# Measuring Data Diversity for Instruction Tuning: A Systematic Analysis and A Reliable Metric

**Yuming Yang[1*], Yang Nan[1*], Junjie Ye[1], Shihan Dou[1], Xiao Wang[1],**
**Shuo Li[1], Huijie Lv[1], Tao Gui[2,3,4†], Qi Zhang[1,3,4†], Xuanjing Huang[1,3,4†]**

[1] College of Computer Science and Artificial Intelligence, Fudan University
[2] Institute of Modern Languages and Linguistics, Fudan University
[3] Institute of Trustworthy Embodied Artificial Intelligence, Fudan University
[4] Shanghai Collaborative Innovation Center of Intelligent Visual Computing
yumingyang23@m.fudan.edu.cn {qz,tgui,xjhuang}@fudan.edu.cn

## Abstract

Data diversity is crucial for the instruction tuning of large language models. Existing studies have explored various diversity-aware data selection methods to construct high-quality datasets and enhance model performance. However, the fundamental problem of precisely defining and measuring data diversity remains underexplored, limiting clear guidance for data engineering. To address this, we systematically analyze 11 existing diversity measurement methods by evaluating their correlation with model performance through extensive fine-tuning experiments. Our results indicate that a reliable diversity measure should properly account for both inter-sample differences and the information density in the sample space. Building on this, we propose *NovelSum*, a new diversity metric based on sample-level "novelty." Experiments on both simulated and real-world data show that *NovelSum* accurately captures diversity variations and achieves a 0.97 correlation with instruction-tuned model performance, highlighting its value in guiding data engineering practices. With *NovelSum* as an optimization objective, we further develop a greedy, diversity-oriented data selection strategy that outperforms existing approaches, validating both the effectiveness and practical significance of our metric. The code is available at https://github.com/UmeanNever/NovelSum.

## 1 Introduction

Instruction tuning (IT) fine-tunes pretrained large language models (LLMs) with annotated instruction data, enabling them to follow human instructions and perform various tasks effectively (Sanh et al., 2022; Zhang et al., 2023). Recent studies indicate that small-scale, high-quality datasets can outperform larger ones in IT performance (Chen et al., 2023a; Zhou et al., 2024), with data diver-

---

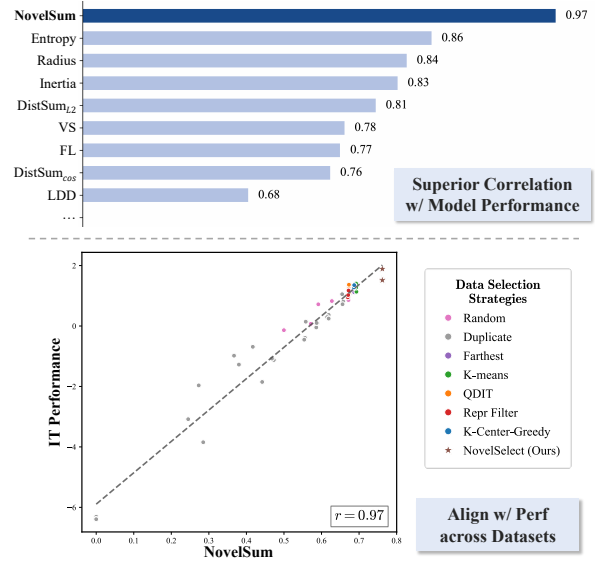*Equal Contribution.
†Corresponding Authors.



Figure 1: Our diversity metric, *NovelSum*, exhibits superior correlation with model performance compared to existing metrics across IT datasets constructed with various data selection strategies.

sity playing a crucial role in achieving optimal results (Liu et al., 2023; Bukharin et al., 2024; Zhang et al., 2024a; Yang et al., 2025). Consequently, various diversity-aware data selection methods have emerged (Qin et al., 2024; Wang et al., 2024a), driven by different interpretations of data diversity.

However, the fundamental problem of precisely defining and measuring data diversity remains underexplored. This ambiguity has turned data engineering for diversity into a black-box process, leading to data selection methods that often fail to generalize and, at times, perform worse than random selection (Xia et al., 2024; Diddee and Ippolito, 2024). While some diversity metrics have been introduced in IT research (Bukharin et al., 2024; Wang et al., 2024b), a comprehensive evaluation and comparative analysis are still needed to identify a reliable metric that strongly correlates with fine-tuning performance in practice.

To this end, we systematically analyze and eval-

uate the reliability of 11 existing diversity metrics through extensive experiments. Using various mainstream diversity-oriented data selection methods, we construct 53 IT datasets and fine-tune models accordingly. We then measure dataset diversity using existing metrics and assess their correlation with model performance. By analyzing the limited correlation of existing metrics, we find that: (1) **A reliable diversity metric must capture differences between samples** to reflect each sample's information uniqueness. Moreover, differences between neighboring samples are more critical for overall diversity but can be overshadowed by variations in distant samples. (2) **Measuring differences between samples should account for both semantic similarity and the uneven distribution of information in space.** In high-density domains like math and code, semantically similar samples can still contain substantial unique information and should therefore be considered more diverse.

Building on these insights, we propose *NovelSum*, a diversity metric that jointly considers intersample differences and uneven information density. Specifically, we define dataset diversity as the sum of each sample's unique contribution to overall information, termed "novelty". Just as a research paper's novelty is judged by its distinction from related work based on field-specific standards, we compute a sample's novelty as the proximity-weighted sum of its differences from other samples in the dataset. These differences are measured using density-aware distances, which capture both semantics and local information density.

To validate the effectiveness of *NovelSum*, we conduct both a visualized simulation study and real-world correlation experiments using two different LLMs. The results show that *NovelSum* accurately captures diversity variations and strongly correlates with instruction-tuned model performance, achieving Pearson's $r = 0.98$ and Spearman's $r = 0.95$, outperforming other metrics. This demonstrates *NovelSum*'s potential to effectively guide data engineering practices. Furthermore, we develop *NovelSelect*, a greedy, diversity-oriented data selection strategy that uses *NovelSum* as the optimization objective. Experimental results confirm its superior performance compared to other approaches.

Our main contributions are three-fold:

- We systematically analyze and evaluate the reliability of existing diversity metrics for instruction tuning by computing their correla-

tion with model performance, thereby unveiling pathways to a more reliable metric.

- We propose *NovelSum*, a diversity metric that captures both inter-sample differences and information density, achieving a strong correlation with instruction-tuning performance, substantially exceeding previous metrics.

- We develop *NovelSelect*, a diversity-oriented data selection strategy based on *NovelSum*, which outperforms existing methods and further validates *NovelSum*'s effectiveness and practical value in instruction tuning.

## 2 Evaluating Existing Diveristy Metrics

We begin by evaluating the correlation between existing diversity metrics and instruction-tuned model performance, identifying limitations to inform the design of a more reliable metric.

Our evaluation follows four steps: (1) Construct multiple IT datasets, each denoted as $\mathcal{X}^{(s)}$, using different data selection strategies from the full data source $\mathcal{X}^{all}$. (2) Measure dataset diversity using existing metrics, denoted as $\mathcal{M}_t(\mathcal{X}^{(s)})$. (3) Fine-tune LLMs on each dataset and evaluate their performance, $\mathcal{P}^{(s)}$, using IT benchmarks. (4) Analyze the correlation between each diversity metric and model performance, denoted as $r_{\mathcal{M}_t, \mathcal{P}}$.

### 2.1 Existing Diversity Metrics

We use 11 existing diveristy metrics for the analysis, categorized into three main types:

**Lexical Diversity** A classical way to measure textual diversity is by analyzing vocabulary usage, where a higher proportion of unique words indicates greater diversity. Two widely used metrics are the **Type-Token Ratio** (TTR) (Richards, 1987) and **vocd-D** (Malvern et al., 2004), with details in the Appendix A.2.

**Distance-based Semantic Diversity** Recent studies primarily measure dataset diversity based on the semantics of individual samples, often represented as embeddings $emb(\cdot)$ from language models like BERT. A common approach quantifies diversity by computing distances between samples using their embeddings, encouraging heterogeneity. For example, a straightforward metric sums the pairwise distances among all samples in a dataset:

$$\mathcal{M}_{DistSum}(\mathcal{X}) = \sum_{x_i, x_j \in \mathcal{X}, i \neq j} \Delta(x_i, x_j), \quad (1)$$

where $\Delta(\cdot, \cdot)$ denotes the distances between two samples. Specifically, **DistSum**$_{cosine}$ uses cosine distance and **DistSum**$_{L2}$ uses Euclidean distance. Beyond simple summation, more refined metrics are proposed. The **KNN distance** (Stasaski et al., 2020; Stasaski and Hearst, 2022) measures the average distance of each sample to its $k$-nearest neighbor, ensuring sample uniqueness:

$$\mathcal{M}_{KNN}(\mathcal{X}) = \frac{1}{|\mathcal{X}|} \sum_{i=1}^{|\mathcal{X}|} \Delta(x_i, N_k(x_i)), \quad (2)$$

where $N_k(x_i)$ denotes the k-th closest neighbor of $x_i$, typically with $k = 1$. We also compute **Cluster Inertia** (Du and Black, 2019), **Vendi Score** (Pasarkar and Dieng, 2023), **Radius** (Lai et al., 2020) and **Log Determinant Distance** (LDD) (Wang et al., 2024b); see Appendix A.2 for details.

**Distribution-based Semantic Diversity** Another notable class of metrics measures diversity from a distributional perspective, assessing how well a selected dataset $\mathcal{X}$ represents the overall sample (semantic) space of $\mathcal{X}^{all}$. One example is the **Facility Location** (FL) function (Farahani and Hekmatfar, 2009), which defines a dataset as diverse if each sample in $\mathcal{X}^{all}$ has a close representative in $\mathcal{X}$, ensuring thorough coverage of space:

$$\mathcal{M}_{FL}(\mathcal{X}) = \sum_{x_j \in \mathcal{X}^{all}} \min_{x_i \in \mathcal{X}} \Delta(x_i, x_j) \quad (3)$$

Another feasible metric, **Partition Entropy**, captures how evenly the selected dataset spans the sample space. It partitions $\mathcal{X}^{all}$ into $K$ clusters using K-means and calculates the entropy of the cluster membership distribution of $\mathcal{X}$.

$$\mathcal{M}_{Entropy}(\mathcal{X}) = -\sum_{k=1}^{K} p_k \log p_k, \quad (4)$$

where $p_k$ is the proportion of selected samples in cluster $k$. Higher entropy indicates greater distributional uncertainty and a more balanced dataset.

## 2.2 IT Dataset Construction and Benchmark

Focusing on general IT, we follow Liu et al., 2023 to construct our IT data source by combining WizardLM (Xu et al., 2024), ShareGPT (Chiang et al., 2023), and UltraChat (Ding et al., 2023), denoted as $\mathcal{X}^{all}$. We extract embeddings for each sample in $\mathcal{X}^{all}$. See Appendix A.1 for preprocessing details.

We then apply several representative diversity-aware data selection strategies to curate IT datasets,

yielding subsets $\mathcal{X}^{(s)} \subset \mathcal{X}^{all}$. To minimize the influence of factors beyond diversity, we control for sample quality differences across datasets by removing anomalous source samples and excluding any data quality filters during selection. We also fix the dataset size at 10,000 samples. The strategies used are: **K-Center-Greedy** (Sener and Savarese, 2017; Du et al., 2023; Wu et al., 2023), which iteratively selects the sample farthest from the current coreset; **Repr Filter** (Liu et al., 2023), which improves $\mathcal{M}_{KNN}$ by applying a minimum distance threshold when adding samples into the coreset; **QDIT** (Bukharin et al., 2024), which optimizes diversity by serially selecting the data point that maximizes $\mathcal{M}_{FL}$; **K-means** (Song et al., 2024; Ge et al., 2024), which partitions samples into clusters and evenly select samples from each; and baselines, including **Random** selection and **Farthest**, which ranks samples by their total distances to others and selects the most distant ones. Additionally, we construct datasets with varying amounts of **Duplicate** samples to simulate low-diversity datasets. Each strategy is run at least three times to ensure robustness, yielding 53 IT datasets. Details on dataset construction are provided in Appendix A.3.

We fine-tune LLaMA-3-8B (Dubey et al., 2024) on these datasets and evaluate model performance using two popular IT benchmarks: MT-bench (Zheng et al., 2023) and AlpacaEval (Li et al., 2023). See Appendix A.4 for details and rationale of the benchmarks. To jointly consider both benchmarks, we normalize the results into Z-scores and compute the aggregated performance as

$$\mathcal{P}^{(s)} = z_{MT-bench}^{(s)} + z_{AlpacaEval}^{(s)} \quad (5)$$

## 2.3 Correlation Analysis

Finally, we compute the correlation between each diversity metric $\mathcal{M}_t$ and model performance $\mathcal{P}$ by averaging their Pearson and Spearman coefficients:

$$r_{\mathcal{M}_t, \mathcal{P}} = (r_{\mathcal{M}_t, \mathcal{P}}^{Pearson} + r_{\mathcal{M}_t, \mathcal{P}}^{Spearman})/2 \quad (6)$$

Since our experiments minimize confounding factors, variations in model performance can be more directly attributed to differences in IT dataset diversity. Thus, the correlation $r_{\mathcal{M}_t, \mathcal{P}}$ indicates how reliably each metric captures "IT-aligned Diversity"[1]—the type of data diversity beneficial for instruction tuning LLMs.

---

[1]Unless otherwise specified, the term "diversity" in this paper generally refers to "IT-aligned Diversity."
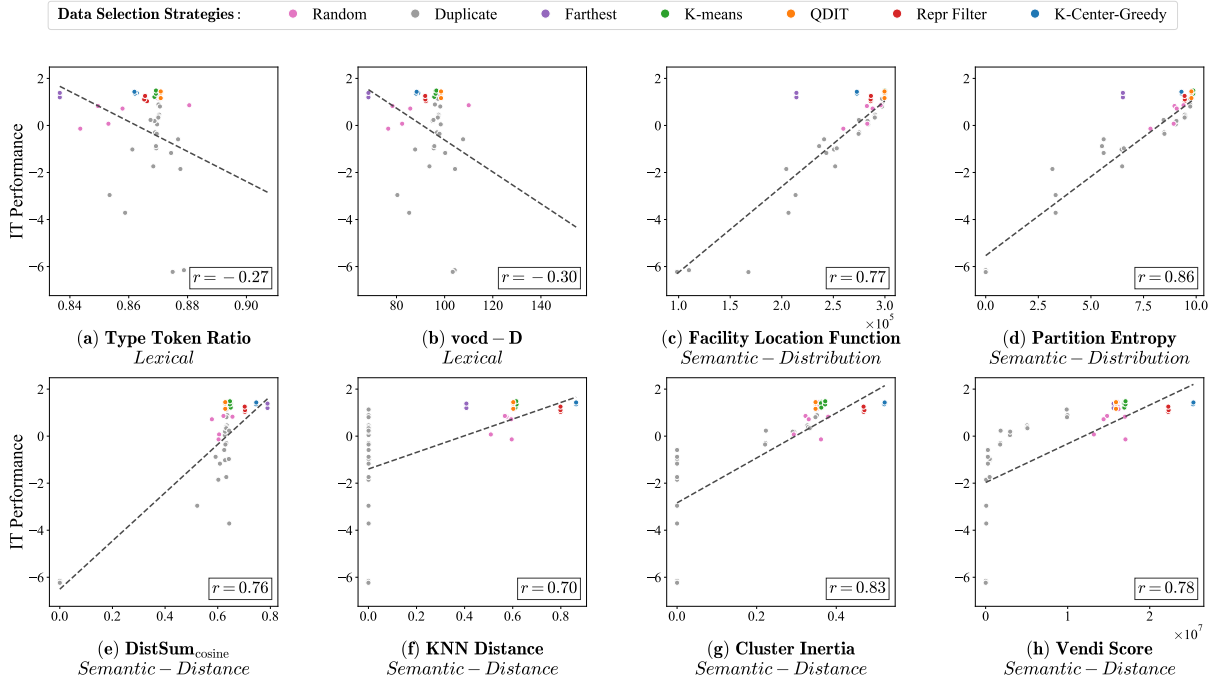
Figure 2: Evaluating existing diversity metrics based on their correlation (Eq. 6) with IT performance (Eq. 5). The X-axis represents diversity measurements. Each point corresponds to a 10k IT dataset constructed using different strategies. Abnormal points highlight the limitations of current metrics and inspire the development of new ones.

The results are shown in Figure 2, with supplementary plots in Appendix E.4. Overall, we observe that each metric tends to favor datasets aligned with its own selection criterion, but may not correlate strongly with performance due to overlooking other aspects of diversity:

**Findings 1** *Lexical diversity metrics fail to distinguish between different samples and datasets, showing weak correlation with model performance.*

As shown in Figure 2(a, b), high- and low-performance datasets exhibit similar lexical diversity. This likely results from the widespread use of diverse vocabulary in IT samples, making lexical diversity an ineffective measure for IT datasets.

**Findings 2** *Since distribution-based semantic diversity metrics neglect sample uniqueness, they often underestimate the diversity of datasets with large inter-sample distances.*

From Figure 2(c, d), we observe that datasets selected by Farthest and K-Center-Greedy (purple and blue points) achieve high IT performance but often receive relatively lower diversity scores from distribution-based diversity metrics, thus weakening their correlation with model performance. This likely occurs because these strategies all prioritize sample uniqueness by selecting samples that are distant from others, a factor not captured by distribution-based metrics. This suggests that over-

looking sample uniqueness diminishes the reliability of diversity metrics.

**Findings 3** *As distance-based semantic diversity metrics neglect information density in semantic space, they often underestimates datasets that are close to the overall sample distribution and overestimates datasets with large inter-sample distances.*

From Figure 2(e, f, g, h), we observe common outliers in the fitting line for datasets selected by QDIT and K-means (orange and green points), which receive low diversity scores despite strong performance according to distance-based diversity metrics. In contrast, K-Center-Greedy and Repr Filter (blue and red points) show the opposite trend, weakening the metrics' correlation with the model performance. This is likely because the former two strategies select more samples from dense semantic regions, which better cover the overall sample distribution but conflicts with distance-based diversity calculations. This suggests that ignoring information density in semantic space reduces the reliability of diversity metrics.

**Findings 4** *Distance-based metrics often fail to accurately measure diversity in datasets containing redundant samples.*

As shown by the duplicated datasets (gray points) in Figure 2(e, f, g, h), DistSum fails to capture redundancy effectively, as total distances are

18533

dominated by variations in distant samples. Meanwhile, other metrics, such as KNN Distance, overly penalize redundant samples by nullifying their contribution to overall diversity.

## 3 Proposed Metric: *NovelSum*

Extending previous findings, we derive some insights on how to design a more reliable metric: (1) **The uniqueness of individual samples should be a key factor in measuring dataset diversity.** This uniqueness stems from sufficient inter-sample distances, providing diverse information that helps the model learn more generalized patterns. (2) **When quantifying a sample's uniqueness, its distance to nearby and distant samples should be balanced.** Differences with nearby samples define uniqueness and should hold greater importance, with weights assigned smoothly. (3) **When calculating inter-sample distances, both semantic differences and local information density should be considered.** In practical applications of instruction fine-tuning, semantic space varies in information density, with scenarios like math and code having denser data and information. Focusing only on semantics overlooks valuable fine-grained information for the model.

Following these principles, we introduce *NovelSum*, a diversity metric that jointly considers distance and distribution. Specifically, we define dataset diversity as the sum of each sample's uniqueness—its unique contribution to overall information, which we later term "novelty":

$$\mathcal{M}_{NovelSum}(\mathcal{X}) = \sum_{x_i \in \mathcal{X}} v(x_i) \qquad (7)$$

Figure 3 and the following paragraphs illustrate how each sample's novelty is computed.

**Proximity-Weighted Sum** In contrast to Dist-Sum (Eq. 1), which calculates a sample's uniqueness as a simple sum of distances to other points, we propose a proximity-weighted sum that assigns higher weights to closer points, giving them a larger influence on the uniqueness score:

$$v(x_i) = \sum_{x_j \in \mathcal{X}, \, x_j \neq x_i} w(x_i, x_j)^\alpha \cdot \Delta(x_i, x_j), \quad (8)$$

where the proximity weight is defined as:

$$w(x_i, x_j) = \phi(\pi_i(j))$$

Here, $\pi_i(j)$ is the rank of $x_j$ in the sorted list of distances from $x_i$ to all other points in $\mathcal{X}$, with



Proximity-Weighted Sum   Density-Aware Distance

$$v(x) = \sum_j \boldsymbol{w_j} \cdot \Delta(x, x_j) \qquad \Delta(x, x_j) = \boldsymbol{\sigma(x_j)} \cdot d(x, x_j)$$
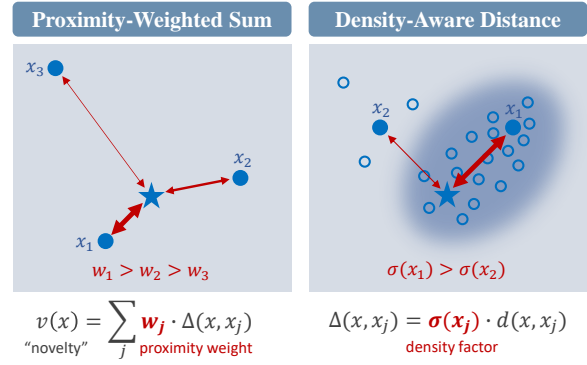"novelty"   proximity weight   density factor

Figure 3: *NovelSum* computes each sample's novelty as a proximity-weighted sum of its density-aware distances to other samples, where closer points have greater influence and high-density regions produce larger distances.

$\pi_i(j) = 1$ indicating that $x_j$ is the nearest neighbor of $x_i$. The function $\phi(\cdot)$ is monotonically decreasing, smoothing the weights according to the proximity, for example, we set $\phi(\pi_i(j)) = \frac{1}{\pi_i(j)}$. The hyperparameter $\alpha$ controls the degree to which proximity impacts the uniqueness score.

**Density-Aware Distance** To account for the local information density when calculating $\Delta(x_i, x_j)$, we introduce a density-aware distance that multiplies the original semantic distance by a density factor $\sigma(x_j)$:

$$\Delta(x_i, x_j) = \sigma(x_j)^\beta \cdot d(x_i, x_j) \qquad (9)$$

Since the probabilistic density of the overall sample distribution is intractable, we approximate the density factor by the inverse of the average distance to the $K$-nearest neighbors of $x_j$ in $\mathcal{X}^{all}$:

$$\sigma(x_j) = \frac{1}{\sum_{k=1}^{K} d(x_j, N_k(x_j))}$$

Here, $d(\cdot, \cdot)$ represents the distance between the embeddings of samples (e.g., cosine distance), and $N_k(x)$ denotes the $k$-th nearest neighbor of $x$. The hyperparameter $\beta$ controls the extent to which density influences the distance. The reference dataset $\mathcal{X}^{all}$ can be replaced to estimate information density under different sample distributions.

This approach mirrors how novelty is assessed in academic papers: a paper's novelty lies in its difference from closely related work, measured within the context of its field for greater accuracy. Accordingly, we treat each sample's quantified uniqueness as its "novelty" and name our method "NovelSum." Appendix E.3 provides a theoretical interpretation, while Appendix C analyzes the computational complexity and underscores *NovelSum*'s efficiency.
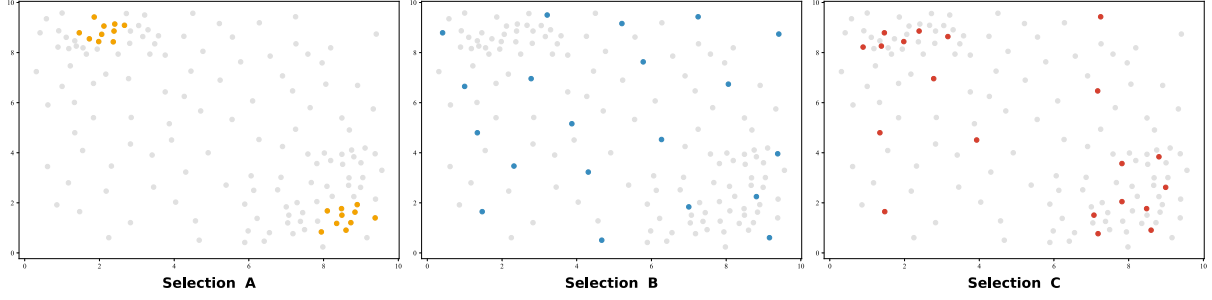
Figure 4: Simulating different data selection scenarios in a 2D sample space: Selection A represents datasets with redundancy, Selection B optimizes inter-sample distances, and Selection C accounts for both distances and density, which prior analysis suggests yields the highest diversity for instruction tuning.
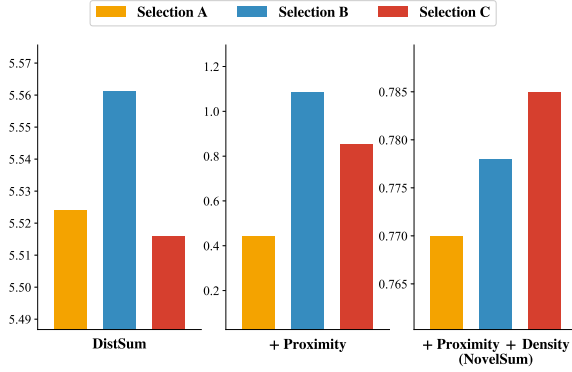


Figure 5: Measuring the diversity of simulated selection A/B/C with various metrics. *NovelSum* accurately captures dataset diversity, exhibiting expected behaviors.

## 4 Simulation Study

To validate whether the proposed metric aligns with our design principles and accurately captures dataset diversity, we create a visualizable simulation environment. We generate 150 points in 2D space as the data source and select 20 samples to form a dataset, simulating the data selection process for instruction tuning. As shown in Figure 4, we analyze three data selection scenarios to examine the behavior of our diversity metric. "Selection A" contains samples from two clusters, with most points close to each other, simulating datasets with redundancy. "Selection B", constructed using K-Center-Greedy, consists of samples far apart, simulating datasets optimized for inter-sample semantic distances. "Selection C" considers both inter-sample distances and information density, simulating datasets that best represent the sample space with unique points. Based on prior analysis, the dataset diversity of the three selections should follow $A < B < C$ order under the IT scenario.

Figure 5 presents the diversity measurement results using DistSum, a proximity-weighted version of DistSum, and *NovelSum*. From left to

right, we see that DistSum counterintuitively considers $\mathcal{M}(A) \simeq \mathcal{M}(C)$, failing to reflect sample uniqueness. Incorporating the proximity-weighted sum improves uniqueness capture but still exhibits $\mathcal{M}(B) > \mathcal{M}(C)$, overlooking information density. *NovelSum* **resolves these issues, accurately capturing diversity variations in alignment with design principles**, yielding $\mathcal{M}(A) < \mathcal{M}(B) < \mathcal{M}(C)$. This study further validates the necessity of the proximity-weighted sum and density-aware distance for precise diversity measurement.

## 5 Experiments

Following the settings in Section 2, we evaluate *NovelSum*'s correlation with the fine-tuned model performance across 53 IT datasets and compare it with previous diversity metrics. Additionally, we conduct a correlation analysis using Qwen-2.5-7B (Yang et al., 2024) as the backbone model, alongside previous LLaMA-3-8B experiments, to further demonstrate the metric's effectiveness across different scenarios. Due to resource constraints, we run each strategy on Qwen for at least two rounds, yielding 25 datasets.

### 5.1 Main Results

*NovelSum* **consistently achieves state-of-the-art correlation with model performance across various data selection strategies, backbone LLMs, and correlation measures.** Table 1 presents diversity measurement results on datasets constructed by mainstream data selection methods (based on $\mathcal{X}^{all}$), random selection from various sources, and duplicated samples (with only $m = 100$ unique samples). Results from multiple runs are averaged for each strategy. Although these strategies yield varying performance rankings across base models, *NovelSum* consistently tracks changes in model performance by accurately measuring dataset diver-

| Metrics | Data Selection Strategies | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | K-means | K-Center -Greedy | QDIT | Repr Filter | Random | | | | | Duplicate |
| | | | | | $\mathcal{X}^{all}$ | ShareGPT | WizardLM | Alpaca | Dolly | |
| *LLaMA-3-8B* | | | | | | | | | | |
| **Model Performance** | 1.32 | 1.31 | 1.25 | 1.05 | 1.20 | 0.83 | 0.72 | 0.07 | -0.14 | -1.35 |
| **NovelSum (Ours)** | 0.693 | 0.687 | 0.673 | 0.671 | 0.675 | 0.628 | 0.591 | 0.572 | 0.50 | 0.461 |
| Vendi Score $_{\times 10^7}$ | 1.70 | 2.53 | 1.59 | 2.23 | 1.61 | 1.70 | 1.44 | 1.32 | 1.44 | 0.05 |
| DistSum$_{cosine}$ | 0.648 | 0.746 | 0.629 | 0.703 | 0.634 | 0.656 | 0.578 | 0.605 | 0.603 | 0.634 |
| Facility Loc. $_{\times 10^5}$ | 2.99 | 2.73 | 2.99 | 2.86 | 2.99 | 2.83 | 2.88 | 2.83 | 2.59 | 2.52 |
| *Qwen-2.5-7B* | | | | | | | | | | |
| **Model Performance** | 1.06 | 1.45 | 1.23 | 1.35 | 0.87 | 0.07 | -0.08 | -0.38 | -0.49 | -0.43 |
| **NovelSum (Ours)** | 0.440 | 0.505 | 0.403 | 0.495 | 0.408 | 0.392 | 0.349 | 0.336 | 0.320 | 0.309 |
| Vendi Score $_{\times 10^6}$ | 1.60 | 3.09 | 2.60 | 7.15 | 1.41 | 3.36 | 2.65 | 1.89 | 3.04 | 0.20 |
| DistSum$_{cosine}$ | 0.260 | 0.440 | 0.223 | 0.421 | 0.230 | 0.285 | 0.211 | 0.189 | 0.221 | 0.243 |
| Facility Loc. $_{\times 10^5}$ | 3.54 | 3.42 | 3.54 | 3.46 | 3.54 | 3.51 | 3.50 | 3.50 | 3.46 | 3.48 |

Table 1: Comparison of fine-tuned model performance (Eq. 5) and diversity measurement results across datasets selected by different strategies. Each row visualizes the relative ranking of metric scores across datasets using color intensity: darker blue indicates higher values per row, and darker orange indicates lower values. *NovelSum* consistently shows a stronger correlation with model performance than other metrics, even as data selection strategies vary in performance between LLaMA-3-8B and Qwen-2.5-7B. More results are provided in Appendix E.4.

| Diversity Metrics | LLaMA | | | Qwen |
|---|---|---|---|---|
| | Pearson | Spearman | Avg. | Avg. |
| TTR | -0.38 | -0.16 | -0.27 | -0.30 |
| vocd-D | -0.43 | -0.17 | -0.30 | -0.31 |
| Facility Loc. | 0.86 | 0.69 | 0.77 | 0.08 |
| Entropy | 0.93 | 0.80 | 0.86 | 0.63 |
| LDD | 0.61 | 0.75 | 0.68 | 0.60 |
| KNN Distance | 0.59 | 0.80 | 0.70 | 0.67 |
| DistSum$_{cosine}$ | 0.85 | 0.67 | 0.76 | 0.51 |
| Vendi Score | 0.70 | 0.85 | 0.78 | 0.60 |
| DistSum$_{L2}$ | 0.86 | 0.76 | 0.81 | 0.51 |
| Cluster Inertia | 0.81 | 0.85 | 0.83 | 0.76 |
| Radius | 0.87 | 0.81 | 0.84 | 0.48 |
| NovelSum | **0.98** | **0.95** | **0.97** | **0.90** |

Table 2: Correlations between different metrics and model performance on LLaMA-3-8B and Qwen-2.5-7B. "Avg." denotes the average correlation (Eq. 6).

| Variants | Pearson | Spearman | Avg. |
|---|---|---|---|
| NovelSum | 0.98 | 0.95 | 0.97 |
| - Use $L2$ distance | 0.97 | 0.83 | $0.90_{\downarrow 0.08}$ |
| - $K = 20$ | 0.98 | 0.96 | $0.97_{\downarrow 0.00}$ |
| - $\alpha = 0$ (w/o proximity) | 0.79 | 0.31 | $0.55_{\downarrow 0.42}$ |
| - $\beta = 0$ (w/o density) | 0.92 | 0.89 | $0.91_{\downarrow 0.07}$ |

Table 3: Ablation Study for *NovelSum*.

sity. For instance, K-means achieves the best performance on LLaMA with the highest NovelSum score, while K-Center-Greedy excels on Qwen, also correlating with the highest NovelSum. Table 2 shows the correlation coefficients between various metrics and fine-tuned model performance for both LLaMA and Qwen experiments, where *NovelSum* achieves state-of-the-art correlation across different models and measures.

***NovelSum* can provide valuable guidance for data engineering practices.** As a reliable indicator of data diversity, *NovelSum* can assess diversity at both the dataset and sample levels, directly guiding data selection and construction decisions. For

example, Table 1 shows that the combined data source $\mathcal{X}^{all}$ is a better choice for sampling diverse IT data than other sources. Moreover, *NovelSum* can offer insights through comparative analyses, such as: (1) ShareGPT, which collects data from real internet users, exhibits greater diversity than Dolly, which relies on company employees, suggesting that IT samples from diverse sources enhance dataset diversity (Wang et al., 2024b); (2) In LLaMA experiments, random selection can outperform some mainstream strategies, aligning with prior work (Xia et al., 2024; Diddee and Ippolito, 2024), highlighting gaps in current data selection methods for optimizing diversity.

## 5.2 Ablation Study

*NovelSum* comprises several tunable components. In our main experiments, we use cosine distance to compute $d(x_i, x_j)$ in Eq. 9, with hyperparameters set to $\alpha = 1$, $\beta = 0.5$, and $K = 10$ nearest neighbors in Eq. 8 and Eq. 9. Here, we conduct an ablation study to investigate the impact of these settings based on LLaMA-3-8B.

In Table 3, $\alpha = 0$ removes the proximity

**Algorithm 1** *NovelSelect*

1: **Input:** Data pool $\mathcal{X}^{all}$, data budget $n$
2: Initialize an empty dataset, $\mathcal{X} \leftarrow \emptyset$
3: **while** $|\mathcal{X}| < n$ **do**
4: $\quad x^{new} \leftarrow \arg\max_{x \in \mathcal{X}^{all}} v(x)$
5: $\quad \mathcal{X} \leftarrow \mathcal{X} \cup \{x^{new}\}$
6: $\quad \mathcal{X}^{all} \leftarrow \mathcal{X}^{all} \setminus \{x^{new}\}$
7: **end while**
8: **return** $\mathcal{X}$

| Strategies | MT-bench | AlpacaEval | Aggregated $\mathcal{P}$ |
|---|---|---|---|
| Random | 6.18 | 75.47 | 1.20 |
| Repr Filter | 6.17 | 72.57 | 1.05 |
| QDIT | 6.21 | 75.91 | 1.25 |
| K-Center-Greedy | 6.33 | 75.30 | 1.31 |
| K-means | 6.33 | 75.46 | 1.32 |
| NovelSelect | **6.47** | **78.07** | **1.55** |

Table 4: Comparisons of different diversity-oriented data selection strategies on IT performance. $\mathcal{P}$ aggregates the performance based on Z-scores (Eq. 5).

weights, and $\beta = 0$ eliminates the density multiplier. We observe that both $\alpha = 0$ and $\beta = 0$ significantly weaken the correlation, validating the benefits of the proximity-weighted sum and density-aware distance. Replacing cosine distance with Euclidean distance and using more neighbors for density approximation have minimal impact, particularly on Pearson's correlation, highlighting *NovelSum*'s robustness to different distance measures. Additionally, Appendix E.1 presents an in-depth analysis of the hyperparameters, demonstrating the reliability of our current configurations and providing guidance for broader application.

## 6 Data Selection Strategy

**Introducing *NovelSelect*** Given *NovelSum*'s accurate diversity measurement and strong correlation with model performance, we investigate its potential as an optimization objective for selecting samples and generating a diverse dataset:

$$\mathcal{X} = \arg\max_{\mathcal{X} \subset \mathcal{X}^{all}} \mathcal{M}_{NovelSum}(\mathcal{X}), \quad (10)$$

where $\mathcal{M}_{NovelSum}(\mathcal{X})$ is defined in Eq. 7. Since directly solving Eq. 10 is NP-hard (Cook et al., 1994), we propose a greedy approach that iteratively selects the most "novel" sample to maximize the *NovelSum* of the final dataset. The "novelty" of a new sample $x$ relative to an existing set $\mathcal{X}$ is defined as:

$$v(x) = \sum_{x_j \in \mathcal{X}} w(x, x_j)^\alpha \cdot \sigma(x_j)^\beta \cdot d(x, x_j), \quad (11)$$

where $w(x, x_j)$ and $\sigma(x_j)$ are the proximity weight and density factor from Eq. 8 and 9. At each step, the sample with the highest novelty is selected: $x^{new} = \arg\max_x v(x), \quad \mathcal{X} \leftarrow x^{new} \cup \mathcal{X}$. This process is repeated from $\mathcal{X} = \emptyset$ until the data budget is reached, resulting in the selected dataset. We refer to this approach as *NovelSelect*.

Algorithm 1 outlines the overall process. Notably, *NovelSelect* is as computationally efficient as existing approaches, with a detailed analysis provided in Appendix C. Furthermore, by incorporating quality scores into $v(x)$, *NovelSelect* can seamlessly integrate with quality-based data selection methods, highlighting its extensibility.

**Data Selection Experiments** We conduct additional data selection experiments on LLaMA-3-8B to evaluate *NovelSelect*'s performance. Following prior settings, we use *NovelSelect* to select 10k samples from $\mathcal{X}^{all}$ and assess the fine-tuned model's performance on MT-bench and AlpacaEval. Results are averaged over three runs.

From Table 4, *NovelSelect* outperforms existing diversity-oriented data selection strategies on both benchmarks, demonstrating superior IT performance. This aligns with the higher *NovelSum* scores achieved by *NovelSelect* (Figure 1), further validating *NovelSum*'s effectiveness and practical value in IT data engineering.

## 7 Discussion

**From General IT to Downstream Tasks** Our study focuses on general instruction tuning, with training and evaluation covering a wide range of downstream tasks, thereby offering insights applicable to broader real-world scenarios. At the same time, extending performance-aligned diversity measurement to specific domains—such as math or code—is also valuable and may warrant dedicated investigation. A promising approach is to adapt *NovelSum* by simply replacing the general dataset $\mathcal{X}^{all}$ used in density estimation with domain-specific data sources. We hope our work lays a solid foundation for future research in this direction.

**Impact of Embedding Extractor** Different embedding extractors may yield different sample distributions in the semantic space, potentially affect-

ing diversity computations. In our study, we use LLaMA-3-8B to compute embeddings for experiments based on LLaMA-3-8B, and Qwen-2.5-7B for those based on Qwen-2.5-7B (see Appendix A.1 for embedding details). This setup is motivated by an interesting observation: extracting embeddings using the same model as the fine-tuning backbone yields metrics with the highest correlation to instruction tuning performance. The corresponding experiment is described in detail in Appendix E.2.

## 8 Related Work

**Measuring Dataset Diversity** Dataset diversity is essential for training generalizable machine learning models, drawing significant research interest (Yu et al., 2022; Sun et al., 2024; Zhang et al., 2024b; Zhao et al., 2024a; Qin et al., 2024). In NLP, numerous lexical diversity metrics have been proposed to measure text diversity through vocabulary usage (Richards, 1987; Malvern et al., 2004). Recently, semantic embeddings have enabled more flexible diversity measurement from distance (Stasaski and Hearst, 2022; Du and Black, 2019; Dang and Verma, 2024) or distribution perspectives (Shao et al., 2024). Focusing on instruction tuning, while some studies have explored the assessment of IT data diversity (Wang et al., 2024b; Bukharin et al., 2024), the proposed metrics lack sufficient validation of their correlation with IT performance; thus, reliable metrics for guiding data engineering remain underexplored.

**Data Selection for Instruction Tuning** Instruction tuning trains LLMs to follow human instructions using instruction-response pairs (Zhang et al., 2023). While earlier work focused on large-scale IT datasets (Longpre et al., 2023; Chiang et al., 2023), recent studies show that small, high-quality data sets can reduce costs and improve performance (Chen et al., 2023a,b; Zhou et al., 2024; Dou et al., 2024; Ye et al., 2024). This has led to the development of data selection strategies to identify subsets that boost IT performance (Liu et al., 2023; Du et al., 2023; Wu et al., 2023; Song et al., 2024; Ge et al., 2024; Kung et al., 2023; Yang et al., 2025). However, the lack of clear definitions and reliable diversity metrics for IT datasets hinders effective optimization. Consequently, some selection methods fail to generalize or perform worse than random selection (Xia et al., 2024; Diddee and Ippolito, 2024). Our work seeks to provide a more reliable diversity metric, based on comprehensive analysis, that accurately reflects the diversity of IT datasets and their instruction tuning performance.

## 9 Conclusion

In this paper, we investigate the fundamental problem of precisely measuring dataset diversity for instruction tuning and propose *NovelSum*, a reliable diversity metric that correlates well with model performance. Inspired by our systematic analysis of existing diversity metrics, *NovelSum* jointly considers inter-sample distances and information density to effectively capture dataset diversity, achieving superior correlations with model performance compared to previous metrics. Based on *NovelSum*, We further develop a data selection strategy, *NovelSelect*, whose remarkable performance validates the practical significance of *NovelSum*.

## Limitations

Although our work systematically analyzes both existing and proposed metrics through extensive fine-tuning experiments, we focus on Qwen-2.5-7B and LLaMA-3-8B as the backbone LLMs, excluding larger models and other series due to resource constraints. Moreover, while we strive to employ comprehensive benchmarks to evaluate instruction tuning performance, the test data may still fall short of fully capturing the diversity of real-world use cases. As a result, the beneficial effects of data diversity on model capabilities may be underrepresented in benchmark results. Finally, as previously noted, diversity measurements on downstream IT tasks may differ from our analysis in the general setting, suggesting the need for further study.

## Acknowledgements

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Alexander Bukharin, Shiyang Li, Zhengyang Wang, Jingfeng Yang, Bing Yin, Xian Li, Chao Zhang, Tuo Zhao, and Haoming Jiang. 2024. Data diversity matters for robust instruction tuning. In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 3411–3425. Association for Computational Linguistics.

Alfonso Cevallos, Friedrich Eisenbrand, and Rico Zenklusen. 2017. Local search for max-sum diversification. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 130–142. SIAM.

Hao Chen, Yiming Zhang, Qi Zhang, Hantao Yang, Xiaomeng Hu, Xuetao Ma, Yifan Yanggong, and Junbo Zhao. 2023a. Maybe only 0.5% data is needed: A preliminary exploration of low training data instruction tuning. *arXiv preprint arXiv:2305.09246*.

Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srinivasan, Tianyi Zhou, Heng Huang, et al. 2023b. Alpagasus: Training a better alpaca with fewer data. *arXiv preprint arXiv:2307.08701*.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6.

Israel Cohen, Yiteng Huang, Jingdong Chen, Jacob Benesty, Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. 2009. Pearson correlation coefficient. *Noise reduction in speech processing*, pages 1–4.

Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. Free dolly: Introducing the world's first truly open instruction-tuned llm.

William J Cook, William H Cunningham, William R Pulleyblank, and Alexander Schrijver. 1994. Combinatorial optimization. *Unpublished manuscript*, 10:75–93.

Vu Minh Hoang Dang and Rakesh M Verma. 2024. Data quality in nlp: Metrics and a comprehensive taxonomy. In *International Symposium on Intelligent Data Analysis*, pages 217–229. Springer.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Harshita Diddee and Daphne Ippolito. 2024. Chasing random: Instruction selection strategies fail to generalize. *arXiv preprint arXiv:2410.15225*.

Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. Enhancing chat language models by scaling high-quality instructional conversations. *arXiv preprint arXiv:2305.14233*.

Shihan Dou, Enyu Zhou, Yan Liu, Songyang Gao, Wei Shen, Limao Xiong, Yuhao Zhou, Xiao Wang, Zhiheng Xi, Xiaoran Fan, et al. 2024. Loramoe: Alleviating world knowledge forgetting in large language models via moe-style plugin. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1932–1945.

Qianlong Du, Chengqing Zong, and Jiajun Zhang. 2023. Mods: Model-oriented data selection for instruction tuning. *arXiv preprint arXiv:2311.15653*.

Wenchao Du and Alan W Black. 2019. Boosting dialog response generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. 2024. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*.

Yuval Eldar, Michael Lindenbaum, Moshe Porat, and Yehoshua Y Zeevi. 1997. The farthest point strategy for progressive image sampling. *IEEE transactions on image processing*, 6(9):1305–1315.

Reza Zanjirani Farahani and Masoud Hekmatfar. 2009. *Facility location: concepts, models, algorithms and case studies*. Springer Science & Business Media.

Sándor P Fekete and Henk Meijer. 2004. Maximum dispersion and geometric maximum weight cliques. *Algorithmica*, 38:501–511.

Yuan Ge, Yilun Liu, Chi Hu, Weibin Meng, Shimin Tao, Xiaofeng Zhao, Hongxia Ma, Li Zhang, Boxing Chen, Hao Yang, et al. 2024. Clustering and ranking: Diversity-preserved instruction selection through expert-aligned quality estimation. *arXiv preprint arXiv:2402.18191*.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.

Po-Nien Kung, Fan Yin, Di Wu, Kai-Wei Chang, and Nanyun Peng. 2023. Active instruction tuning: Improving cross-task generalization by training on prompt sensitive tasks. *arXiv preprint arXiv:2311.00288*.

Yi-An Lai, Xuan Zhu, Yi Zhang, and Mona Diab. 2020. Diversity, density, and homogeneity: Quantitative characteristic metrics for text collections. *arXiv preprint arXiv:2003.08529*.

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023. Alpacaeval: An automatic evaluator of instruction-following models.

Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. 2023. What makes good data for alignment? a comprehensive study of automatic data selection in instruction tuning. *arXiv preprint arXiv:2312.15685*.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR.

James MacQueen. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, volume 5, pages 281–298. University of California press.

David Malvern, Brian Richards, Ngoni Chipere, and Pilar Durán. 2004. *Lexical diversity and language development*. Springer.

Amey P Pasarkar and Adji Bousso Dieng. 2023. Cousins of the vendi score: A family of similarity-based diversity metrics for science and machine learning. *arXiv preprint arXiv:2310.12952*.

Yulei Qin, Yuncheng Yang, Pengcheng Guo, Gang Li, Hang Shao, Yuchen Shi, Zihan Xu, Yun Gu, Ke Li, and Xing Sun. 2024. Unleashing the power of data tsunami: A comprehensive survey on data assessment and selection for instruction tuning of language models. *arXiv preprint arXiv:2408.02085*.

Sekharipuram S Ravi, Daniel J Rosenkrantz, and Giri Kumar Tayi. 1994. Heuristic and special case algorithms for dispersion problems. *Operations research*, 42(2):299–310.

Brian Richards. 1987. Type/token ratios: What do they really tell us? *Journal of child language*, 14(2):201–209.

Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. 2022. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*.

Ozan Sener and Silvio Savarese. 2017. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*.

Yunfan Shao, Linyang Li, Zhaoye Fei, Hang Yan, Dahua Lin, and Xipeng Qiu. 2024. Balanced data sampling for language model training with clustering. *arXiv preprint arXiv:2402.14526*.

Jielin Song, Siyu Liu, Bin Zhu, and Yanghui Rao. 2024. Iterselecttune: An iterative training framework for efficient instruction-tuning data selection. *arXiv preprint arXiv:2410.13464*.

Katherine Stasaski and Marti A Hearst. 2022. Semantic diversity in dialogue with natural language inference. *arXiv preprint arXiv:2205.01497*.

Katherine Stasaski, Grace Hui Yang, and Marti A Hearst. 2020. More diverse dialogue datasets via diversity-informed data collection. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 4958–4968.

Peng Sun, Bei Shi, Daiwei Yu, and Tao Lin. 2024. On the diversity and realism of distilled dataset: An efficient dataset distillation paradigm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9390–9399.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Jiahao Wang, Bolin Zhang, Qianlong Du, Jiajun Zhang, and Dianhui Chu. 2024a. A survey on data selection for llm instruction tuning. *arXiv preprint arXiv:2402.05123*.

Peiqi Wang, Yikang Shen, Zhen Guo, Matthew Stallone, Yoon Kim, Polina Golland, and Rameswar Panda. 2024b. Diversity measurement and subset selection for instruction tuning datasets. *arXiv preprint arXiv:2402.02318*.

Shengguang Wu, Keming Lu, Benfeng Xu, Junyang Lin, Qi Su, and Chang Zhou. 2023. Self-evolved diverse data sampling for efficient instruction tuning. *arXiv preprint arXiv:2311.08182*.

Tingyu Xia, Bowen Yu, Kai Dang, An Yang, Yuan Wu, Yuan Tian, Yi Chang, and Junyang Lin. 2024. Rethinking data selection at scale: Random selection is almost all you need. *arXiv preprint arXiv:2410.09335*.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. Wizardlm: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Yuming Yang, Wantong Zhao, Caishuang Huang, Junjie Ye, Xiao Wang, Huiyuan Zheng, Yang Nan, Yuran Wang, Xueying Xu, Kaixin Huang, Yunke Zhang, Tao Gui, Qi Zhang, and Xuanjing Huang. 2025. Beyond boundaries: Learning a universal entity taxonomy across datasets and languages for open named entity recognition. In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 10902–10923. Association for Computational Linguistics.

Junjie Ye, Yuming Yang, Qi Zhang, Tao Gui, Xuanjing Huang, Peng Wang, Zhongchao Shi, and Jianping Fan. 2024. Empirical insights on fine-tuning large language models for question-answering. *CoRR*, abs/2409.15825.

Yu Yu, Shahram Khadivi, and Jia Xu. 2022. Can data diversity enhance learning generalization? In *Proceedings of the 29th international conference on computational linguistics*, pages 4933–4945.

Jerrold H Zar. 2005. Spearman rank correlation. *Encyclopedia of biostatistics*, 7.

Dylan Zhang, Justin Wang, and Francois Charton. 2024a. Instruction diversity drives generalization to unseen tasks. *arXiv preprint arXiv:2402.10891*.

Dylan Zhang, Justin Wang, and Francois Charton. 2024b. **Only-IF**: Revealing the decisive effect of instruction diversity on generalization. *arXiv preprint arXiv:2410.04717*.

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, et al. 2023. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*.

Dorothy Zhao, Jerone TA Andrews, AI Sony, Tokyo Orestis Papakyriakopoulos, and Alice Xiang. 2024a. Measuring diversity in datasets. In *International Conference on Learning Representations*.

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

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

Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.

# A Details of Correlation Evaluation

## A.1 Data Processing and Semantic Embeddings

We apply basic preprocessing to remove anomalous samples from the data sources, ensuring more stable results while preserving generality. In the early stage of our work, we observe that short samples often exhibit low quality and tend to be outliers in the semantic space, potentially distorting experimental results. To address this, we filter out samples shorter than 256 tokens using the BERT (Devlin et al., 2019) tokenizer, ensuring consistency for experiments across different LLMs. Furthermore, to ensure the dataset's relevance for English-language tasks and math problems, we exclude samples with a non-English-or-number ratio exceeding 0.8.

When computing sample embeddings, we set the maximum sequence length to 256 to mitigate bias from varying text lengths. This applies only to embedding computation; fine-tuning uses a much larger maximum length. We extract the last hidden layer of the language model and apply mean pooling, excluding padding tokens, to generate robust sample-level embeddings. We analyze and discuss the choice of embedding extractors in Appendix E.2. Based on this analysis, we use LLaMA-3-8B to compute embeddings in experiments using LLaMA-3-8B as the fine-tuning backbone. Similarly, for experiments using Qwen-2.5-7B as the fine-tuning backbone, we use Qwen-2.5-7B to compute embeddings.

## A.2 Details of Existing Diversity Metrics

For lexical diversity, the **Type-Token Ratio** (TTR) quantifies the lexical diversity of a text sequence $x_i$ as the ratio of distinct tokens to the total number of tokens. The overall lexical diversity of a dataset $\mathcal{X} = \{x_1, x_2, ..., x_N\}$ is computed as the average TTR across all samples:

$$\mathcal{M}_{TTR}(\mathcal{X}) = \frac{1}{N} \sum_{i=1}^{N} \frac{|Unique(x_i)|}{|x_i|}. \quad (12)$$

To mitigate the influence of text length on TTR, we randomly sample 30 tokens from each data point to compute the TTR.

To address the sensitivity of TTR to text length, **vocd-D** extends this measure by computing $TTR_i^k$ over sampled sub-sequences of varying lengths $k$ and fitting the following curve:

$$T\hat{T}R_i^k = \frac{D}{k} \left( (1 + 2\frac{k}{D})^{\frac{1}{2}} - 1 \right), \quad (13)$$

where $D$ is the estimated parameter representing lexical diversity. The vocd-D metric is defined as $\mathcal{M}_{vocd-D} = D_{\text{best fit}}$, with larger values indicating greater lexical diversity. In our experiments, we compute $TTR_i^k$ for $k = 10, 20, 30, 40, 50$ and take the average of the resulting values as the final lexical diversity score.

For distance-based semantic diversity, **Cluster Inertia** (Du and Black, 2019) quantifies diversity by partitioning the dataset into $K$ clusters using K-means and summing the squared distances between each sample and its cluster centroid:

$$\mathcal{M}_{Inertia}(\mathcal{X}) = \sum_{j=1}^{K} \sum_{x_i \in C_j} \|emb(x_i) - \mu_j\|^2, \quad (14)$$

where $\mu_j$ is the centroid of cluster $C_j$. A higher inertia value suggests a greater spread of samples. Additionally, **Vendi Score** (VS) (Pasarkar and Dieng, 2023) measures diversity based on the eigenvalues of the similarity kernel matrix. The generalized VS metric is defined as:

$$\mathcal{M}_{VS}(\mathcal{X}) = \exp \left( \frac{1}{1-\alpha} \log_2 \sum_{i=1}^{|\mathcal{X}|} \bar{\lambda}_{i|\theta}^{\alpha} \right), \quad (15)$$

where $\bar{\lambda}_{i|\theta}$ represents the normalized eigenvalues. We set $\alpha = 0.5$ to enhance measurement under severe class imbalance. **Radius** (Lai et al., 2020) characterizes the dispersion of the sample space by approximating embeddings as a multi-variate Gaussian distribution. It computes the geometric mean of the standard deviations along each dimension:

$$\mathcal{M}_{Radius}(\mathcal{X}) = \sqrt[H]{\prod_{j=1}^{H} \sigma_j}, \quad (16)$$

where $H$ is the embedding dimension, and $\sigma_j$ denotes the radius of the ellipsoid along the $j$-th axis. Larger values indicate a greater spread of samples in the embedding space. **Log Determinant Distance** (Wang et al., 2024b) utilizes the determinant of the similarity matrix as a measure of dataset diversity. In our work, we employ the cosine similarity function to compute the similarity matrix.

Note that for **DistSum**$_{cosine}$, we use cosine distance $\Delta(x_i, x_j) = 1 - \cos(emb(x_i), emb(x_j))$. For **DistSum**$_{L2}$, we use Euclidean distance $\Delta(x_i, x_j) = \|emb(x_i) - emb(x_j)\|_2^2$.

For **Partition Entropy**, we cluster $\mathcal{X}^{all}$ into 1,000 clusters, while for **Cluster Inertia** (Du and Black, 2019), we cluster $\mathcal{X}^s$ into 200 clusters for subsequent computations.

## A.3 Details of Data Selection Strategies

All IT datasets in our experiments are selected from $\mathcal{X}^{all}$ and sampled over three rounds (two for Qwen) per strategy variant, unless stated otherwise. We assume these datasets have similar average sample quality, as they come from the same source without any quality filters. Additionally, the dataset size is standardized to 10,000 samples. Thus, our experiments can more accurately reflect the correlation between dataset diversity and model performance, without introducing significant confounders.

**K-Center-Greedy** (Sener and Savarese, 2017; Chen et al., 2023a; Du et al., 2023; Wu et al., 2023) This strategy begins by randomly selecting a data point from the dataset $\mathcal{X}^{all}$ as the initial point of the subset $\mathcal{X}^{(s)}$. Subsequently, it iteratively computes the closest distance between the remaining points in $\mathcal{X}^{all} \setminus \mathcal{X}^{(s)}$ and selected samples in $\mathcal{X}^{(s)}$. The point with the maximum minimum distance (i.e., the farthest point) is added to $\mathcal{X}^{(s)}$. This process continues until the desired subset size is achieved.

**Repr Filter** (Liu et al., 2023) Unlike the K-Center-Greedy strategy, which selects the farthest point from the remaining data pool, the Repr Filter randomly selects a data point whose similarity with all embeddings in $\mathcal{X}^{(s)}$ is below a predefined threshold. Due to the unique distribution of embeddings across different models, it is necessary to set distinct thresholds for each similarity function and model embedding. To ensure diversity across different experimental rounds, we employ cosine similarity and set the threshold to 0.3 for LLaMA-3-8B and 0.1 for Qwen-2.5-7B.

**QDIT** (Bukharin et al., 2024) QDIT sampling combines diversity and quality scores for data selection; however, in our work, we focus exclusively on its diversity score. This method computes the sum of similarities between each sample in $\mathcal{X}^{all} \setminus \mathcal{X}^{(s)}$ and its closest data point in $\mathcal{X}^{(s)}$. For each candidate data point, we calculate the similarity sum as if it were added to $\mathcal{X}^{(s)}$, defining its Facility Location (FL) score. The algorithm then iteratively selects the data point with the highest FL score. For the initial selection, it chooses the data point that exhibits the highest overall similarity to all other embeddings. In our experiments, we employ cosine similarity for computing these scores. Since the Facility Location function yields a fixed subset $\mathcal{X}^{(s)}$ for a given $\mathcal{X}^{all}$, and to maintain consistency with other strategies, we utilize the same subset of data but vary the training random seeds across three rounds of experiments.

**K-means Clustering** (Song et al., 2024; Ge et al., 2024) For this strategy, we apply the K-means clustering algorithm (MacQueen, 1967) to partition all sample embeddings in $\mathcal{X}^{all}$ into $K$ clusters. Subsequently, given a target data budget $n$, we randomly sample $\frac{n}{K}$ data points from each cluster. For our experiments, we use both 1000 and 100 clusters for LLaMA-3-8B, and 100 clusters for Qwen-2.5-7B.

**Random Selection** In this baseline strategy, we randomly sample 10,000 data points from $\mathcal{X}^{all}$. To explore the impact of data sources, we also sample from individual datasets, including Alpaca (Taori et al., 2023), Dolly (Conover et al., 2023), WizardLM, UltraChat, and ShareGPT, with similar preprocessing. Although we assume that the average sample quality of these sources does not significantly differ from that of $\mathcal{X}^{all}$, we use only a single round of results from each source in the overall correlation analysis as supplementary data to avoid potential quality differences affecting the outcome.

**Duplicate Selection** To address the challenge of defining low-diversity datasets, which is crucial for our study, we construct datasets with redundant samples. Given a target data budget $n$, the dataset is constructed by selecting $m$ unique data points, each duplicated $\frac{n}{m}$ times. We set $m$ to 1, 10, 50, 100, 500, 1000, 2000, and 5000. This approach allows us to systematically control and analyze the impact of diversity on model performance.

## A.4 Details of Performance Evaluation

We follow prior work in adopting LLM-based judges—MT-Bench and AlpacaEval—as they demonstrate strong alignment with human preferences and offer broad coverage of downstream tasks (Zhou et al., 2024; Zhao et al., 2024b). AlpacaEval evaluates single-turn dialogue ability through pairwise preference comparisons between model responses and those of a strong baseline, judged by GPT-4 (Achiam et al., 2023). MT-Bench, by contrast, assesses multi-turn conversational ability via GPT-4-based evaluations. Together, these benchmarks cover a wide range of diverse user queries and representative IT tasks, including mathematics and code generation. In our evaluation, we use GPT-4-0613 as the judge for both MT-Bench and AlpacaEval. For AlpacaEval, we follow the original setup to adopt text-davinci-003 responses

as the baseline.

To account for the length bias inherent in LLM-based evaluation, we adopt the Length-Controlled Win Rates metric (Dubois et al., 2024) for Alpaca-Eval, which has demonstrated stronger alignment with human judgments. We also verify that the average response lengths across the evaluated instruction-tuned (IT) models show minimal variation. This consistency is achieved by controlling for length when sampling from each model's IT training dataset—for example, by using a fixed context length when computing embeddings and removing quality filters that might favor longer samples. Based on these measures, we believe our evaluation methodology offers a more reliable assessment of model performance.

Conversely, we didn't incorporate additional evaluation methods, such as multiple-choice QA benchmarks, because they may introduce biases toward specific domains and may not align well with human preferences. Thus, while we acknowledge the inherent limitations of LLM-based evaluation, it appears to remain the most accepted method for fairly evaluating open-ended responses in the absence of better alternatives.

### A.5 Details of Correlation Measures

We compute the correlation between each diversity metric and model performance using both Pearson (Cohen et al., 2009) and Spearman (Zar, 2005) correlation measures. For example, Pearson's $r$ for a metric $\mathcal{M}_t$ is computed as:

$$r_{\mathcal{M}_t,\,\mathcal{P}}^{Pearson} = \frac{\sum_s (\mathcal{M}_t^{(s)} - \bar{\mathcal{M}}_t)(\mathcal{P}^{(s)} - \bar{\mathcal{P}})}{\sigma_{\mathcal{M}_t}\sigma_{\mathcal{P}}} \quad (17)$$

### A.6 Details of Model Fine-Tuning

In our experiments, we leverage four or eight NVIDIA H800 GPUs for training the LLaMA-3-8B and Qwen-2.5-7B models. To enable efficient parallel training, we implement DeepSpeed Zero-Stage 2. Across all experiments conducted in this study, the training parameters are configured as follows: a maximum input length of 4096 tokens, a batch size of 128, 3 training epochs, a learning rate of 2e-5, and a warm-up ratio of 0.1 utilizing cosine warm-up. We use the official chat templates of LLaMA-3 and Qwen-2.5, respectively, to fine-tune each model. All models are trained with BF16 precision to optimize computational efficiency and memory usage. A single run of fine-tuning on a 10k dataset typically takes about one hour.

## B Implementation Details

### B.1 Implementation Details of *NovelSum*

As described in Section 3, our approach to computing data diversity incorporates both proximity-weighted and density-aware considerations. In practice, we begin by embedding the samples in the given dataset $\mathcal{X}^{(s)}$ as vectors and computing a similarity matrix that captures pairwise distances. We then apply proximity and density weights to achieve the desired outcomes.

To estimate sample density, we utilize FAISS (Johnson et al., 2019), which efficiently leverages GPU capabilities for vector similarity searches. Specifically, for each sample in $\mathcal{X}^{(s)}$, we identify its $k = 10$ nearest neighbors within the overall sample space $\mathcal{X}^{all}$ to compute its density factor, which we then broadcast to match the dimensions of the similarity matrix. Next, we perform element-wise multiplication between the density matrix and the similarity matrix to obtain density-aware distances in the embedding space.

Subsequently, we sort each row of the resulting matrix to determine the proximity ranks of all samples relative to the corresponding sample in that row. Finally, we compute the proximity-weighted sum for each row to derive each sample's "novelty" score and sum these scores to obtain $\mathcal{M}_{NovelSum}(\mathcal{X})$.

As noted earlier, we set the hyperparameters to $\alpha = 1$, $\beta = 0.5$, and $K = 10$ for experiments with both LLaMA-3-8B and Qwen-2.5-7B.

### B.2 Implementation Details of *NovelSelect*

Since selecting a subset from $\mathcal{X}^{all}$ that maximizes *NovelSum* is an NP-Hard problem, similar to selecting a subset with maximum Euclidean distance, we implement a greedy strategy (Section 6).

In our implementation, we iteratively compute the sample-level "novelty" $v(x)$ (Eq. 11) for each unselected candidate point with respect to the currently selected set, following the same computation process as *NovelSum*. At each step, the candidate with the highest $v(x)$ is added to the subset. Notably, the density factor used for distance computation is $\sigma(x_j)$—that of the selected point—rather than the candidate's own $\sigma(x)$, to remain consistent with *NovelSum*'s definition. That said, as an alternative greedy strategy, replacing $\sigma(x_j)$ with $\sigma(x_j) + \sigma(x)$ to jointly account for both samples' densities may also be reasonable, and we leave this for future exploration.

## C Computational Complexity

In practice, both *NovelSum* and *NovelSelect* incur acceptable computational costs—approximately 10 seconds and under one hour, respectively—relative to the overall fine-tuning process, and are comparable to or more efficient than many existing methods. Crucially, our approaches avoid pairwise computations over the entire large-scale source dataset of size $N$, operating instead on the selected subset, which is typically small (e.g., $n = 10,000$ in our experiments). And for density estimation, which considers the distribution of source samples, we leverage modern vector search libraries such as FAISS (Appendix B). FAISS supports near-constant-time nearest neighbor queries independent of $N$, with only a one-time $O(N)$ cost to index all source samples—both negligible in the overall computation.

*NovelSum* has a time complexity of $O(n^2)$ as it computes pairwise distances among the selected samples. This is as efficient as most existing embedding-based diversity measures. For density estimation that accounts for source sample distribution, we use FAISS, which incurs an approximate cost of $O(N)$ for indexing all source samples and $O(n)$ for querying the selected samples' density factor—both negligible in the overall computation. In practice, computing *NovelSum* (including density precomputation) on 10k samples with 4096-dimensional embeddings takes only 10 seconds on a single H800 GPU. Additionally, the one-time cost of building the FAISS index for 396k source samples is also under 10 seconds.

For our data selection strategy, *NovelSelect* runs in $O(N \cdot n^2)$ time—significantly more efficient than QDIT's $O(N^3)$. This involves distance computation between candidate samples (of size $N$) and selected samples (of size $n$) across $n$ selection iterations. In practice, selecting 10k samples from a 396k-sample pool takes under one hour using a single H800 GPU, which is faster than fine-tuning on 10k samples and negligible compared to fine-tuning on the full 396k dataset.

Embedding extraction, a shared step across embedding-based methods, takes under two hours for 396k samples in $\mathcal{X}^{all}$ using LLaMA-3-8B and vLLM on 8×H800 GPUs. As a one-time cost, this remains acceptable.

## D Data Statistics

Our data sources are detailed in Table 5. After filtering out short data and non-English data, ap-

| Data Pool | Dataset Source | Sample Size |
|---|---|---|
| $\mathcal{X}^{all}$ | ShareGPT | 103 K |
| | UltraChat | 207 K |
| | WizardLM | 196 K |
| $\mathcal{X}^{other}$ | Alpaca | 52 K |
| | Dolly | 15 K |

Table 5: Statistics of Data Pools $\mathcal{X}^{all}$ and $\mathcal{X}^{other}$. The column "Dataset Source" indicates the origin of the data used for sampling, while "Sample Size" denotes the number of samples in each dataset. This table provides an overview of the data used in our experiments.

proximately 396K samples remain in $\mathcal{X}^{all}$ for use in our experiments. Note that we use the latest versions of these datasets, which may have a larger size than the initial versions. These datasets encompass samples from a wide range of domains.

## E More Results and Analysis

### E.1 Hyperparameter Analysis

We conduct a more fine-grained hyperparameter analysis to study the effects of varying $\alpha$ and $\beta$, and investigate the sensitivity to hyperparameters for potential broader application of *NovelSum*. The results are shown in Table 6:

| Variants | LLaMA-3-8B | Qwen-2.5-7B |
|---|---|---|
| NovelSum ($\alpha = 1$, $\beta = 0.5$) | **0.97** | **0.90** |
| - $\alpha = 0$ | 0.55 | 0.51 |
| - $\alpha = 0.5$ | 0.77 | 0.64 |
| - $\alpha = 0.8$ | 0.91 | 0.85 |
| - $\alpha = 0.9$ | 0.94 | 0.88 |
| - $\alpha = 1.1$ | 0.95 | **0.91** |
| - $\alpha = 1.2$ | 0.93 | 0.89 |
| - $\alpha = 1.5$ | 0.86 | 0.86 |
| - $\alpha = 2$ | 0.81 | 0.82 |
| - $\beta = 0$ | 0.91 | 0.73 |
| - $\beta = 0.2$ | 0.93 | 0.80 |
| - $\beta = 0.3$ | 0.94 | 0.83 |
| - $\beta = 0.4$ | 0.94 | 0.86 |
| - $\beta = 0.6$ | 0.94 | **0.91** |
| - $\beta = 0.7$ | 0.92 | 0.82 |
| - $\beta = 0.8$ | 0.88 | 0.69 |
| - $\beta = 1$ | 0.76 | 0.37 |

Table 6: Hyperparameter analysis of *NovelSum* with varying $\alpha$ and $\beta$ configurations on LLaMA-3-8B and Qwen-2.5-7B.

These results show that NovelSum consistently achieves strong correlation with model performance across a relatively wide range of hyperparameters, without drastic fluctuations. This suggests that the sensitivity issue may not be particularly severe in practice.

From a theoretical perspective, we view the proximity decay coefficient $\alpha$ as related to the semantic richness of the source data. For richer dataset, it's better to consider more neighbors in distance computations, corresponding to a smaller $\alpha$. Given that most current IT tasks and datasets are semantically rich, the current choice of $\alpha$ is likely to remain effective as long as the domain is not overly narrow. On the other hand, the density coefficient $\beta$ controls the balance between distance and density components. We believe this balance is not specific to a particular dataset, but rather general across IT tasks. While the exact optimal value of $\beta$ may vary slightly depending on the implementation of distance and density calculations, the use of cosine distance and nearest-neighbor density estimation—as adopted in our work—provides a stable basis. Therefore, re-tuning $\beta$ is unlikely to be necessary in most cases.

Based on the above discussion, we believe our current hyperparameter configuration is robust for general instruction tuning and can exhibit a considerable degree of generalizability across broader scenarios. This helps reduce the need for costly hyperparameter tuning. Therefore, the level of hyperparameter sensitivity observed here may not be a major obstacle to the broader applicability of our method.

### E.2 Analysis of Embedding Extractors

To investigate the impact of the embedding model choice on diversity measurements, we conduct an additional ablation study using four different models to generate embeddings for *NovelSum* computation: LLaMA-3-8B, LLaMA-2-7B, Qwen-2.5-7B, and BERT-base. We then measure the correlation between the resulting *NovelSum* scores and instruction tuning (model) performance under two fine-tuning backbones: LLaMA-3-8B and Qwen-2.5-7B.

| Embedding Model | Fine-tuning Backbone | |
| --- | --- | --- |
| | LLaMA-3-8B | Qwen-2.5-7B |
| LLaMA-3-8B | **0.97** | 0.81 |
| Qwen-2.5-7B | 0.92 | **0.90** |
| LLaMA-2-7B | 0.94 | 0.87 |
| BERT-base | 0.90 | 0.64 |

Table 7: Correlation between *NovelSum* computed using different embedding models and instruction tuning performance under two fine-tuning backbones.

The results, presented in Table 7, indicate that using the same model for both embedding extraction and fine-tuning yields the strongest correlation between diversity metrics and instruction tuning performance, likely due to shared representation spaces. In contrast, employing a general-purpose encoder such as BERT-base leads to weaker correlations compared to other LLMs.

Following prior work, we use pretrained base models directly (Appendix A.1) for sample embedding computation, primarily for research purposes. For practical applications, however, one may consider using state-of-the-art LLM-based embedding models fine-tuned specifically for embedding tasks, which may offer improved performance.

### E.3 Theoretical Analysis

We begin by situating our work within broader literatures. The problem of selecting representative objects from a given set has been extensively studied in Operations Research (Ravi et al., 1994; Fekete and Meijer, 2004; Cevallos et al., 2017), often through formulations such as maximum dispersion and facility location. These approaches share similar motivations with our method and help explain the effectiveness of the density-aware distance. In parallel, prior work on sampling strategies (Eldar et al., 1997) conceptualizes sampling as a stochastic process for reconstruction and highlights the effectiveness of maximizing inter-sample distances in progressive image sampling. Building on similar insights, our diversity metric *NovelSum* accounts for both inter-sample distances and information density in the sample space.

To facilitate a deeper understanding of *NovelSum*, we offer the following theoretical interpretation as a possible perspective: The basic sum of semantic distances among all samples (Eq. 1) represents the semantic variance of the selected samples. A larger variance implies significant differences among some samples but does not necessarily guarantee a diverse semantic distribution across the entire semantic space. In contrast, the **proximity-weighted sum** (Eq. 8) specifically measures semantic diversity by focusing on the gaps between neighboring samples rather than on distant ones. A larger proximity-weighted sum indicates that samples are more distinct from their immediate neighbors, reflecting higher overall semantic diversity. Further, the **density-aware distance** (Eq. 9) incorporates an additional density factor into the distance calculation, explicitly considering information density within the semantic space.

By extracting information from the gaps between neighboring samples—multiplying the semantic distance between each sample and its neighbors by information density factors and performing a proximity-weighted sum—we effectively quantify each sample's unique informational contribution. Consequently, *NovelSum* measures the total unique information of all samples, given their semantic embeddings and scenario-specific information density (e.g., general instruction tuning). This aggregate value corresponds to the "IT-aligned Diversity" we aim to measure.

### E.4 Additional Results

Additional scatter plots for the analysis in Section 2 are provided in Figure 6, Figure 7 and Figure 8 , illustrating the correlation for $DistSum_{L2}$, Radius, and Log Determinant Distance, respectively.

The full results of the correlation experiments on LLaMA-3-8B and Qwen-2.5-7B are presented in Table 8 and Table 9, respectively. These tables provide a comprehensive comparison of diversity metrics across different experimental configurations.

## F Others

### F.1 License for Artifacts and Data Consent

In this paper, the artifacts used are all available for academic research work, including ShareGPT, WizardLM, UltraChat, Alpaca and Dolly. The diversity metrics and data selection methods compared in this paper can all be used for academic research. All data originates from the original authors' open-source releases and can be used for academic research and publication.

### F.2 Data Statement

The training datasets may contain offensive content; however, they do not include any personal information. Furthermore, our training approach is designed to align the model with human preferences without producing harmful content.

### F.3 AI Assistant Usage Statement

We utilized ChatGPT for writing refinement and minor coding assistance. AI assistants were not employed for research innovation, and all core contributions were solely developed by the authors.

### F.4 Budgets

We instruction-tune (train) each model for approximately one hour on a single node with eight H800-80G GPUs, totaling around 80 hours across 80 runs.

Additionally, we spend around $1,000 on the GPT API to evaluate our models using MT-bench and AlpacaEval.

| Data Selection Strategy | Runs | Alpaca | MT - bench | Aggregated | NovelSum | DistSum$_{cosine}$ | DistSum$_{L2}$ | KNN | Inertia | Radius $10^{-2}$ | VS $10^7$ | Entropy | FL $10^5$ | LDD $10^4$ | TTR | vocd - D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Random - alpaca | 1 | 58.8 | 4.40 | 0.07 | 0.572 | 0.605 | 1.09 | 0.509 | 0.293 | 1.13 | 1.32 | 8.93 | 2.83 | -1.70 | 0.853 | 82.3 |
| Random - dolly | 1 | 47.5 | 4.34 | -0.14 | 0.500 | 0.603 | 1.09 | 0.596 | 0.362 | 1.12 | 1.70 | 7.83 | 2.59 | -1.31 | 0.844 | 76.6 |
| Random - sharegpt | 1 | 74.6 | 6.46 | 0.83 | 0.628 | 0.656 | 1.14 | 0.574 | 0.380 | 1.16 | 1.70 | 8.97 | 2.83 | -1.23 | 0.850 | 78.3 |
| Random - ultrachat | 1 | 72.5 | 6.76 | 0.86 | 0.672 | 0.622 | 1.11 | 0.567 | 0.323 | 1.15 | 1.48 | 9.40 | 2.96 | -1.52 | 0.881 | 110 |
| Random - wizardlm | 1 | 76.9 | 5.82 | 0.72 | 0.591 | 0.578 | 1.07 | 0.594 | 0.331 | 1.11 | 1.44 | 9.08 | 2.88 | -1.52 | 0.858 | 85.7 |
| Random - $\mathcal{X}^{all}$ | 3 | 75.5 | 6.18 | 1.20 | 0.675 | 0.634 | 1.12 | 0.606 | 0.353 | 1.15 | 1.61 | 9.80 | 2.99 | -1.38 | 0.870 | 97.2 |
| Farthest | 3 | 74.0 | 6.30 | 1.22 | 0.687 | 0.789 | 1.25 | 0.407 | 0.350 | 1.20 | 1.56 | 6.52 | 2.14 | -1.25 | 0.837 | 68.3 |
| Duplicate m = 1 | 3 | 0.57 | 1.01 | -6.36 | 0.000 | 0.000 | 0.00 | 0.000 | 0.000 | 0.00 | 1.08 | 0.00 | 1.25 | -inf | 0.887 | 121 |
| Duplicate m = 10 | 3 | 31.2 | 3.54 | -2.97 | 0.268 | 0.589 | 1.02 | 0.000 | 0.000 | 1.03 | 7.16 | 3.27 | 2.08 | -inf | 0.863 | 90.0 |
| Duplicate m = 50 | 3 | 51.0 | 4.38 | -1.35 | 0.388 | 0.608 | 1.09 | 0.001 | 0.000 | 1.12 | 2.74 | 5.58 | 2.40 | -inf | 0.873 | 101 |
| Duplicate m = 100 | 3 | 63.6 | 5.42 | 0.05 | 0.461 | 0.634 | 1.12 | 0.001 | 0.000 | 1.15 | 4.95 | 6.50 | 2.52 | -inf | 0.869 | 92.3 |
| Duplicate m = 500 | 3 | 65.0 | 5.67 | 0.30 | 0.556 | 0.635 | 1.12 | 0.001 | 0.222 | 1.15 | 1.79 | 8.47 | 2.75 | -inf | 0.869 | 96.7 |
| Duplicate m = 1000 | 3 | 71.3 | 6.00 | 0.86 | 0.587 | 0.630 | 1.12 | 0.001 | 0.292 | 1.15 | 2.99 | 9.06 | 2.83 | -inf | 0.869 | 96.1 |
| Duplicate m = 2000 | 3 | 59.6 | 5.31 | -0.23 | 0.618 | 0.633 | 1.12 | 0.001 | 0.330 | 1.15 | 5.07 | 9.46 | 2.90 | -inf | 0.870 | 97.3 |
| Duplicate m = 5000 | 3 | 51.5 | 4.85 | -0.98 | 0.656 | 0.634 | 1.12 | 0.001 | 0.349 | 1.15 | 9.92 | 9.72 | 2.97 | -inf | 0.871 | 97.1 |
| K - Center - Greedy | 3 | 75.3 | 6.33 | 1.31 | 0.687 | 0.746 | 1.22 | 0.864 | 0.522 | 1.24 | 2.53 | 9.30 | 2.73 | -7.44 | 0.862 | 88.5 |
| Kmeans Clustering$_{1000}$ | 3 | 76.5 | 6.31 | 1.35 | 0.692 | 0.646 | 1.13 | 0.615 | 0.372 | 1.17 | 1.70 | 9.87 | 2.99 | -1.32 | 0.869 | 96.4 |
| Kmeans Cluster$_{100}$ | 3 | 74.4 | 6.36 | 1.28 | 0.693 | 0.650 | 1.13 | 0.610 | 0.362 | 1.16 | 1.69 | 9.78 | 2.99 | -1.33 | 0.869 | 96.1 |
| QDIT | 3 | 75.9 | 6.21 | 1.25 | 0.673 | 0.629 | 1.12 | 0.602 | 0.348 | 1.15 | 1.59 | 9.77 | 2.99 | -1.41 | 0.871 | 98.5 |
| Repr Filter | 3 | 72.6 | 6.17 | 1.05 | 0.671 | 0.703 | 1.18 | 0.799 | 0.470 | 1.21 | 2.23 | 9.45 | 2.86 | -9.12 | 0.866 | 92.0 |
| NoveSelect | 3 | 78.1 | 6.47 | 1.55 | 0.762 | 0.821 | 1.28 | 0.704 | 0.534 | 1.30 | 2.55 | 9.23 | 2.73 | -6.27 | 0.862 | 87.9 |

Table 8: Comprehensive experimental results on LLaMA-3-8B. Each data selection strategy variant is evaluated over three independent runs (except for random selection) to ensure the robustness and reliability of the findings. The results from multiple runs are averaged. Note that *NovelSelect* results are only included in Section 6 and are not part of the correlation calculations. Details of the data selection strategies are provided in Appendix A.3.

| Data Selection Strategy | Runs | Alpaca | MT-bench | Aggregated | NovelSum | DistSum$_{cosine}$ | DistSum$_{L2}$ | KNN | Inertia | Radius $10^{-3}$ | VS $10^7$ | Entropy | FL $10^5$ | LDD $10^4$ | TTR | vocd-D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Random - alpaca | 1 | 71.5 | 5.52 | -0.38 | 0.336 | 0.189 | 0.596 | 0.223 | 0.066 | 4.06 | 1.89 | 8.66 | 3.50 | -4.40 | 0.853 | 82.4 |
| Random - dolly | 1 | 55.2 | 6.24 | -0.49 | 0.320 | 0.221 | 0.651 | 0.293 | 0.098 | 4.52 | 3.04 | 7.92 | 3.46 | -3.62 | 0.844 | 76.6 |
| Random - sharegpt | 1 | 82.4 | 7.91 | 0.07 | 0.392 | 0.285 | 0.731 | 0.289 | 0.110 | 4.64 | 3.36 | 8.87 | 3.51 | -3.34 | 0.850 | 78.3 |
| Random - ultrachat | 1 | 78.0 | 7.66 | -0.02 | 0.389 | 0.200 | 0.620 | 0.252 | 0.074 | 4.11 | 2.09 | 9.30 | 3.52 | -4.17 | 0.881 | 110 |
| Random - wizardlm | 1 | 77.1 | 7.28 | -0.08 | 0.349 | 0.211 | 0.631 | 0.296 | 0.093 | 4.42 | 2.65 | 9.03 | 3.50 | -3.84 | 0.858 | 85.7 |
| Random $\mathcal{X}^{all}$ | 2 | 81.9 | 7.57 | 0.87 | 0.408 | 0.230 | 0.661 | 0.286 | 0.092 | 4.36 | 1.41 | 9.77 | 3.54 | -3.81 | 0.869 | 97.0 |
| Duplicate m=50 | 2 | 69.3 | 7.41 | -1.46 | 0.252 | 0.215 | 0.638 | 0.001 | 0.000 | 4.29 | 0.13 | 5.64 | 3.44 | -inf | 0.871 | 98.1 |
| Duplicate m=100 | 2 | 75.1 | 7.46 | -0.43 | 0.309 | 0.243 | 0.664 | 0.001 | 0.000 | 4.32 | 0.20 | 6.54 | 3.48 | -inf | 0.870 | 97.7 |
| Duplicate m=500 | 2 | 72.9 | 7.51 | -0.69 | 0.357 | 0.240 | 0.672 | 0.001 | 0.057 | 4.41 | 0.51 | 8.50 | 3.54 | -inf | 0.868 | 94.9 |
| Duplicate m=1000 | 2 | 78.6 | 7.59 | 0.37 | 0.364 | 0.229 | 0.658 | 0.001 | 0.076 | 4.36 | 0.74 | 9.05 | 3.56 | -inf | 0.869 | 97.0 |
| Duplicate m=5000 | 2 | 82.0 | 7.53 | 0.81 | 0.399 | 0.230 | 0.661 | 0.001 | 0.091 | 4.36 | 1.27 | 9.68 | 3.57 | -inf | 0.870 | 97.8 |
| K-Center-Greedy | 2 | 81.6 | 7.90 | 1.45 | 0.505 | 0.440 | 0.923 | 0.501 | 0.214 | 6.13 | 3.09 | 8.50 | 3.42 | -2.29 | 0.837 | 94.9 |
| K-means Clustering | 2 | 79.8 | 7.84 | 1.06 | 0.440 | 0.260 | 0.698 | 0.301 | 0.106 | 4.54 | 1.60 | 9.86 | 3.54 | -3.63 | 0.868 | 94.9 |
| QDIT | 2 | 80.0 | 7.81 | 1.00 | 0.403 | 0.223 | 0.650 | 0.283 | 0.091 | 4.33 | 2.60 | 9.74 | 3.54 | -3.87 | 0.871 | 99.1 |
| Repr Filter | 2 | 81.8 | 7.83 | 1.35 | 0.495 | 0.421 | 0.901 | 0.476 | 0.199 | 5.94 | 7.15 | 8.59 | 3.46 | -2.42 | 0.839 | 69.8 |

Table 9: Comprehensive experimental results on Qwen-2.5-7B. Each data selection strategy variant is evaluated over two independent runs (except for random selection).
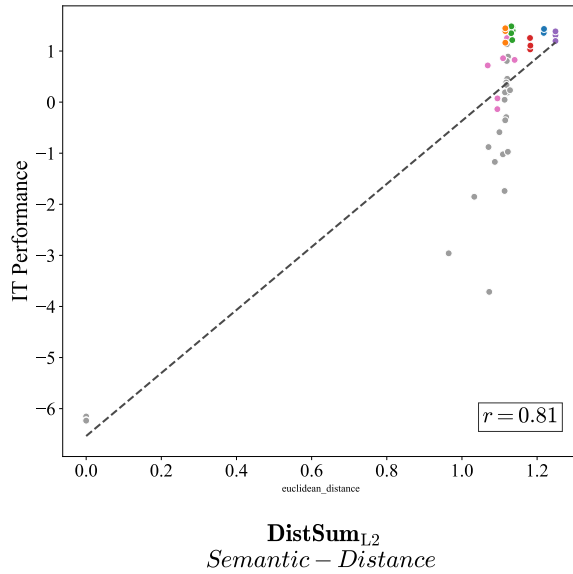


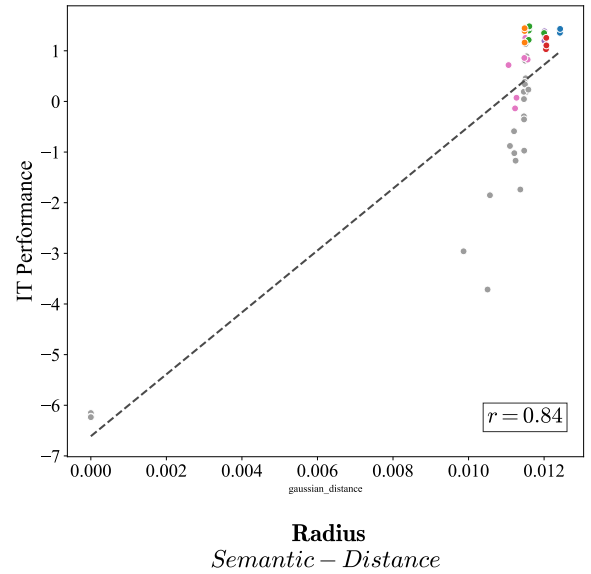Figure 6: Evaluation of DistSum$_{L2}$ metric by their correlation with IT performance.



Figure 7: Evaluation of Radius metric by their correlation with IT performance.

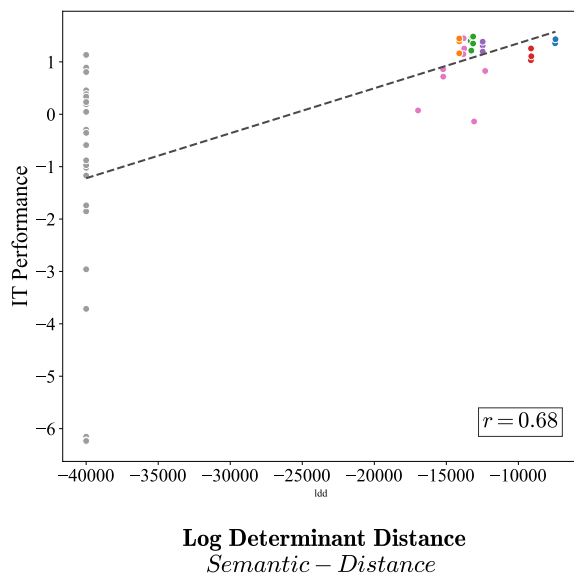**Log Determinant Distance**
*Semantic − Distance*

Figure 8: Evaluation of Log Determinant Distance metric by their correlation with IT performance.