat templating works!

A.2 TRAINING DETAILS

All experiments utilized the parameter-efficient LoRA (Hu et al., 2021) approach. A linear warm-up with 0.03 warmup ratio was employed, peaking at a learning rate of 2×10^{-5} . We trained for 4 epochs on each dataset with a batch size of 128. The LoRA module had a rank of 128, an α of 512, a dropout rate of 0.1, and learned matrices for all attention layers. Specifically, the Llama2-7b model used is meta-llama/Llama-2-7b-hf, Llama2-13b used is meta-llama/Llama-2-13b-hf, Mistral-7B is mistralai/Mistral-7B-v0.1, Llama3-8B is meta-llama/Meta-Llama-3-8B.

Each main experiment was repeated three times with random seeds 0, 1, and 2. Under random selection methods, three different random subsets from the training set were chosen for each seed. For LESS, TSDS, NICE, and NICEAMC, we first performed warmup training on subsets chosen by each seed, then selected distinct subsets from the resulting warmup model for each trial. We used the same

	Task-agnostic								
Dataset	Size	Sourced from	# Turns	Prompt Len.	Response Len.				
FLAN V2	100,000	Based on Existing NLP	1	355.7	31.2				
		Datasets							
COT	100,000	Based on Existing NLP	1	266.0	53.2				
		Datasets							
DOLLY	15,011	Human-written from	1	118.1	91.3				
		scratch							
OASST	55,668	Human-written from	1.6	34.8	212.5				
		scratch							
		Task-awa	re						
RLHF	41,643	Human feedback data	2.46	145.0	517.3				
Code Alpaca	20,022	Code-related human-	1	294.8	197.0				
		written instructions							

Table 5: Detailed information about the training set. The task-agnostic training set is the same as in Xia et al. (2024) and Wang et al. (2023).

optimization seeds as those used for the warmup model. For experiments without reported standard deviation, we used seed 0. Full training was conducted on seed 0 only, due to heavy computation.

B EVALUATION

B.1 EVALUATION DATASET

AlpacaEval sources its data from self-instruct, OASST, Anthropic's helpful dataset, Vicuna, and Koala, widely used for understanding model behavior in structured, instruction-driven settings. The evaluation metric is a length-controlled win rate, adjusting for biases in response length using a regression model to ensure fair and accurate assessments (Dubois et al., 2024). Each response is compared to a baseline model, *text_davinci_003*, using *weighted_alpaca_eval_gpt4_turbo* as annotator. We use *text_davinci_003* as a baseline because its relatively lower baseline performance can more clearly highlight the performance difference between selected subsets. By contrast, using a stronger baseline (e.g., *gpt4_turbo*) could mask the differences among the model trained on the subset selected by different selection strategies. For annotation, we employ *weighted_alpaca_eval_gpt4_turbo*, chosen for its high agreement with human annotations, large context capacity, and cost-effectiveness.

HumanEval uses a set of programming challenges to evaluate code correctness and functionality, measured by the pass@k metric which is how many correct solutions appear within a specified number of attempts k (e.g., we use 100 in the main experiment). This metric assesses the ability of models to generate correct solutions within a limited number of attempts, reflecting the model's efficiency in code generation. We additionally use the unbiased estimator of pass@k to avoid the high variance of vanilla pass@k (Chen et al., 2021).

TLDR uses cleaned data from Stiennon et al. (2020), focusing on well-structured input-output pairs for summarization tasks. It is evaluated against a reward model (OpenAssistant, 2023) that is trained based on human feedback to ensure the generation of high-quality summaries.

RLHF is designed for training and evaluating language models using human feedback to optimize response generation. It consists of prompt-response pairs and we only use the "chosen" response whose last-turn response aligns best with human preferences, according to specific criteria such as relevance and safety. The evaluation metric is reward model (Ray2333, 2024) which can measure the helpfulness of model responses.

B.2 EVALUATION DETAILS

For AlpacaEval, two samples are drawn from each of the five subtasks (self-instruct, OASST, Anthropic's helpful dataset, Vicuna, Koala), resulting in a total of ten samples. For HumanEval, since pass@k aims to check if there is at least one functional code within the k generations of each test data point, we select training data points that enhance performance across all validation data. Each validation data point is treated as an individual task. Hence, AlpacaEval and HumanEval are

multi-subtask scenarios with q = 5 and q = 10, while TLDR and RLHF are single-subtask scenarios with q = 1. The aggregation of each training data point's influence score w.r.t validation set follows from Sec. 2.2.

For TLDR and RLHF, the reward function r used during the computation of the policy gradient is the reward model, and the evaluation metric for the validation set is the average reward for each validation point. For AlpacaEval, the reward function r for each generated response is the annotator's decision for that response, and the evaluation metric for the validation set is not simply an average but an average of a debiased version of each annotator's decision (Dubois et al., 2024). For HumanEval, the reward function r for each generated response (code) is the boolean result of all unit tests for that code, while the evaluation metric for the validation set is pass@k. This metric measures functional correctness by generating k codes for each test problem and considering the problem solved if any code passes all unit tests of that problem. The pass@1 score can be viewed as the average of each boolean result, whereas pass@10 and pass@100 are computed via a problem-level "OR" across the k generated codes, followed by an average across all problems. Additionally, to address numerical instability and reduce variance, we use the unbiased estimator version of pass@k following Chen et al. (2021).

We evaluate the validation set D_V (the same set used for data selection) at the end of each epoch, and the best-performing checkpoint (measured by the validation performance) is assessed on the test set.

C NICE AND NICEAMC DETAILS

NICE and NICEAMC rely solely on the probability of generated responses and the score from a reward function, as enabled by the policy gradient mechanism. For AlpacaEval, labels are provided to all approaches because the reward function requires ground-truth labels. For the other three tasks, including TLDR, RLHF, and HumanEval, the ground-truth labels of the validation dataset are not used by our method, although they are available to the other baselines (such as LESS). This is because in these three tasks, the reward function does not require ground-truth labels. Specifically:

- For RLHF and TLDR, the reward function is a learned reward model that outputs scores based on the prompt and generated response.
- For HumanEval, the reward is defined by whether the generated code passes unit tests, not requiring reference solutions.

C.1 NICE DETAILS

We generate 20 MC samples for all evaluation datasets, except for HumanEval, where we generate 500 samples due to the task's difficulty (i.e., a lower code pass rate). Generally, we set the sampling temperature to 1.2 to promote diversity, except for AlpacaEval, where we use 1.0. The generated responses of the validation set under temperature of 1.0 yield a higher win rate on the final checkpoint compared to 1.2. We employ multinomial sampling with top_k = 50 and top_p = 0.95.

C.2 NICEAMC DETAILS

For NICEAMC, we use gpt-4-turbo-2024-04-09. Regarding the GPT API hyperparameters, we set the frequency_penalty to 0, presence_penalty to 0, and temperature to 0.8 for all tasks.

For TLDR tasks, we add a prompt A brief summary of my post is (TL;DR): after the prompt of the data point, before generating the response, to enhance generation quality.

For HumanEval, we prepend a prompt:

```
Complete the following python function to return only the function
body (completion).
Do not include the function header or docstring.
```

before the coding question. This ensures the model outputs only the necessary code, avoiding chain-of-thought content that could fail unit tests.

D BASELINE DETAILS

For LESS and TSDS, the hyperparameter settings are the same as in their official repos. For BM25, DSIR, and RDS, these methods are warmup model-agnostic, meaning the selection process does not rely on the warmup models' randomness. Hence, we only run them on seed 0. To avoid the instruction format's effect on the representation or retrieval, we use the format of:

Question: [Question] \n\nAnswer: [Answer]

for the majority of the training data and validation data. We use TULU format for RLHF training and validation data, as RLHF contains many turns.

For BM25, we use the rank_bm25 package (https://github.com/dorianbrown/rank_bm25). We treat each validation data point as a query to retrieve the BM25 scores of each training data point. For RDS, we adopt the sentence-transformers/all-MiniLM-L6-v2 model (https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2). We compute the cosine similarity of embeddings between each validation data point and each training data point. Overall, for each validation point, BM25 and RDS will have a score vector of dimension n (the size of the training dataset is n). We then follow the same aggregation steps as in Sec. 2.2 to aggregate the scores. For DSIR, we use the official Github repo and match the distribution between training data and validation data (https://github.com/p-lambda/dsir).

E PRELIMINARY

E.1 DERIVATION OF THE CHANGE IN THE VALIDATION LOSS

. . .

The detailed derivation of the change in the validation loss is explained here: When a training data point z_i is included in the training step t, the model parameters are updated accordingly, leading to a change in the validation loss. Assuming a small learning rate η_t is used in the parameter updates with the Stochastic Gradient Descent (SGD) optimizer, this one-step change at step t can be approximated using a first-order Taylor expansion (Pruthi et al., 2020):

$$L(z_{v}; \theta^{t+1}) - L(z_{v}; \theta^{t})$$

$$= \nabla_{\theta} L(z_{v}; \theta^{t}) \cdot (\theta^{t+1} - \theta^{t}) + O(\|\theta^{t+1} - \theta^{t}\|^{2})$$

$$\approx \nabla_{\theta} L(z_{v}; \theta^{t}) \cdot (\theta^{t+1} - \theta^{t})$$

$$= \nabla_{\theta} L(z_{v}; \theta^{t}) \cdot (-\eta_{t} \nabla_{\theta} L(z_{i}; \theta^{t}))$$

$$= -\eta_{t} \langle \nabla_{\theta} L(z_{v}; \theta^{t}), \nabla_{\theta} L(z_{i}; \theta^{t}) \rangle.$$
(1)

E.2 DERIVATION OF INFLUENCE FUNCTION

Influence Function (IF) (Koh & Liang, 2017) measures the influence of down-weighting the training data point z_i by some small ϵ , on the new parameter $\theta_{\epsilon,z_i}^E := \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n L(z_i; \theta) + \epsilon L(z_i; \theta)$. The parameter change is given by:

$$\mathcal{I}_{\text{down,params}}(z_i) = \frac{\partial \theta_{\epsilon, z_i}^E}{\partial \epsilon} \bigg|_{\epsilon=0} = H_{\theta^E}^{-1} \nabla_{\theta} L(z_i; \theta^E) ,$$

where $H_{\theta^E} = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta}^2 L(z_i; \theta^E)$ is the Hessian matrix of the average loss over the training set and is positive definite by assumption. Then, by applying the chain rule, IF can measure the influence of down-weighting z_i on the loss of the validation data point z_v :

$$\begin{aligned}
\operatorname{Inf}_{\operatorname{IF}}(z_{i}, z_{v}) &= \frac{\partial L(z_{v}; \theta_{\epsilon, z_{i}}^{E})}{\partial \epsilon} \Big|_{\epsilon=0} \\
&= \nabla_{\theta} L(z_{v}; \theta^{E})^{\top} \frac{\partial \theta_{\epsilon, z_{i}}^{E}}{\partial \epsilon} \Big|_{\epsilon=0} \\
&= \nabla_{\theta} L(z_{v}; \theta^{E})^{\top} H_{\theta^{E}}^{-1} \nabla_{\theta} L(z_{i}; \theta^{E}) .
\end{aligned} \tag{2}$$

E.3 EFFICIENT INFLUENCE FUNCTION

DataInf makes use of Bartlett's second identity (Bartlett, 1953) to approximate the Hessian. That is, the Hessian can be replaced with the second moment of the first-order gradients: $G(\theta^E) = \frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta} L(z_i; \theta^E) \nabla_{\theta} L(z_i; \theta^E)^{\top}$, which simplifies the computation of the Hessian matrix. The θ^E here is the last checkpoint of the warmup model. In the later section, the computation of DataInf is all w.r.t. to the last checkpoint of the warmup model. We denote the gradient of the loss of training data point z_i w.r.t. ϕ (ϕ can be a single layer's parameter) by $\nabla_{\phi} L_i := \nabla_{\phi} L(z_i; \theta^E)$. To further address computational challenges, DataInf adopts another two techniques: 1. **Damping**: A small positive constant, λ , is added to the diagonal elements of $G(\theta^E)$, enhancing its positive definiteness and invertibility (Martens et al., 2010). 2. **Block Diagonal Matrix Representation**: $G(\theta^E)$ is further approximated using its block diagonal matrix, where each block is a layer of the neuron network (Grosse et al., 2023). The inverse of Hessian then becomes:

$$\left(\frac{1}{n}\sum_{i=1}^{n}\nabla_{\theta_{l}}L_{i}\nabla_{\theta_{l}}L_{i}^{\top}+\lambda I_{d_{l}}\right)^{-1},$$

where $\theta_l \in \mathbb{R}^{d_l}$ is the model parameter in *l*-the layer and $I_{d_l} \in \mathbb{R}^{d_l \times d_l}$ is identify matrix of size d_l . Following these transformations, DataInf simplifies the inverse of the average of the gradient outer products regularized by λI_{d_l} using the Sherman-Morrison formula:

$$\left(\frac{1}{n} \sum_{i=1}^{n} \nabla_{\theta_l} L_i \nabla_{\theta_l} L_i^{\top} + \lambda I_{d_l} \right)^{-1} \approx \frac{1}{n} \sum_{i=1}^{n} \left(\nabla_{\theta_l} L_i \nabla_{\theta_l} L_i^{\top} + \lambda I_{d_l} \right)^{-1}$$
$$= \frac{1}{n\lambda} \sum_{i=1}^{n} \left(I_{d_l} - \frac{\nabla_{\theta_l} L_i \nabla_{\theta_l} L_i^{\top}}{\lambda + \nabla_{\theta_l} L_i^{\top} \nabla_{\theta_l} L_i} \right) .$$

The computation of the Hessian and the later computation of the influence function can require retrieving the training gradient twice. To additionally optimize the time, we apply random projections to each gradient vector $\nabla_{\phi} L_i$ and store them. Later, we can retrieve the stored projected gradients and compute Hessian and the influence function score.

Besides *DataInf*, other methods have been proposed for efficiently computing the influence function for large language models. Grosse et al. (2023) improved the computation of the block-diagonal Hessian using the Kronecker product of uncentered forward and backward covariances of each layer (*EK-FAC*). *TRAK* projects gradients into a low-dimensional space and calculates influence scores within the subspace (Park et al., 2023). *LORGA* further improves the projection step with an efficient gradient projection strategy that leverages the gradient structure in backpropagation (Choe et al., 2024). However, neither *EK-FAC* nor *LORGA* discussed their applicability to *LoRA* fine-tuned models, leading us to exclude them for experiment. Additionally, *TRAK*'s approach of treating the multi-class classification problem as a single binary logistic regression may result in significant information loss. Therefore, we also did not apply it, as our task is purely generative and each token's prediction is a multi-class classification prediction whose prediction space is vocabulary size. For the hyperparameter in terms of Datainf, the projected gradient is 2048 dimension. The smaller dimension is due to the fact that we need to project gradients for each layer. The λ is 0.0001.

F EXPLICIT FORM OF NICE DURING IMPLEMENTATION

As discussed in Sec. 3.1, Xia et al. (2024) replaces the SGD with Adam gradient and replaces the inner product with the cosine similarity of the original TracIn for performance consideration. We

integrate these two enhancements on NICE as well, leading to:

$$\begin{aligned} \operatorname{Inf}_{\operatorname{NICE}}(z_{i}, z_{v}) &= \sum_{e=1}^{E} \eta_{\bar{e}} \frac{\left\langle \mathbb{E}_{\hat{y}_{v} \sim f(y|x_{v};\theta^{e})} \left[-\nabla_{\theta} \log(f(\hat{y}_{v}|x_{v};\theta^{e}))r(z_{v},\hat{y}_{v}) \right], \Gamma_{\theta}(z_{i};\theta^{e}) \right\rangle}{\left\| \mathbb{E}_{\hat{y}_{v} \sim f(y|x_{v};\theta^{e})} \left[-\nabla_{\theta} \log(f(\hat{y}_{v}|x_{v};\theta^{e}))r(z_{v},\hat{y}_{v}) \right] \right\| \left\| \Gamma_{\theta}(z_{i};\theta^{e}) \right\|} \\ \Gamma_{\theta}(z_{i};\theta^{e}) &\triangleq \frac{\mathbf{m}^{e+1}}{\sqrt{\mathbf{v}^{e+1} + \epsilon}} \\ \mathbf{m}^{e+1} &= \frac{\beta_{1}\mathbf{m}^{e} + (1-\beta_{1})\nabla_{\theta}L(z_{i};\theta^{e})}{1-\beta_{1}^{e}} \\ \mathbf{v}^{e+1} &= \frac{\beta_{2}\mathbf{v}^{e} + (1-\beta_{2})\left(\nabla_{\theta}L(z_{i};\theta^{e})\right)^{2}}{1-\beta_{2}^{e}} \,, \end{aligned}$$

where every operation is applied elementwise. Here, β_1 and β_2 represent the hyperparameters for the first and second moments, respectively, with ϵ serving as a small constant.

G RELATED WORK: DATA CURATION FOR INSTRUCTION TUNING

Curating high-quality, diverse, and complex instruction tuning data has been shown to improve the instruction-following ability of LLMs. Researchers have proposed different strategies to measure and improve different aspects of instruction-following ability: Cao et al. (2023) utilize natural language indicators to evaluate quality; Zhao et al. (2024) employ GPT-transformed instructions to measure complexity; and Chen et al. (2023) adopt LLM-annotated scores to assess both quality and complexity. Bukharin & Zhao (2023); Du et al. (2023); Lu et al. (2024); Wang et al. (2024b) optimize instruction data by emphasizing diversity in tandem with quality or complexity. Liu et al. (2024a) offers a comprehensive comparison of existing methods regarding these three properties. They further train a model to predict complexity and quality and iteratively filter out the most diverse points using embeddings. These approaches typically select data without a validation set. Our work extends this line of research to both task-agnostic and task-aware settings, selecting data that aligns the most with downstream tasks to enhance specific model capabilities. Our setting is more similar to Xia et al. (2024), Liu et al. (2024b) and Li et al. (2024b). However, Xia et al. (2024) selects data based on the influence on validation loss, which can lead to discrepancies between minimizing loss and maximizing performance, while we select data based on their influence on validation performance. Liu et al. (2024b) uses loss-based gradients to measure data distance when optimizing for distribution alignment and diversity, which can also suffer from the aforementioned discrepancy to some extent. Li et al. (2024b) utilizes a perplexity-based scoring system to select the most advantageous data for a defined anchor set, but their methodology is limited to single-turn training data. Another line of work, exemplified by Bhatt et al. (2024), frames curation as active learning by selecting the most informative prompts for predicting the label (i.e., generate responses) through uncertainty or diversity maximization. This active learning paradigm differs from our setting. A more comprehensive review of data curation methods is provided by Albalak et al. (2024).

H ADDITIONAL ANALYSIS

Unless otherwise specified, the experiments in the section below are conducted on the Llama2-7B model.

H.1 ADDITIONAL RESULTS OF PASS@k ON HUMANEVAL

We provide additional metrics, specifically pass@1 and pass@10, in Tab. 6. Overall, they align with our main findings: NICE and NICEAMC generally outperform loss-based influence estimation (LESS) and baselines, and the subsets selected by NICE and NICEAMC have the potential to outperform the full dataset. Note that BM25 performs well in terms of pass@1 and pass@10 on the Mistral-7B model and even outperforms other baselines in the task-agnostic setting with the Mistral-7B model. Empirically, we observe that BM25 is more likely to select data points that contain codes. This is likely because BM25 is based on TF-IDF, which assigns higher scores to training data that is more relevant in terms of word frequency, and certain words appear more frequently than

Table 6: Additional Results of Pass@k on HumanEval for both *task-agnostic* and *task-aware* settings on Llama2-7B and Mistral-7B. **Bold** numbers indicate the top-performing selected subset. A purple cell suggests that NICE outperforms LESS which uses loss-based influence estimation. <u>Underlined</u> numbers show that the subset selected by our approach exceeds the performance of the full dataset. Numbers in small font represent standard deviations.

Task-agnostic betting								
		Llama2-7B			Mistral-7B			
	Pass@1	Pass@10	Pass@100	Pass@1	Pass@10	Pass@100		
Full	7.61	25.52	47.44	29.47	59.68	83.63		
Random	$8.34_{\pm 0.34}$	23.85 ± 0.40	$44.30_{\pm 2.36}$	$29.99_{\pm 1.37}$	$62.04_{\pm 1.57}$	85.56 ± 1.27		
RDS	10.00	25.55	45.29	30.31	62.00	84.15		
BM25	8.27	24.51	46.19	31.81	62.43	84.09		
DSIR	9.53	24.02	42.22	27.71	56.81	79.17		
TSDS	$10.30_{\pm 1.58}$	$25.47_{\pm 1.27}$	$43.68_{\pm 1.82}$	$27.50_{\pm 1.47}$	$59.78_{\pm 1.86}$	$82.78_{\pm 1.25}$		
LESS	$9.24_{\pm 0.77}$	$26.12_{\pm 0.17}$	47.50 ± 1.57	$26.85_{\pm 0.58}$	$60.66_{\pm 0.39}$	85.24 ± 0.45		
NICE	$10.35_{\pm 1.72}$	$27.37_{\pm 1.56}$	$48.59_{\pm 2.08}$	$29.48_{\pm 0.93}$	$62.05_{\pm 2.23}$	$85.59_{\pm 1.41}$		
NICEAMC	$9.04_{\pm 2.35}$	$25.11_{\pm 1.72}$	$45.10_{\pm 2.84}$	$29.96_{\pm 1.95}$	$62.10_{\pm 1.82}$	$85.67_{\pm 0.34}$		

Task-agnostic Setting

Fask-aware Setting	
---------------------------	--

	8							
	Pass@1	Llama2-7B Pass@10	Pass@100	Pass@1	Mistral-7B Pass@10	Pass@100		
Full	13.27	30.30	51.27	33.14	64.09	84.27		
Random	11.99 ± 0.22	29.86 ± 0.42	$51.91_{\pm 1.61}$	33.15 ± 0.76	$63.62_{\pm 1.92}$	$83.34_{\pm 2.54}$		
RDS	12.40	31.34	54.74	33.32	63.23	86.75		
BM25	13.66	31.22	52.23	33.58	64.35	84.81		
DSIR	11.98	30.43	53.10	32.30	59.07	79.91		
TSDS	12.85 ± 0.50	28.15 ± 0.91	$49.85_{\pm 3.17}$	$31.74_{\pm 1.63}$	$63.04_{\pm 1.06}$	$85.51_{\pm 1.28}$		
LESS	13.55 ± 0.28	$30.53_{\pm 0.57}$	$52.67_{\pm 0.71}$	$34.05_{\pm 1.28}$	$64.12_{\pm 0.37}$	85.26 ± 1.13		
NICE	$13.43_{\pm 0.33}$	$31.70_{\pm 0.66}$	$55.09_{\pm 1.66}$	$33.61_{\pm 1.29}$	$65.56_{\pm 1.32}$	$87.35_{\pm 1.03}$		
NICEAMC	$12.87_{\pm 0.53}$	$30.39_{\pm 0.26}$	$50.67_{\pm 1.24}$	$34.13_{\pm 0.88}$	$63.91_{\pm 0.84}$	$84.18_{\pm 1.63}$		

Table 7: Performance on the HumanEval task under the task-aware setting for models trained on NICE-selected subsets with different temperatures for generating MC samples. We can decrease the temperature to improve pass@1 metric, while at the cost of decreased performance on pass@10 and pass@100.

Model	1	Temperature	1.0	Temperature 1.2		
WIGHEI	pass@1	pass@10	pass@100	pass@1	pass@10	pass@100
Llama2-7B	14.13	31.21	49.95	13.12	31.41	53.96
Mistral-7B	35.70	63.60	81.04	35.23	63.85	85.69

others in the codes. Consequently, BM25 performs well, especially when the training data contains data from multiple different domains that are not coding-related. However, this good result does not transfer to either the Llama2-7B model or the task-aware setting due to two main reasons: 1) BM25 does not use the information from the models, meaning that the same data subset will be selected for different models. Intuitively, different models require different data to achieve better performance. Therefore, selecting data using BM25 is sub-optimal; 2) Selection based on the word frequency is not enough for the task-aware setting. In a task-aware setting, training data points are more relevant to the task, possibly resulting in comparable BM25 scores for all data points. Consequently, a more careful selection based on other criteria (i.e., not just word frequency) is needed. Additionally, while BM25 achieves a 31.81 pass@1 for Mistral-7B, pass@1 is not the sole evaluation criterion in practice, because there are situations where multiple responses can be generated from the LLM and checked by a verifier (e.g., test cases). Therefore, for HumanEval, focusing on pass@k with larger k (and higher accuracy) is more desirable.

Additionally, we argue that NICE can improve pass@1 performance with lower temperature while sacrificing the performance of pass@10 and pass@100. We analyze the performance on the HumanEval task for models trained on NICE-selected subsets when using different temperatures to generate MC samples. As shown in Tab. 7, the pass@1 performance on downstream tasks is improved by using a lower temperature, while at the cost of reduced performance on pass@10 and pass@100 metrics. A lower temperature reduces uncertainty during generation. If the model is good at certain problems, it increases the probability of answering these problems correctly. However, this reduction

Table 8: Comparison of asymptotic time complexity and wall-clock time (in GPU hours) for each stage in data selection. The time for Warmup training with LoRA is measured on H100, and the others are measured on L40.

Stage	Warmup LoRA Training	Training Grad Comp	Validation Grad Comp	Data Selection			
Remark Asymptotic Compute	$NICE = LESS O(D_W E) 3h$	$NICE = LESS O(D_N E) 48h$	$\begin{split} \text{NICE} > \text{LESS} \\ \text{LESS: } O(D_V E); \text{NICE: } O(D_V EM) \\ \text{LESS: } 0.11\text{h on avg}; \text{NICE: } 14.67\text{h on avg} \end{split}$	$ \begin{array}{l} \text{NICE} = \text{LESS} \\ O(D_N D_V d) \\ < 0.02 \text{h} \end{array} $			

in uncertainty comes at the expense of diversity, as the generated responses tend to be very similar to one another. Consequently, for difficult questions, if all generated responses are incorrect, pass@k (for larger k) suffers. This trade-off implies that using a lower temperature to generate responses improves pass@1 performance. When these responses generated under a lower temperature are used to compute policy gradients for data selection, the resulting selected subset also favors the pass@1 metric, while at a cost of decreased performance on pass@k (for larger k).

H.2 ROBUSTNESS ACROSS VALIDATION SPLITS



Figure 4: Performance of models trained on different NICE-selected subsets using different validation sets as references. Models trained on NICE-selected subsets consistently outperform those trained on randomly selected subsets, regardless of the validation set used.

We demonstrate the robustness of NICE across different validation splits and address concerns about potential overfitting to a specific validation set. In the task-agnostic setting, we randomly selected an alternative validation set D'_V as a reference and re-selected a subset D'_S . The performance of models retrained on D_S (selected based on the original validation split D_V) and D'_S (selected based on the new validation split D'_V) are shown in Fig. 4. Importantly, models trained on the selected subsets, whether D_S or D'_S , consistently outperform models trained on randomly selected subsets. This verifies the robustness of our approach and confirms that its effectiveness does not depend on a specific validation split.

H.3 TIME COMPLEXITY ANALYSIS

We provide a comparative analysis of the computational costs between NICE and LESS, an approach that adopts loss-based influence estimation, showing that NICE remains within a practical computational range. Tab. 8 lists the asymptotic complexity and wall-clock runtime (the time for warmup training with LoRA is measured in single H100 GPU hours, others are measured in single L40 GPU hours) for each stage in the data selection procedure. Tab. 9 highlights the validation gradient computation where NICE differs from LESS. Let *E* denote the number of epochs (saved checkpoints), *d* the dimension of the projected gradients, and *M* the number of Monte Carlo (MC) samples. Let $|D_W|$, $|D_N|$, and $|D_V|$ denote the warmup, training, and validation set sizes, respectively. When $|D_V|$ and *M* are small, NICE adds only marginal overhead to LESS (e.g., AlpacaEval).

While NICEAMC utilizing Monte Carlo sampling can indeed increase the computational cost, this trade-off is justified by our approach not needing validation labels—a key motivation of our work. NICE fills a gap left by existing loss-based baselines by supporting data selection with unlabeled validation data in cases where the evaluation metrics are label-independent. Furthermore, we can observe the performance improvement over other methods in Tab. 3.

Task	$\mid M$	NICE (MC Sampling)	NICE (Val Grad)	LESS (Val Grad)
$\begin{array}{l} \mbox{AlpacaEval} \left(\left D_V \right = 10 \right) \\ \mbox{TLDR} \left(\left D_V \right = 322 \right) \\ \mbox{RLHF} \left(\left D_V \right = 2192 \right) \\ \mbox{HumanEval} \left(\left D_V \right = 10 \right) \end{array}$	20	0.17h	0.05h	<0.02h
	20	8h	1.47h	0.08h
	20	32h	10h	0.33h
	500	5h	2h	<0.02h

Table 9: Validation gradient computation time across tasks for NICE and LESS in single L40 GPU hours.

Table 10: Additional Results on RLHF dataset for Llama2-13B and Llama3-8B. **Bold** numbers indicate the top-performing selected subset. A purple cell suggests that NICE outperforms LESS which uses loss-based influence estimation.

RIHE	Llama2-13B		Llama3-8B		
KLIII	Task-agnostic	Task-aware	Task-agnostic	Task-aware	
Random	$2.06_{\pm 0.04}$	$1.20_{\pm 0.07}$	$1.97_{\pm 0.07}$	$1.12_{\pm 0.06}$	
RDS	1.77	0.70	1.75	0.81	
BM25	2.72	1.34	2.84	1.43	
LESS	$1.52_{\pm 0.09}$	1.65 ± 0.04	$1.64_{\pm 0.14}$	$1.65_{\pm 0.08}$	
NICE	$2.87_{\pm 0.04}$	$1.76_{\pm 0.04}$	$3.22_{\pm 0.02}$	$1.99_{\pm 0.06}$	

H.4 ADDITIONAL RESULTS ON LLAMA3-8B AND LLAMA2-13B

We evaluate NICE against various data selection baselines using the state-of-the-art model, Llama3-8B (Dubey et al., 2024), and a larger model, Llama2-13B, on the RLHF dataset (see Tab. 10). The superiority of NICE underscores our method's generalizability across different model sizes and state-of-the-art models.

H.5 THE DISCREPANCY BETWEEN NTP LOSS AND EVALUATION METRICS.

We additionally include the NTP loss and performance (i.e., measured by each task's evaluation metric here) of the validation set for the remaining three tasks in several training checkpoints. The results in Fig. 5 are similar to those in Fig. 1: checkpoints with minimal loss (highest negative losses) do not correspond to checkpoints with the best performance; the performance can continue to increase even if the loss increases (negative loss decreases).



Figure 5: Discrepancy observed between the NTP loss and performance (as measured by each task's evaluation metric) on the checkpoints. The checkpoint with the lowest loss (i.e., most negative) can exhibit relatively poorer performance. Notably, performance can continue to improve even as the loss worsens (i.e., the negative loss becomes higher).

H.6 COMPARISON OF VALIDATION PERFORMANCE ACROSS FINAL CHECKPOINTS FOR NICE AND LESS

We plot the validation performance, measured by the reward of the last few checkpoints, for models trained on NICE-selected subsets and LESS-selected subsets in Fig. 6 for the RLHF dataset. NICE-selected data optimizes in the direction of increasing validation performance, whereas LESS-selected data prioritizes loss reduction, which may not necessarily lead to improved validation performance.



Figure 6: Validation performance, measured by the reward model of the last few checkpoints, for models trained on subsets selected by NICE and LESS for the RLHF dataset. NICE-selected data optimizes for improved validation performance, whereas LESS-selected data focuses on loss reduction, which may not always enhance validation performance.



Figure 7: Standard deviation of different runs w.r.t the number of MC samples. Increasing the number of MC samples generally lowers the standard deviation across runs, indicating better stability.

H.7 STABILITY OF MC SAMPLING

We provide an ablation study with results in Fig. 7, varying MC samples from 5 to 20 on the RLHF dataset in the task-agnostic setting for a more in-depth discussion on stability. Results show that increasing the number of MC samples generally lowers the standard deviation across runs with different seeds, indicating better stability. The benefit of reduced standard deviation diminishes as it increases. This validates that our chosen MC (MC=20 for RLHF task) provides a good trade-off, offering sufficient stability without excessive computation.

H.8 ON THE ADDITIONAL COST OF NICEAMC

Note that NICEAMC is an optional enhancement—NICE itself does not require GPT-4. We list the projected GPT-4 cost for NICEAMC in Tab. 11. The costs are low for the majority of the tasks, except for RLHF, due to its large validation set (which can be addressed by using alternative models as discussed in the next paragraph).

Table 11: Projected GPT-4 cost for NICEAMC across different tasks. The cost for RLHF is high due to a large validation set.

Task	AlpacaEval	TLDR	RLHF	HumanEval	Avg
GPT Cost (\$)	1.70	14.26	291.17	6.34	78.37

Use of Open-Source/Smaller LLMs. To reduce cost, we can use high-performing open-source models. On the RLHF dataset, we use Qwen 2.5-3B/7B-Instruct (Yang et al., 2024) for AMC. Both outperform NICE. Notably, even a small model like Qwen 2.5-3B-Instruct performs better due to its better alignment training, despite its smaller size. These models offer comparable performance to GPT-4 without incurring the additional API cost.

Table 12: Performance on the RLHF dataset using different models for NICEAMC. Qwen models offer competitive performance without the API cost of GPT-4.

Model	NICE	NICE AMC (GPT-4)	NICE AMC (Qwen2.5 7B)	NICE AMC (Qwen2.5 3B)
RLHF	$2.82_{\pm 0.10}$	$3.03_{\pm 0.02}$	$3.00_{\pm 0.03}$	$2.97_{\pm 0.03}$

I ABLATION STUDIES

Unless otherwise specified, the ablation studies in the section below are conducted on the Llama2-7B model.

I.1 USING AN ALTERNATIVE WAY TO COMPUTE POLICY GRADIENT

Table 13: Performance comparison between loss-based influence estimation (LESS) and NICE, which uses different approaches to compute policy gradient in the task-aware setting. *PG* refers to using the MC policy gradient as described in Sec. 3.1, while *PPO* denotes using Proximal Policy Optimization. Employing policy gradient computed from either policy optimization approach during data selection results in a better-selected subset compared to loss-based influence estimation.

dataset	RLHF	HumanEval
LESS	$1.62_{\pm 0.05}$	$52.67_{\pm 0.71}$
PG	$1.69_{\pm 0.05}$	$55.09_{\pm 1.66}$
PPO	$1.73_{\pm 0.02}$	$52.08_{\pm 1.31}$

Besides the vanilla Monte-Carlo policy gradient, we conducted an ablation study of using another policy optimization methodology to compute the policy gradient used in NICE. We tried computing the policy gradient using Proximal Policy Optimization (PPO) (Schulman et al., 2017), with results presented in Tab. 13 for the task-aware setting. These results demonstrate that compared to the loss-based influence estimation (LESS), which selects training data by optimizing in the direction of decreasing validation loss, integrating either the MC policy gradient or PPO gradient to select training data by optimizing in the direction of improving downstream task performance leads to a better-selected subset. A "better-selected subset" refers to a subset of training data that, when used for model training, results in improved performance on downstream tasks.

I.2 INTRODUCE ADDITIONAL KNOWLEDGE ON LOSS-BASED INFLUENCE ESTIMATION

Tab. 14 shows the performance of LESS using GPT-generated labels (LESS+GPT), which is generally worse than our approaches and can even be worse than LESS + true labels (LESS). Hence, simply

using GPT-generated labels with loss-based approaches cannot always address the unavailability of labels for validation data.

Table 14: Performance of LESS using GPT-generated labels (2nd row) in the task-agnostic setting, which is generally worse than our approaches and can even be worse than LESS + true labels (1st row).

Table D	Alpaca	TLDR	RLHF	HumanEval
LESS	$26.94_{\pm 2.37}$	$3.37_{\pm 0.78}$	$1.44_{\pm 0.07}$	$47.50_{\pm 1.57}$
LESS+GPT	27.35 ± 1.86	$3.41_{\pm 0.26}$	$3.03_{\pm 0.01}$	$43.04_{\pm 1.39}$
NICE	$27.61_{\pm 2.12}$	$3.61_{\pm 0.78}$	$2.82_{\pm 0.10}$	$48.59_{\pm 2.08}$
NICEAMC	$30.45_{\pm 2.39}^{-}$	$3.55_{\pm 0.40}$	$3.03_{\pm0.02}$	$45.10_{\pm 2.84}$

I.3 EFFECT OF REWARD SCORE

We conduct a simple ablation study to evaluate the effectiveness of the reward score within the context of the policy gradient methodology. In the task-agnostic setting for the RLHF task, we compare the performance between the continuous reward score (our current methodology) and a discrete reward score, where the reward is set to 1 if positive and 0 if non-positive. This discrete reward setup can also be interpreted as a form of rejection sampling. The subset selected using the continuous reward yields a model with a performance score of $2.82_{\pm 0.10}$, whereas the subset selected using the discrete reward yields a performance of $2.25_{\pm 0.12}$. Although less effective than the original policy gradient method, the policy gradient with a discrete reward still demonstrates some effectiveness, outperforming both LESS $(1.44_{\pm 0.07})$ and random $(2.05_{\pm 0.11})$ selection strategies. This ablation study highlights that subsets selected with the guidance of reward scores can produce models with better performance than those trained on subsets selected by LESS or random sampling.

J QUALITATIVE ANALYSIS

J.1 GENERATED RESPONSES

We present three generated responses from GPT-4 on an AlpacaEval validation data point in Table 15 and Table 16. By comparing the original labels with the GPT-4 responses, particularly in Table 15, we observe that the generated responses exhibit higher quality and more accurate results. By observing Tab. 16, we find that the generated responses are not only better in quality, provide richer detail, and more vivid imagery for monologues, but they also offer diverse responses. These monologues differ in style and narrative focus, ranging from a grand epic tone to reflective introspection. The first explores ancient ruins and lost civilizations, the second emphasizes nature and mysticism, and the last tells mythic tales.

J.2 HIGHEST/LOWEST SCORING TRAINING DATA FOR RLHF TASK IN TASK-AGNOSTIC SETTING

In a task-agnostic setting, we present the highest- and lowest-scoring training data points in Table 17. The highest-scoring training data points complete the task effectively and offer additional knowledge, while the lowest-scoring ones tend to be vague, minimal, or fail to address the task meaningfully. Considering the RLHF task targets helpfulness, the high-scoring data points can potentially boost the model's helpfulness. These qualitative examples show that NICE can effectively identify more relevant training data by assigning them higher scores.

J.3 HIGHEST/LOWEST SCORING TRAINING DATA FOR RLHF TASK IN TASK-SPECIFIC SETTING

In a task-specific setting, we present the highest- and lowest-scoring training data points in Table 18. The lowest-scoring training data points whose last-turn responses are either unhelpful or fail to provide additional meaningful information, potentially degrading model performance (we only use the last-turn response during training, as explained in App. A.1). Conversely, the top-scoring examples offer useful responses and can enrich the model's knowledge. Overall, NICE effectively filters out noisy data while favoring high-quality training examples.

Table 15: We compare three GPT-4-generated responses with the original labels of an AlpacaEval validation data point, observing that GPT-4 consistently produces higher-quality and more accurate outputs.



. Now, plug these values into the formula:

$$d = \sqrt{(10-2)^2 + (4-(-2))^2}$$
$$d = \sqrt{8^2 + 6^2}$$
$$d = \sqrt{64+36}$$
$$d = \sqrt{100}$$
$$d = 10$$

Thus, the length of the segment is 10 units.

Table 16: We present three diverse and vividly detailed generated monologues which answers an AlpacaEval validation question, demonstrating how GPT-4 can vary tone, style, and narrative focus in generated responses.

Question: "Write a monologue. Do not offer any explanation or out of character notes. The setting is in the Elder Scrolls universe."

Baseline Response: "Ah, what a glorious day! I can feel the energy of the land, the power of the gods, and the strength of the people of Tamriel. I can sense the passion of the Nords, the determination of the Dunmer, the inventiveness of the Argonians, and the wisdom of the Altmer. The air is filled with the sounds of battle, adventure, and hope. I am so proud to be part of this world, to be a part of the Elder Scrolls. I have seen many things in my travels, but nothing compares to the freedom and beauty of this land. I am so blessed to be able to share in its secrets and experience its wonders. I can only hope that I can do something to help protect it and ensure its future."

Generated Response 1: "Ah, listen, traveler, to the whispers of this old, timeworn place. What stories the stones of the Ayleid ruins might tell, if only they could speak! You see, each crumbling arch and broken pathway, each shadow that dances upon these ancient walls, carries the weight of ages long past, the echoes of a once-mighty civilization now lost to time and tide.

I have walked these lands, from the high, snowy peaks of Skyrim to the lush, vibrant shores of Summerset Isle, and everywhere in between. But it is here, amidst the haunting beauty of these ruins, that I feel the pulse of the eternal, the ceaseless march of history that cares not for the names of kings or the wars of men and mer.

Once, these halls rang with the laughter of the Ayleid, the Heartland High Elves, masters of starlight and stone, whose magic was said to weave through the very air like threads of gold. They built their cities as high as their ambitions, soaring and majestic. But pride, ah, pride—it festers, and it corrupts. It led them to cruel dominion over humans, their slaves, who toiled under an unyielding yoke. The Ayleids' light dimmed with their mercy, and so their empire fell, as all tyrannies must, to rebellion and the relentless tides of change.

Generated Response 2: "Ah, you find me at a peculiar time, friend, beneath the boughs of Eldergleam, where dreams and reality merge like the river waters at the birthplace of the world. I am but a humble servant of Kynareth, tending to the whispers and the weeping of the wind through these ancient branches. You seek wisdom, or perhaps power? Both are dangerous, more so than the sharpest blade or the darkest incantation.

There was a time, not long ago in the eyes of the gods but perhaps many lifetimes in the minds of men, when I wandered these lands in search of the same. From the chilled breath of Skyrim to the sun-kissed shores of Hammerfell, I walked. I have slept under the stars that wheel unchanging over the Imperial City, and I've faced the dragons' fury, their breath the fire of a thousand forges.

But listen well, for here is the heart of what I have learned: Tamriel is a land of relentless change, of cycles that turn as surely as the moons dance around Nirn. Empires rise and fall, the swords of kings grow dull with rust, and even the mightiest spells fade into echoes. Yet, the trees, the stones, and the rivers—they remember. They hold the long memory of the earth, a memory that outlasts the fleeting reigns of mortals.

Do you feel the roots of the Eldergleam? Do you hear their slow, ceaseless whisper? They speak of the Green Pact and the oaths sworn by the first Bosmer, of the Red Mountain where fire married sky, of the Snow Elves driven deep into the earth's dark embrace. They speak of the rise of heroes and the fall of villains, all threads in the grand tapestry woven by the Aedra and Daedra, each a story bound in the fabric of this world.

Generated Response 3: "Ah, traveler! Welcome to the warmth of my hearth. Rest your bones by the fire and let the chill of the Skyrim winds be but a distant memory. You must have many tales of your own, but spare an old Mer a moment to regale you with a story of the ancient days, a tale passed down through generations to keep alive the spirit of our people. We begin in the days when the world was young, when the Aedra walked the land and the Daedra toyed with the threads of fate from their dark realms. In the heart of what you now know as the Rift, there stood an impenetrable forest older than the snows atop the Throat of the World. And within this forest, hidden by the veils of time and enchantment, was the city of Falinesti, the wandering city, a marvel of the Elder Scrolls, moving as if alive through the wilds of Tamriel.

Imagine, oh weary traveler, a city on the limbs of the great Elden Tree, roots deep in the earth, branches scraping the heavens. It was there, in that moving city, that the tale of Aelarion was born. Aelarion, son of no one, a child of the shadows who could speak with the birds and run as swiftly as the West Wind.

Now, Aelarion was a creature of curiosity and mischief. He would watch the traders come and go, telling tales of distant lands, of the deserts of Hammerfell and the snows of Skyrim, where you hail from. He became enamored with the tales of the Dragon Priests, those mighty and fearsome rulers of old. Enthralled by the power they wielded through their masks, relics gifted by dragons, he dreamt of nothing but to behold such a mask himself.

Table 17: We list the training data points with the highest or lowest NICE score in a task-agnostic setting for RLHF task, illustrating that highest-scoring responses are helpful and can convey additional knowledge. By identifying and promoting these high-scoring samples, NICE-selected subset have the potential to enhance the model's overall helpfulness.

Dataset	Dialogue				
	Training data points with the highest NICE score				
Flan V2	User: Detailed Instructions: In this task, you're given an ambiguous question (which can be answered in				
	more than one way). Your task is to provide one question which clarifies the input question and it has one unique answer, and also provide an answer to the generated question. Generated question and answer should				
	be separated with a new line. Problem: When was free health care introduced in Canada? Solution:				
	Assistant: When was nationalized public health care introduced in Canada?				
	1947				
Flan V2	User: Q: In this task, you're given an ambiguous question (which can be answered in more than one way).				
	Your task is to provide one question which clarifies the input question and it has one unique answer, and also				
	provide an answer to the generated question. Generated question and answer should be separated with a new				
	line.				
	Who got the biggest NBA contract right now?				
	Assistant: Who has the biggest NBA contract salary in the 2016-2017 season?				
	LeBron James				
DOLLY	User: Task: What team did Bobby Orr play for in his 5th season?				
	Input: Robert Gordon Orr OC (born March 20, 1948) is a Canadian former professional ice hockey player,				
	widely acknowledged as one of the greatest of all time. Orr used his ice skating speed, scoring, and play-				
	making abilities to revolutionize the position of defenceman. He played in the National Hockey League				
	(NHL) for 12 seasons, the first 10 with the Boston Bruins, followed by two with the Chicago Black Hawks.				
	Orr remains the only detencement to have won the league scoring title with two Art Ross Trophes. He holds				
	the record for most points and assists in a single season by a defenceman. Orr won a record eight consecutive				
	Norms frophies as the NHL's best detendement and three consecutive Hart frophies as the league's most valuable player (MVP). Ore use inducted into the Hockey Hall of Fame in 1070 at age 31, the youngest to be				
	valuable player (MVP). Our was inducted into the Hockey Hall of Pame in 1979 at age 31, the youngest to be				
	NUL Deven" in history				
	NHL Players in instory.				
	Assistant: Bobby Orr played for the Boston Bruins in his 5th season in the National Hockey League (NHL)				
	He played in the NHI for 12 seasons the first 10 with the Boston Bruis followed by two with the Chicago				
	Black Hawks				
	Training data points with the lowest NICE score				
Flan V2	User: Translate "," to Finnish?				
	Assistant:				
DOLLY	User: Is tap water safe to drink in Japan?				
	Assistant: Yes, tap water is safe to drink in Japan.				
OASST	User: I have trouble staying on track when I need to do schoolwork. What are some ways that I can keep				
	myself on track?				
	Assistant: Getting rid of anything that distracts you could be helpful.				

Table 18: We show the training data points with the highest or lowest NICE score in a task-specific setting for RLHF task, demonstrating that top-scoring samples enrich model knowledge while low-scoring ones can be noizy. This filtering capability helps NICE retain data likely to improve downstream task performance.

Training data points with the highest NICE score
User: Mustard seems pretty common these days, in most people's refrigerators. When did people start using
mustard as a condiment on foods, I'm curious about its history.
Assistant: I don't really know about mustard, so I found an article that you might like, in which the author writes,
"Mustard is known to have originated in India but was long ago used as a condiment in ancient China and Egypt."
Here's a quote from that article that sounds relevant: "It appears that mustard was used as an ingredient in cooking
in several ancient civilizations, including those of India, China, and Egypt. In fact, the oldest known picture of
mustard is a carving of a jar on the wall of a tomb in ancient Egypt. This artifact, which dates to the 15th century
B.C., shows the Egyptians enjoying a dish that may well have included mustard as a condiment." I'd be happy to
read that article for you if you want.
User: What are some good stretches I should do before going to the gym?
Assistant: Stretches are an important part of a workout regimen. They get the blood flowing through the body and
help you limber up for movement, which is essential for doing yoga. What types of stretching are you looking to do?
User: I was hoping you could tell me.
Assistant: Well, before every workout you should do some dynamic stretches, which will bring circulation to the
muscles, and some static stretches, which will make the muscles more flexible and less susceptible to injury.
User: Ok, which stretches specifically?
Assistant: Ok, what about calf raises?
User: Can you teach me how to make candles?
Assistant: Sure, just follow these directions and it should be easy!
Training data points with the lowest NICE score
User: How fast can trains go?
Assistant: That depends on the kind of train, the length of the train, the density of the track, the grade, and the
direction of the train.
User: What about coal trains?
Assistant: I'm not sure what you mean by "coal trains."
User: What's the best brand of TV?
Assistant: I'm not sure what you mean by "best brand."
User: How can I make my clothes white?
Assistant: You can use vinegar to "bleach" the clothes.
User: vinegar can bleach clothes?
Assistant: Yes