

# LM-MIXUP: TEXT DATA AUGMENTATION VIA LANGUAGE MODEL BASED MIXUP

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Instruction tuning is crucial for aligning Large Language Models (LLMs), yet the quality of instruction-following data varies significantly. While high-quality data is paramount, it is often scarce; conversely, abundant low-quality data is frequently discarded, leading to substantial information loss. Existing data augmentation methods struggle to augment this low-quality data effectively, and the evaluation of such techniques remains poorly defined. To address this, we formally define the task of *Instruction Distillation*: distilling multiple low-quality and redundant inputs into high-quality and coherent instruction-output pairs. Specifically, we introduce a comprehensive data construction pipeline to create MIXTURE, a 144K-sample dataset pairing low-quality or semantically redundant imperfect instruction clusters with their high-quality distillations. We then introduce *LM-Mixup*, by first performing supervised fine-tuning on MIXTURE and then optimizing it with reinforcement learning. This process uses three complementary reward signals: quality, semantic alignment, and format compliance, via Group Relative Policy Optimization (GRPO). We demonstrate that *LM-Mixup* effectively augments imperfect datasets: fine-tuning LLMs on its distilled data, which accounts for only about 3% of the entire dataset, not only surpasses full-dataset training but also competes with state-of-the-art high-quality data selection methods across multiple benchmarks. Our work establishes that low-quality data is a valuable resource when properly distilled and augmented with *LM-Mixup*, significantly enhancing the efficiency and performance of instruction-tuned LLMs.

## 1 INTRODUCTION

In recent years, large language models (LLMs) have achieved notable progress in natural language processing and multimodal understanding (Team et al., 2023; Singhal et al., 2023; Deng et al., 2025; Li et al., 2024b; 2025a). This progress stems not only from improved architectures and larger scales but also from more efficient ways for models to learn and apply knowledge (Patil & Jadon, 2025; Dredze, 2025). While the conventional view holds that high-quality human alignment requires massive annotated data (Kim et al., 2024; Köpf et al., 2023), recent studies show that LLMs acquire most knowledge during pre-training (Brown et al., 2020; Roberts et al., 2020). Only a small, carefully curated dataset is sufficient for effective alignment in instruction tuning or supervised fine-tuning (SFT) (He et al., 2024; Wei et al., 2023), so many works now focus on selecting high-quality data, demonstrating that fine-tuning on such subsets alone can already yield strong performance (Pang et al., 2024; Fu et al., 2025; Jha et al., 2023). This shifts the research focus from “more data” to “better data”, emphasizing efficient high-quality data selection for model improvement.

However, high-quality samples are scarce and costly, while real-world datasets contain abundant redundant or low-quality data, leading to significant information waste. This gap mainly arises from data characteristics: low-quality samples are often simple or repetitive with limited learning signals, while high-quality samples involve complex reasoning or rich knowledge, making them more beneficial for training (Morishita et al., 2024), as shown in Figure 1. **Moreover, in many specialized domains or low-resource settings (e.g., low-resource machine translation and domain-specific tasks such as law or medicine), the scarcity of high-quality data is widely regarded as a key bottleneck that limits further progress (Abdalla et al., 2025; Dehouck & Gómez-Rodríguez, 2020; Alzubaidi et al., 2023).** Recently, some works have begun exploring ways to enhance low-quality data to unlock their potential; however, most existing approaches still rely on heuristic rules or handcrafted

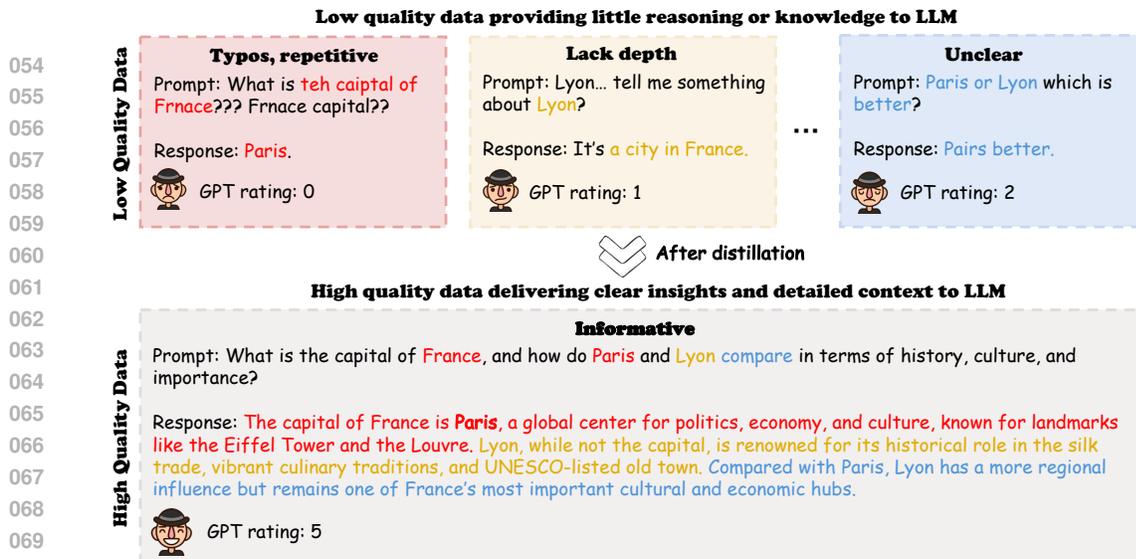


Figure 1: The goal of *Instruction Distillation*. Low-quality samples show issues such as typos, lack of depth, and unclear intent, each receiving low GPT ratings. After distillation, they are combined and refined into a single high-quality sample with clear, informative, and context-rich content. Ratings are on a 5-point scale. Additional case studies are provided in Appendix J.

templates, struggling to substantially enrich their information content or complexity (Chai et al., 2025; Zhu et al., 2025; Lee et al., 2024). This raises a key question: *can we fully exploit low-quality data and transform it into a valuable resource for improving LLM training?*

In this work, we study how to efficiently leverage low-quality data and introduce the *Instruction Distillation* paradigm: given topic-related but sparse and incomplete inputs, the goal is to aggregate and rewrite them into a single information-dense target. To facilitate this paradigm, we construct MIXTURE, a Wikipedia-based dataset with about 144K instances across five task types, providing hierarchical mappings from multiple low-quality inputs to a single high-quality output, as shown in Figure 2. Each high quality data pair with 2 to 20 controlled low quality variants and optional chain-of-thought supervision. To further improve diversity in the dataset and robustness during the training process, cross-topic mixing and noise injection are added.

Since SFT concentrates on memorizing answers (Li et al., 2025b; Chu et al., 2025) and fails to explore diverse strategies for distilling low-quality samples into high-quality outputs, we adopt GRPO (Guo et al., 2025) to optimize the generation process. Building on MIXTURE, we further train *LM-Mixup* with GRPO to fully leverage its potential. Concretely, we perform cold-start pre-training on a subset of MIXTURE to equip the model with the basic ability to generate high-quality outputs; then, we apply GRPO-based reinforcement learning to jointly optimize output quality along three dimensions: quality, semantic alignment, and format compliance. *LM-Mixup* significantly outperforms SFT and selective baselines across multiple tasks, with small-scale models even surpassing strong instruction models under direct prompting; moreover, a small amount of original high-quality data combined with distilled results from the low-quality data (totally 10K) can achieve or exceed the performance of large-scale datasets (300K) and advanced data selection methods, demonstrating excellent data efficiency and generalization.

Our contributions are summarized as follows:

- We introduce the *Instruction Distillation task*, which aims to transform sparse, incomplete, and low-quality inputs into a single information-dense output; to support this paradigm, we construct MIXTURE, a **144K-instance Wikipedia-based dataset** with hierarchical mappings from multiple low-quality variants to high-quality targets.
- We introduce *LM-Mixup*, initialized through cold-start pretraining and optimized with **GRPO-based reinforcement learning** using **multi-dimensional rewards** (quality, semantic alignment, and format compliance), achieving **superior performance on the MIXTURE test set** compared to SFT and strong baselines.
- Experiments show that training downstream models on the distilled data together with the original high-quality data, totaling only **≈3% of the full dataset**, matches or surpasses full-dataset training

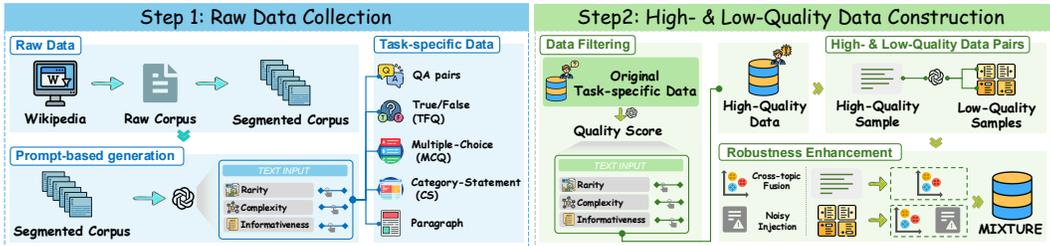


Figure 2: Overview of the MIXTURE construction pipeline. The process consists of two stages: (1) Raw data collection from Wikipedia, including segmentation and prompt-based generation of five task types; (2) High- and low-quality data construction via LLM-based scoring, multi-variant degradation, cross-topic fusion, and noisy injection, producing the final dataset for training.

and advanced data selection methods on public benchmarks, demonstrating the **value of low-quality data after distilling**.

## 2 RELATED WORK

**Data-centric AI.** Recent work on data-centric AI emphasizes improving model performance through high-quality data (Ng et al., 2021; 2022). These efforts can be broadly categorized as follows: data cleaning methods (Geerts et al., 2013; Krishnan et al., 2016; 2017; Zhang et al., 2021; Sen et al., 2022; Côté et al., 2024; Skjerve et al., 2025), such as the probabilistic model Holo-clean (Rekatsinas et al., 2017) and the automated framework by Mavrogiorgos et al. (2022); Data augmentation expands training sets through strategies such as linear interpolation, geometric transformations of images (Zhang et al., 2017), or introducing latent variables (Jiang & Mei, 2019) to enrich the parameter space; Human-led data labeling and annotation, often assisted by large language models, resulting in high-quality datasets (Wang et al., 2018; Rajpurkar et al., 2016). Recent data quality scoring systems like DS2 (Pang et al., 2025) correct LLM rating errors effectively; and data healing, which employs techniques such as model proposals, regularization adjustments and so on (Han et al., 2018; Yao et al., 2018; Tanno et al., 2019; Hu et al., 2019; Liu & Guo, 2020; Ma et al., 2020). Wang et al. (2024) developed NoiseGPT to detect and correct mislabeled instances, while Yang et al. (2024) utilized cubic regularization to efficiently identify noise. **Concurrently, MergeIT (Cai et al., 2025) focuses on merging similar medium-quality instructions, whereas we study distillation from multiple noisy low-quality responses. Although low-quality data is pervasive in real-world applications, research on systematically improving its quality remains limited, and our work aims to fill this gap.**

**Mixup methods.** Linear interpolation-based data augmentation has demonstrated significant advantages in enhancing model robustness and generalization across various domains (Zhang et al., 2020; Cao et al., 2024). These approaches have been adapted to textual data at multiple representation levels (Zhang & Vaidya, 2021; Chen et al., 2022). Previous work has explored interpolation in the word embedding space or sentence embedding space (Guo et al., 2019; Guo, 2020; Kong et al., 2022). Recent work has developed diverse methods to improve augmented data quality, including structural approaches like subtree decomposition in syntactic and semantic trees (Zhang et al., 2022), input-level tuning (Yoon et al., 2021), representation-level mixing of embeddings and latent features (Chen et al., 2020; Jindal et al., 2020) and so on (Yoon et al., 2021; Yang & Xiang, 2024; Zheng et al., 2023). Notably, Sun et al. (2020) were the first to directly apply mixup to textual data by fine-tuning Transformer models, enabling linear interpolation between text samples for effective augmentation. In this work, we leverage the mixup concept to systematically combine multiple low-quality, same-topic data samples into higher-quality instances, thereby enabling more effective utilization of abundant but noisy data resources.

## 3 MIXTURE: A DATASET FOR *Instruction Distillation*

### 3.1 TASK FORMULATION

The task of *Instruction Distillation* aims to **aggregate and** distill multiple potentially imperfect inputs (e.g., redundant or low-quality data) into a single, high-quality **instruction–response pair**. Let

$\mathcal{X}$  denote the universe of low-quality samples and  $\mathcal{Y}$  the set of high-quality texts. Each training instance consists of a multi-source input  $X = \{\ell_1, \dots, \ell_k\} \subset \mathcal{X}$  describing the same topic/task and a reference high-quality target  $Y \in \mathcal{Y}$  that is a single high-quality instruction–response pair drawn from the same task. We denote the training corpus as  $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^n$ . The objective is to learn a generator  $f_\theta$  that semantically fuses the multi-source inputs and produces a high-quality text  $\hat{Y} = f_\theta(X)$  which (i) preserves the salient information in  $X$  while denoising conflicts and increasing information density, (ii) aligns semantically with  $X$ , and (iii) adheres to the task-specific format. Compared with standard instruction tuning, the mapping is one-to-many: there exist multiple valid fusions for the same  $X$ , reflecting an underspecified target space that requires modeling output diversity while enforcing quality.

### 3.2 DATASET STATISTICS

We introduce MIXTURE, a dataset specifically designed for *Instruction Distillation*, comprising five task types: QA pairs, True/False (TFQ), Paragraph, Multiple-Choice Question (MCQ), and Category-Statement (CS). The overall pipeline is illustrated in Figure 2. Overall, Mixture comprises 144,884 samples spanning these five task types, with a balanced distribution across normal, cross-topic, and noisy variants, as shown in Appendix K.

### 3.3 RAW DATA COLLECTION

**Source Selection.** As shown in Step 1 of Figure 2, we use the Wikipedia dataset<sup>1</sup> as the initial source and sample entries across topical categories to construct the original corpus with broad coverage. To improve quality, we remove overly short entries, extract only plain text, and apply basic deduplication to eliminate redundant content. After filtering, about 10,000 Wikipedia entries are retained.

**Paragraph Segmentation.** Since Wikipedia articles are often long, directly feeding them into LLMs may cause inefficiency and instability (Liu et al., 2023a). We segment each article into semantically coherent blocks by first splitting into sentences and then greedily concatenating them until a token limit  $T^2$  is reached. To balance coherence and boundary effects, we allow optional overlaps and merge very short segments, while over-length sentences are further split at punctuation marks. The final blocks preserve the original order to ensure narrative consistency and traceability.

### 3.4 HIGH- AND LOW-QUALITY DATA CONSTRUCTION

**High-Quality Sample Generation.** As shown in Step 2 of Figure 2, to transform generic Wikipedia paragraphs into task-specific samples, we use a prompt-based rewriting approach with the ChatGPT-4o-mini (Achiam et al., 2023). All tasks follow the principle of information density to produce knowledge-rich outputs, with MCQ, TFQ, and CS pairs differing only in prompt formats (see Appendix I). For paragraph-level tasks, segmented text blocks that are correctly parsed and meet template constraints are directly used as high-quality samples; invalid ones are discarded.

Following previous work (Li et al., 2024a; Chen et al., 2025; Pang et al., 2024), we further use ChatGPT-4o-mini to perform quality scoring along multiple dimensions, including rarity, complexity, and informativeness. The scores are aggregated into a single overall rating, which is then discretized to a 1-5 scale. To ensure consistency across tasks and sessions, we retain only samples with a score of 4 or above as the final high-quality subset, while those below 4 are regarded as low-quality data.

**Low-Quality Sample Generation.** To enrich alignment signals, we use ChatGPT-4o-mini rewriting to generate multiple degraded variants for each high-quality sample, reducing information density or reasoning completeness while preserving the topic. For each target  $Y$ , we construct  $k \in \{2, \dots, 20\}$  variants  $X$ , with the distribution shown in Appendix (Figure 9), forming hierarchical mappings that teach the model to aggregate information and complete reasoning. We further add chain-of-thought traces as intermediate supervision to improve interpretability.

<sup>1</sup>[https://huggingface.co/datasets/lucadiliello/english\\_wikipedia](https://huggingface.co/datasets/lucadiliello/english_wikipedia)

<sup>2</sup>We use `cl1100k_base` tokenizer.

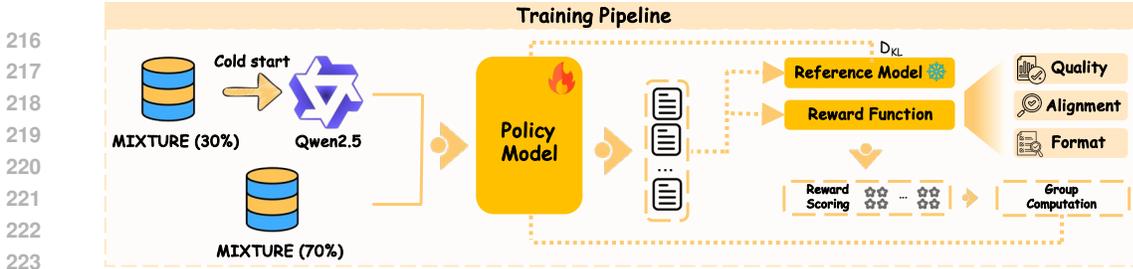


Figure 3: Overview of the training pipeline. The process involves cold-start pretraining on a subset of MIXTURE, followed by policy optimization using multi-dimensional rewards.

**Robustness Enhancement.** To enhance data diversity and robustness, we perform cross-topic synthesis by selecting semantically similar sample pairs with constrained entity overlaps, followed by GPT-based topic fusion rewriting, retaining only samples with quality scores above 4. To further improve generalization to noisy inputs, we inject surface-level perturbations such as spelling variations, synonym substitutions, and minor formatting shifts during training, while preserving some clean samples to balance robustness and fidelity.

#### 4 LM-Mixup: TRAINING FRAMEWORK FOR *Instruction Distillation*

This section is organized into three parts: (i) Cold Start Pretraining, (ii) Multi-Dimensional Reward Design, and (iii) Reinforcement Learning with GRPO. The overall framework of the proposed *LM-Mixup* training pipeline, built upon the Qwen-2.5-1.5B-Instruct (Team, 2024), is illustrated in Figure 3.

##### 4.1 COLD START

Directly starting reinforcement learning from randomly initialized parameters often leads to training instability (Guo et al., 2025; Wei et al., 2025). Therefore, we first perform cold start pretraining on the subset of MIXTURE. Specifically, given a high-quality sample  $Y$  and its corresponding  $k$  low-quality samples  $\{\ell_1, \ell_2, \dots, \ell_k\}$ , we linearize them into a conditional input sequence

$$X = \text{Linearize}(\ell_1, \dots, \ell_k), \quad (1)$$

and minimize the conditional likelihood with a standard autoregressive language modeling objective:

$$\mathcal{L}_{\text{CE}}(\theta) = - \sum_{t=1}^T \log p_{\theta}(y_t | y_{<t}, X), \quad (2)$$

where  $y_t$  denotes the  $t$ -th token of the target output  $Y$ . Through this stage, the model acquires basic language generation and information fusion capabilities, providing a stable initial policy distribution for subsequent reinforcement learning.

##### 4.2 REWARD DESIGN

To encourage the model to produce outputs with stronger information aggregation, semantic alignment, and structural conformity when mapping multiple low-quality samples to high-quality outputs, we design three complementary reward components. Given a model output  $\hat{Y}$  and the corresponding high-quality reference  $Y$ , the total reward is defined as

$$R(\hat{Y}, Y) = \lambda_q R_q(\hat{Y}) + \lambda_a R_a(\hat{Y}, Y) + \lambda_f R_f(\hat{Y}), \quad (3)$$

where  $\lambda_q, \lambda_a, \lambda_f$  are the normalized weights. In our experiments, we set  $\lambda_q = 0.5, \lambda_a = 0.4,$  and  $\lambda_f = 0.1$ .

**(1) Quality Reward  $R_q$ :** To efficiently approximate LLM ratings during training, we introduce a KNN-Bayes scoring scheme. Given a generated output  $\hat{Y}$ , we retrieve its  $k$  nearest neighbors from

a large reference set with pre-computed LLM scores and estimate the posterior distribution of the true quality label via a score transition matrix  $T$  (Zhu et al., 2021; Pang et al., 2024):

$$P(y = i \mid \mathbf{h}(\hat{Y})) \propto p_i \cdot \exp\left(\sum_j h_j(\hat{Y}) \log T_{ij}\right), \quad (4)$$

where  $\mathbf{h}(\hat{Y})$  is the neighbor rating histogram. The expected quality score  $\hat{s}(\hat{Y})$  from this posterior is then mapped into a parameterized piecewise reward:

$$R_q(\hat{Y}; \lambda, \kappa, \alpha, \beta) = \begin{cases} \alpha & \hat{s}(\hat{Y}) \geq \lambda, \\ \beta & \hat{s}(\hat{Y}) = \kappa, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

By default we set  $\lambda=4$ ,  $\kappa=3$ , and  $\alpha=1$ ,  $\beta=0.3$ . The offline construction of the reference set and transition matrix estimation is provided in Appendix C.2.

**(2) Semantic Alignment Reward  $R_a$ :** To ensure semantic consistency between generated outputs and reference answers, we encode both using the embedding model<sup>3</sup> and compute the normalized cosine similarity

$$R_a(\hat{Y}, Y) = \mathbb{1}\left(\text{cosine}(e(\hat{Y}), e(Y)) \geq \tau\right), \quad (6)$$

where  $e(\cdot)$  denotes the SentenceBERT encoder,  $\tau$  is the similarity threshold, and  $\mathbb{1}(\cdot)$  is the indicator function that returns 1 if the condition holds and 0 otherwise.

**(3) Format Compliance Reward  $R_f$ :** To enforce structural consistency with the `<think>...</think><answer>...</answer>` template, we use regular expressions to verify the output format. Outputs fully matching the template receive  $R_f(\hat{Y}) = 1$ , otherwise 0.

Finally, the total reward in Eq. equation 3 integrates quality, semantic alignment, and format compliance into a unified multi-dimensional signal.

### 4.3 REINFORCEMENT LEARNING WITH GRPO

Building on cold-start pretraining and designed rewards, we adopt GRPO for reinforcement-learning fine-tuning. Unlike standard SFT, which forces the model to imitate a single reference answer, our *Instruction Distillation* task allows infinitely many valid aggregation or generation strategies. Sole reliance on SFT risks overfitting to one canonical form and ignoring the diverse space of high-quality outputs. In contrast, reinforcement learning enables optimizing directly against reward signals, encouraging exploration of diverse outputs and progressively improving generation quality.

Specifically, GRPO is a variant of PPO (Schulman et al., 2017) that removes the need for a learned value (critic) function by replacing the baseline with group-wise statistics. For each input  $X$ , the model samples multiple candidate outputs  $\{\hat{Y}_1, \dots, \hat{Y}_m\}$ , which are scored by the multi-dimensional reward  $R(\hat{Y}_i, Y)$ . To reduce variance and mitigate scale inconsistency across candidates, GRPO computes a normalized reward within each group:

$$\tilde{R}_i = \frac{R(\hat{Y}_i, Y) - \mu_X}{\sigma_X + \epsilon}, \quad (7)$$

where  $\mu_X$  and  $\sigma_X$  are the mean and standard deviation of  $\{R(\hat{Y}_j, Y)\}_{j=1}^m$ , and  $\epsilon$  is a small constant to ensure numerical stability. The policy optimization objective becomes:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_X \left[ \min(r_i(\theta) \tilde{R}_i, \text{clip}(r_i(\theta), 1 - \epsilon_{\text{clip}}, 1 + \epsilon_{\text{clip}}) \tilde{R}_i) \right] - \beta \text{KL}(\pi_\theta(\cdot \mid X) \parallel \pi_{\theta_0}(\cdot \mid X)), \quad (8)$$

where  $r_i(\theta) = \frac{\pi_\theta(\hat{Y}_i \mid X)}{\pi_{\theta_0}(\hat{Y}_i \mid X)}$  is the importance ratio between the current policy  $\pi_\theta$  and the reference (old) policy  $\pi_{\theta_0}$ ;  $\epsilon_{\text{clip}}$  is the PPO clipping parameter; and  $\beta$  controls the strength of the KL regularization to ensure stability (Christiano et al., 2017).

<sup>3</sup><https://huggingface.co/BAAI/bge-m3>

Table 1: Performance comparison across different models on five tasks. The best results per column are highlighted in **bold**.

Model	cs	mcq	para	qa	tfq	Avg
LLaMA-3.1-8B-Instruct	3.61	2.57	3.57	3.71	2.10	3.27
LLaMA-3.2-3B-Instruct	3.58	2.66	3.52	3.78	2.49	3.21
DeepSeek-R1-Distill-Qwen-7B	3.61	2.46	3.40	3.23	2.41	3.02
Qwen-2.5-7B-Instruct	3.70	2.77	3.58	3.53	2.57	3.28
Qwen-2.5-1.5B-Instruct	3.39	2.44	3.34	3.33	1.34	2.86
GPT-4o-mini	3.81	2.86	<b>3.69</b>	3.64	2.61	3.37
Qwen-2.5-1.5B-SFT	3.54	3.25	2.82	3.73	3.05	3.28
Qwen-2.5-7B-SFT	3.53	3.31	3.41	3.78	3.10	3.46
<i>LM-Mixup</i>	<b>3.85</b>	<b>3.55</b>	3.31	<b>4.17</b>	<b>3.32</b>	<b>3.66</b>

Table 2: Data pool statistics.

Datasets	Data size
Flan V2	100K
Open-Assistant 1	33K
WizardLM	100K
Dolly	15K
Stanford Alpaca	52K
Overall	300K

#### 4.4 CAPACITY-CONSTRAINED CLUSTERING

After GRPO training, the model can distill multiple low-quality samples into high-quality ones. For downstream tasks, we introduce a Capacity-Constrained Clustering method to automatically collect low-quality inputs with flexible control over cluster number and size, which also mitigates the severe imbalance or over-fragmentation issues often observed in standard clustering methods. Given a text collection  $\mathcal{D} = \{x_i\}_{i=1}^N$ , we encode each sample into  $\mathbf{h}_i \in \mathbb{R}^d$  using a pre-trained encoder. A target capacity vector  $\mathbf{c} = (c_1, \dots, c_K)$  is drawn from a truncated normal distribution with  $c_k \in [c_{\min}, c_{\max}]$ . We then perform two-stage clustering: (i) run MiniBatchKMeans to obtain  $k$  initial cluster centers  $\{\mathbf{c}_k\}$ ; (ii) iteratively assign samples to the most similar clusters under capacity constraints, with a few refinement steps to ensure semantic compactness and balanced partitioning.

## 5 EXPERIMENTS

### 5.1 MIXTURE EXPERIMENTAL RESULTS

**Experimental Setup.** To comprehensively evaluate the performance of *LM-Mixup* on the MIXTURE dataset, we conducted standardized experiments on the test set using a variety of models. Specifically, the experiments involved the following models: ChatGPT-4o-mini (Achiam et al., 2023), Qwen-2.5-1.5B-Instruct, Qwen-2.5-7B-Instruct (Team, 2024), LLaMA-3.1-8B-Instruct, LLaMA-3.2-3B-Instruct (Dubey et al., 2024), DeepSeek-R1-Distill-Qwen-7B (Guo et al., 2025), Qwen-2.5-1.5B-SFT, Qwen-2.5-7B-SFT (obtained via supervised fine-tuning on the full MIXTURE). All models were evaluated on the same test set, constructed by holding out a non-overlapping 20% split from MIXTURE, under identical prompt conditions to ensure fair comparison. For automated evaluation, we employed ChatGPT-4o-mini as the rating model to assess the quality of the generated outputs.

**Results.** Table 1 presents the performance comparison across different models on the MIXTURE test set. It can be observed that *LM-Mixup* consistently outperforms all baseline models, achieving the best overall results. Compared with standard supervised fine-tuning, GRPO training with multi-dimensional quality rewards enables the model to learn generation patterns that produce higher-quality answers rather than merely mimicking the ground truth.

### 5.2 OPENLLM LEADERBOARD EVALUATION RESULTS

**Experimental Setup.** Following previous work (Pang et al., 2024), to evaluate *LM-Mixup*'s performance on OOD datasets, we construct an additional data pool consisting of Flan\_v2 (Longpre et al., 2023), Open Assistant 1 (Köpf et al., 2023), WizardLM (Xu et al., 2023), Dolly (Conover et al., 2023), and Stanford Alpaca (Taori et al., 2023). Detailed statistics of the data pool are provided in Table 2. To identify low-quality samples within this pool, we employ ChatGPT-4o-mini for quality rating following the same protocol described in Sec 3.4, where samples with a score below 4 are regarded as low-quality data. We then apply *LM-Mixup* to perform mixup on the low-quality samples within the data pool and compute long-tail scores using embeddings. The top-ranked samples from both the original high-quality data and the mixup data from low-quality data are then selected for instruction fine-tuning.

Table 3: Results on the OpenLLM leaderboard using LLaMA-3.1-8B as the base model. The top-performing scores are shown in **bold**, while the second-best scores are marked with underlines. Unless otherwise specified, the size of the fine-tuning dataset is 10K. \* indicates that the values are sourced from Pang et al. (2024).

Model	MMLU (factuality)	TruthfulQA (truthfulness)	GSM (reasoning)	BBH (reasoning)	TyDiQA (multilinguality)	Average
VANILLA BASE MODEL*	64.1	33.5	56.5	55.4	23.3	46.6
COMPLETION LENGTH*	64.2	41.4	62.5	60.7	23.0	50.4
PERPLEXITY*	63.1	40.4	55.5	60.2	62.1	56.3
$k$ -NN-10*	62.4	44.3	57.0	59.1	63.8	57.3
RANDOM SELECTION*	63.4	39.1	62.2	61.3	61.1	57.4
LESS*	63.0	39.0	57.5	63.1	67.2	58.0
FULL DATA (300K)*	63.5	42.0	61.0	59.1	62.8	57.7
ALPAGASUS*	63.4	42.6	66.0	59.1	59.4	58.1
DEITA*	64.5	50.1	60.0	60.3	63.7	59.7
DS2 W/O CURATION*	63.3	51.5	62.0	59.7	64.3	60.2
DS2*	64.0	50.3	67.5	59.0	66.1	<u>61.4</u>
BACK-TRANSLATION	62.0±0.4	46.5±2.9	61.2±0.8	58.8±2.2	60.2±0.8	57.7±0.1
EDA	61.6±0.9	43.7±2.0	56.2±1.0	59.7±0.3	62.0±1.6	56.6±0.6
REPHRASING	61.4±0.7	36.0±2.5	63.2±1.0	59.6±0.2	62.2±0.8	56.5±0.6
BASE 70% + ORI 30%	62.2±0.9	40.7±0.5	54.3±0.2	55.1±1.1	23.2±0.3	47.1±0.1
BASE 50% + ORI 50%	62.1±0.1	37.4±0.7	50.8±0.6	54.1±0.7	22.9±0.4	45.4±0.2
BASE 30% + ORI 70%	61.3±0.7	38.2±0.7	51.2±1.4	54.2±0.9	23.0±0.5	45.6±0.4
LOW 70% + ORI 30%	62.7±0.7	17.8±2.0	62.5±4.0	60.3±0.9	65.6±1.1	53.6±1.2
MIXUP 70% + ORI 30%	63.0±0.2	47.9±0.3	63.3±0.6	61.1±0.3	64.2±0.6	<b>59.9±0.1</b> <sup>†6.3</sup>
LOW 50% + ORI 50%	62.4±0.6	39.0±9.8	62.7±1.1	61.0±2.2	64.0±0.3	57.9±2.0
MIXUP 50% + ORI 50%	63.3±0.3	52.6±0.1	65.5±0.2	61.3±0.3	64.6±0.3	<b>61.5±0.1</b> <sup>†3.6</sup>
LOW 30% + ORI 70%	60.9±2.1	41.1±5.6	62.7±1.9	59.9±1.7	60.4±1.7	57.0±0.8
MIXUP 30% + ORI 70%	<b>63.1±0.3</b>	46.8±0.6	61.2±1.5	58.0±0.2	63.4±1.1	<b>58.5±0.1</b> <sup>†1.5</sup>

**Metrics.** We report task-specific metrics, including accuracy on MMLU (Hendrycks et al., 2020), BBH (Suzgun et al., 2022), and GSM8K (Cobbe et al., 2021), the Informative-Truthful Rate on TruthfulQA (Lin et al., 2021), and F1 scores on TyDiQA (Clark et al., 2020).

**Training Settings.** We fine-tune three base models, LLaMA-2-7B (Touvron et al., 2023), LLaMA-3.1-8B (Dubey et al., 2024), and Mistral-7B-v0.3 (Jiang et al., 2023), on 10K samples under six settings: three with *mixup* data, three with *direct* low-quality samples (without mixup), each combined with original high-quality data at 70%, 50%, or 30%. **For all experiments newly conducted by us, we report the average results over three independent runs to ensure statistical reliability, whereas the baseline results sourced from prior work are reported as originally published.** Notably, for the full original data setting, we adopt DS2 (Pang et al., 2024) as the baseline for comparison. Additionally, we also provide the results for both the full mixup data and the full low-quality data in Appendix G. **Further details of the additional training settings are provided in Appendix F.**

**Baselines.** We compare our method against several representative data selection baselines commonly used in LLM fine-tuning, including *Random Selection*, *Completion Length*, *Perplexity*, *k-NN*, *LESS* (Xia et al., 2024), *AlpaGasus* (Chen et al., 2023b), *DEITA* (Liu et al., 2023b), *DS2* (Pang et al., 2024), and *Full Data*. **We also include widely used data augmentation baselines, including back-translation (Edunov et al., 2018), paraphrasing (Abaskohi et al., 2023), and EDA (Wei & Zou, 2019).** Additionally, to more comprehensively assess the gains achieved by *LM-Mixup*, we also report the zero-shot performance of Qwen-1.5B-Instruct without any further training. Detailed descriptions of these baselines are provided in Appendix E.

**Low-quality data matters: After being processed by LM-Mixup, it can even outperform high-quality-only baselines.** Table 3 shows that combining low-quality data (after mixup) with original high-quality samples can surpass baselines that rely solely on high-quality data selection. In particular, the 50% mixup + 50% original configuration achieves the top average score across all five OpenLLM Leaderboard benchmarks, with Tables 5 and 6 showing similar trends on Mistral-7B and LLaMA-2-7B. This demonstrates that even low-quality data, when fused into high-quality samples, can enhance diversity and complement real data to boost performance. Additional results on more models are provided in the Appendix G.

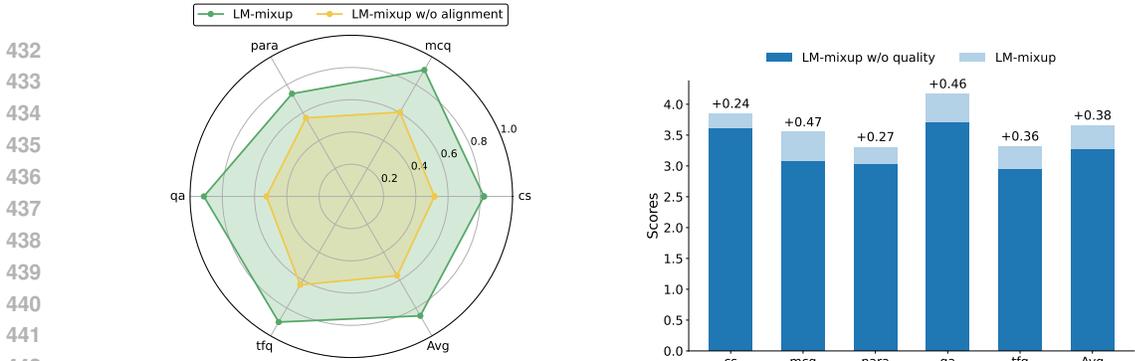


Figure 4: Ablation study on reward components in *LM-Mixup*. The left figure evaluates the effect of removing the alignment reward, while the right figure shows the impact of removing the quality reward across different tasks.

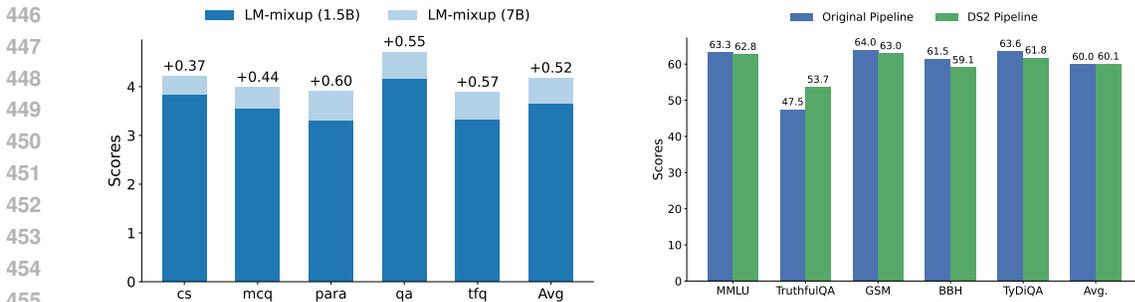


Figure 5: Effect of model scaling on performance.

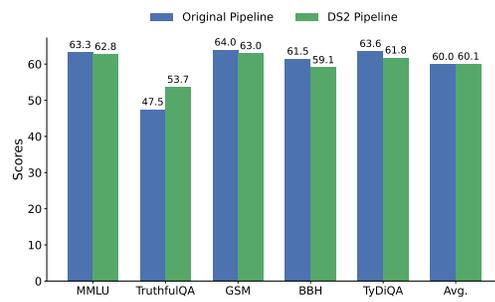


Figure 6: Comparison of the *LM-Mixup* using DS2 pipelines under the Mixup 70% + Ori 30%.

***LM-Mixup* demonstrates strong performance.** As shown in Table 3, *LM-Mixup* significantly outperforms the zero-shot Qwen-1.5B-Instruct model used for mixup without any training. Furthermore, our method consistently surpasses standard data augmentation baselines, indicating that it truly elevates low-quality samples into high-quality supervision signals rather than merely increasing surface-level diversity.

**Mixup outperforms the full data pool with only 3.3% of the data.** Notably, on both LLaMA-3.1-8B (Table 3) and Mistral-7B (Table 5), our best-performing mixup configuration using only 10K training samples even surpasses the 300K full data pool baseline, demonstrating that mixup not only enhances data diversity but also enables a highly compact training set to outperform large-scale unfiltered data.

**Effect of mixup on low quality data.** As shown in Table 3, applying *LM-Mixup* to low-quality data consistently improves performance across all mixture ratios. E.g., in the 70% low-quality + 30% high-quality setting, *LM-Mixup* raises the score from 54.2 to 60.0 ( $\uparrow 5.8$ ), with similar gains in the 50% ( $\uparrow 1.6$ ) and 30% ( $\uparrow 2.2$ ) settings. This highlights that properly modeling low-quality data can yield substantial benefits for model training.

### 5.3 ABALTION STUDY

**Ablation on Reward Components.** We conduct ablation studies to investigate the contribution of each reward component in our GRPO-based *LM-Mixup* training framework, which incorporates quality and alignment rewards alongside the base objective. As shown in Figure 4, removing the alignment reward causes the model to exhibit reward hacking behavior: it tends to memorize answers from the reference set regardless of the input, leading to significantly lower semantic similarity with the ground truth. On the other hand, removing the quality reward makes the model behave similarly to standard SFT, producing outputs with limited quality improvement. These results highlight that both rewards are essential: the alignment reward ensures semantic faithfulness to the input, while the quality reward drives the generation of high-quality outputs beyond simple imitation.

486 **Effect of Model Scaling.** To investigate the impact of scaling up model parameters, we extend  
487 our training pipeline from Qwen-2.5-1.5B-Instruct to Qwen-2.5-7B-Instruct using the same GRPO-  
488 based optimization described in Sec.4. As shown in Fig.5, the larger 7B model consistently outper-  
489 forms its 1.5B counterpart across all tasks on the MIXTURE test set, achieving an average score of  
490 4.18 compared to 3.66 on the smaller model. These results demonstrate both the effectiveness and  
491 the scalability of our approach when applied to models with larger parameter sizes.

492 **Revisiting LLM rating bias.** Recent work has noted that LLM-as-judge scores can be biased (Ye  
493 et al., 2024; Chen et al., 2024). In our pipeline we use ChatGPT-4o-mini for rating, which may  
494 introduce such bias. To assess sensitivity, we conducted experiments using the DS2 pipeline (Pang  
495 et al., 2024). We conduct the same experiments described in Sec.5.2 experiments under the Mixup  
496 70% + Ori 30% setting, where the overall performance shows only marginal changes, as shown in  
497 Fig.6. We hypothesize two reasons: (i) *LM-Mixup*’s GRPO with multi-dimensional rewards and  
498 many-to-one mixup supervision provides strong signals that attenuate upstream rating noise; and  
499 (ii) diversity is governed by embedding-based long-tail selection, largely independent of the rating  
500 scale. Overall, while LLM rating bias is real, our design appears tolerant to moderate bias; further  
501 de-biasing (e.g., multi-judge ensembling, cross-model adjudication, or light human spot-checks)  
502 may be needed to unlock additional gains.

## 503 504 6 CONCLUSION

506 In this work, we introduce *Instruction Distillation* and present a comprehensive data construction  
507 pipeline to create MIXTURE, a large-scale dataset pairing low-quality, noisy and redundant instruc-  
508 tion clusters with their high-quality distillations. Building on this MIXTURE, we propose *LM-Mixup*,  
509 a model trained with supervised fine-tuning followed by reinforcement learning using customized  
510 rewards. Our results demonstrate that: *LM-Mixup* can efficiently distill plenty of imperfect data  
511 samples into condensed high-quality ones, significantly compress the training data size, fully ex-  
512 tract the information value of neglected low quality data, and meanwhile effectively enhance the  
513 efficiency and performance of instruction-tuned LLMs.

514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

## ETHICS STATEMENT

We use both publicly available datasets and our constructed MIXTURE derived from Wikipedia text under its open license, ensuring no sensitive or private information is included. LM-Mixup distills low-quality or redundant samples into high-quality data while filtering harmful content, reducing data scale and computational cost to support responsible and sustainable AI development.

## REPRODUCIBILITY STATEMENT

Our experimental settings are detailed in Section 5 and Appendix D, and both the code and dataset will be released upon paper acceptance.

## REFERENCES

- Amirhossein Abaskohi, Sascha Rothe, and Yadollah Yaghoobzadeh. Lm-cppf: Paraphrasing-guided data augmentation for contrastive prompt-based few-shot fine-tuning. *arXiv preprint arXiv:2305.18169*, 2023.
- Hemn Barzan Abdalla, Yulia Kumar, Jose Marchena, Stephany Guzman, Ardalan Awlla, Mehdi Gheisari, and Maryam Cheraghy. The future of artificial intelligence in the face of data scarcity. *Computers, Materials & Continua*, 84(1), 2025.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Laith Alzubaidi, Jinshuai Bai, Aiman Al-Sabaawi, Jose Santamaría, Ahmed Shihab Albahri, Bashar Sami Nayyef Al-Dabbagh, Mohammed A Fadhel, Mohamed Manoufali, Jinglan Zhang, Ali H Al-Timemy, et al. A survey on deep learning tools dealing with data scarcity: definitions, challenges, solutions, tips, and applications. *Journal of Big Data*, 10(1):46, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Hongyi Cai, Yuqian Fu, Hongming Fu, and Bo Zhao. Mergeit: From selection to merging for efficient instruction tuning. *arXiv preprint arXiv:2503.00034*, 2025.
- Chengtai Cao, Fan Zhou, Yurou Dai, Jianping Wang, and Kunpeng Zhang. A survey of mix-based data augmentation: Taxonomy, methods, applications, and explainability. *ACM Computing Surveys*, 57(2):1–38, 2024.
- Yaping Chai, Haoran Xie, and Joe S Qin. Text data augmentation for large language models: A comprehensive survey of methods, challenges, and opportunities. *arXiv preprint arXiv:2501.18845*, 2025.
- Guiming Hardy Chen, Shunian Chen, Ziche Liu, Feng Jiang, and Benyou Wang. Humans or llms as the judge? a study on judgement biases. *arXiv preprint arXiv:2402.10669*, 2024.
- Hui Chen, Wei Han, Diyi Yang, and Soujanya Poria. Doublemix: Simple interpolation-based data augmentation for text classification. *arXiv preprint arXiv:2209.05297*, 2022.
- Jiaao Chen, Zichao Yang, and Diyi Yang. Mixtext: Linguistically-informed interpolation of hidden space for semi-supervised text classification. *arXiv preprint arXiv:2004.12239*, 2020.
- Jiaao Chen, Derek Tam, Colin Raffel, Mohit Bansal, and Diyi Yang. An empirical survey of data augmentation for limited data learning in nlp. *Transactions of the Association for Computational Linguistics*, 11:191–211, 2023a.
- Junkai Chen, Zhijie Deng, Kening Zheng, Yibo Yan, Shuliang Liu, PeiJun Wu, Peijie Jiang, Jia Liu, and Xuming Hu. Safeeraser: Enhancing safety in multimodal large language models through multimodal machine unlearning. *arXiv preprint arXiv:2502.12520*, 2025.

- 594 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa Gunaratna, Vikas Yadav, Zheng Tang, Vijay  
595 Srinivasan, Tianyi Zhou, Heng Huang, et al. Alpagasus: Training a better alpaca with fewer data.  
596 *arXiv preprint arXiv:2307.08701*, 2023b.
- 597 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep  
598 reinforcement learning from human preferences. *Advances in neural information processing sys-*  
599 *tems*, 30, 2017.
- 601 Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V.  
602 Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation  
603 model post-training, 2025. URL <https://arxiv.org/abs/2501.17161>.
- 604 Jonathan H Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev,  
605 and Jennimaria Palomaki. Tydi qa: A benchmark for information-seeking question answering in  
606 ty pologically diverse languages. *Transactions of the Association for Computational Linguistics*,  
607 8:454–470, 2020.
- 609 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
610 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
611 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 612 Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick  
613 Wendell, Matei Zaharia, and Reynold Xin. Free dolly: Introducing the world’s first truly open  
614 instructiontuned llm. *arXiv preprint arXiv:2304.12244*, 2023.
- 616 Pierre-Olivier Côté, Amin Nikanjam, Nafisa Ahmed, Dmytro Humeniuk, and Foutse Khomh. Data  
617 cleaning and machine learning: a systematic literature review. *Automated Software Engineering*,  
618 31(2):54, 2024.
- 619 Mathieu Dehouck and Carlos Gómez-Rodríguez. Data augmentation via subtree swapping for  
620 dependency parsing of low-resource languages. In *28th international conference on computa-*  
621 *tional linguistics*, pp. 3818–3830. International Committee on Computational Linguistics; Inter-  
622 national . . . , 2020.
- 624 Zhijie Deng, Chris Yuhao Liu, Zirui Pang, Xinlei He, Lei Feng, Qi Xuan, Zhaowei Zhu, and Jiaheng  
625 Wei. Guard: Generation-time llm unlearning via adaptive restriction and detection. *arXiv preprint*  
626 *arXiv:2505.13312*, 2025.
- 627 Kaiser Sun Mark Dredze. Amuro & char: Analyzing the relationship between pre-training and fine-  
628 tuning of large language models. In *10th Workshop on Representation Learning for NLP*, pp. 131,  
629 2025.
- 631 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
632 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
633 *arXiv e-prints*, pp. arXiv–2407, 2024.
- 634 Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. Understanding back-translation at  
635 scale. *arXiv preprint arXiv:1808.09381*, 2018.
- 637 Steven Y Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mita-  
638 mura, and Eduard Hovy. A survey of data augmentation approaches for nlp. *arXiv preprint*  
639 *arXiv:2105.03075*, 2021.
- 640 Yanjun Fu, Faisal Hamman, and Sanghamitra Dutta. T-shirt: Token-selective hierarchical data se-  
641 lection for instruction tuning. *arXiv preprint arXiv:2506.01317*, 2025.
- 643 Floris Geerts, Giansalvatore Mecca, Paolo Papotti, and Donatello Santoro. The llunatic data-  
644 cleaning framework. *Proceedings of the VLDB Endowment*, 6(9):625–636, 2013.
- 645 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
646 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
647 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

- 648 Hongyu Guo. Nonlinear mixup: Out-of-manifold data augmentation for text classification. In  
649 *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 4044–4051, 2020.  
650
- 651 Hongyu Guo, Yongyi Mao, and Richong Zhang. Augmenting data with mixup for sentence classifi-  
652 cation: An empirical study. *arXiv preprint arXiv:1905.08941*, 2019.  
653
- 654 Bo Han, Jiangchao Yao, Gang Niu, Mingyuan Zhou, Ivor Tsang, Ya Zhang, and Masashi Sugiyama.  
655 Masking: A new perspective of noisy supervision. *Advances in neural information processing*  
656 *systems*, 31, 2018.
- 657 Yexiao He, Ziyao Wang, Zheyu Shen, Guoheng Sun, Yucong Dai, Yongkai Wu, Hongyi Wang, and  
658 Ang Li. Shed: Shapley-based automated dataset refinement for instruction fine-tuning. *Advances*  
659 *in Neural Information Processing Systems*, 37:99382–99403, 2024.  
660
- 661 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and  
662 Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint*  
663 *arXiv:2009.03300*, 2020.
- 664 Wei Hu, Zhiyuan Li, and Dingli Yu. Simple and effective regularization methods for training on  
665 noisily labeled data with generalization guarantee. *arXiv preprint arXiv:1905.11368*, 2019.  
666
- 667 Aditi Jha, Sam Havens, Jeremy Dohmann, Alex Trott, and Jacob Portes. Limit: Less is more for  
668 instruction tuning across evaluation paradigms. *arXiv preprint arXiv:2311.13133*, 2023.  
669
- 670 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap-  
671 lot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
672 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril,  
673 Thomas Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.  
674
- 675 Ruoqiao Jiang and Shaohui Mei. Polar coordinate convolutional neural network: From rotation-  
676 invariance to translation-invariance. In *2019 IEEE International Conference on Image Processing*  
677 *(ICIP)*, pp. 355–359. IEEE, 2019.
- 678 Amit Jindal, Arijit Ghosh Chowdhury, Aniket Didolkar, Di Jin, Ramit Sawhney, and Rajiv Shah.  
679 Augmenting nlp models using latent feature interpolations. In *Proceedings of the 28th Interna-*  
680 *tional Conference on Computational Linguistics*, pp. 6931–6936, 2020.  
681
- 682 Davood Karimi, Haoran Dou, Simon K Warfield, and Ali Gholipour. Deep learning with noisy  
683 labels: Exploring techniques and remedies in medical image analysis. *Medical image analysis*,  
684 65:101759, 2020.
- 685
- 686 Dongyoung Kim, Kimin Lee, Jinwoo Shin, and Jaehyung Kim. Aligning large language models  
687 with self-generated preference data. *arXiv e-prints*, pp. arXiv–2406, 2024.
- 688 Fanshuang Kong, Richong Zhang, Xiaohui Guo, Samuel Mensah, and Yongyi Mao. Dropmix: A  
689 textual data augmentation combining dropout with mixup. In *Proceedings of the 2022 Conference*  
690 *on Empirical Methods in Natural Language Processing*, pp. 890–899, 2022.  
691
- 692 Andreas K opf, Yannic Kilcher, Dimitri Von R utte, Sotiris Anagnostidis, Zhi Rui Tam, Keith  
693 Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Rich ard Nagyfi, et al. Openassistant  
694 conversations-democratizing large language model alignment. *Advances in neural information*  
695 *processing systems*, 36:47669–47681, 2023.
- 696 Sanjay Krishnan, Jiannan Wang, Eugene Wu, Michael J Franklin, and Ken Goldberg. Activeclean:  
697 Interactive data cleaning for statistical modeling. *Proceedings of the VLDB Endowment*, 9(12):  
698 948–959, 2016.  
699
- 700 Sanjay Krishnan, Michael J Franklin, Ken Goldberg, and Eugene Wu. Boostclean: Automated  
701 error detection and repair for machine learning. corr abs/1711.01299 (2017). *arXiv preprint*  
*arXiv:1711.01299*, 2017.

- 702 Nicholas Lee, Thanakul Wattanawong, Sehoon Kim, Karttikeya Mangalam, Sheng Shen, Gopala  
703 Anumanchipalli, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. Llm2llm: Boosting  
704 llms with novel iterative data enhancement. *arXiv preprint arXiv:2403.15042*, 2024.
- 705  
706 Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun  
707 Liu. Llms-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint*  
708 *arXiv:2412.05579*, 2024a.
- 709 Ling Li, Yu Ye, Bingchuan Jiang, and Wei Zeng. Georeasoner: Geo-localization with reasoning  
710 in street views using a large vision-language model. In *Forty-first International Conference on*  
711 *Machine Learning*, 2024b.
- 712  
713 Ling Li, Yao Zhou, Yuxuan Liang, Fugee Tsung, and Jiaheng Wei. Recognition through rea-  
714 soning: Reinforcing image geo-localization with large vision-language models. *arXiv preprint*  
715 *arXiv:2506.14674*, 2025a.
- 716 Ziniu Li, Congliang Chen, Tian Xu, Zeyu Qin, Jiancong Xiao, Zhi-Quan Luo, and Ruoyu Sun.  
717 Preserving diversity in supervised fine-tuning of large language models, 2025b. URL <https://arxiv.org/abs/2408.16673>.
- 718  
719 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human  
720 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 721  
722 Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and  
723 Percy Liang. Lost in the middle: How language models use long contexts. *arXiv preprint*  
724 *arXiv:2307.03172*, 2023a.
- 725  
726 Wei Liu, Weihao Zeng, Keqing He, Yong Jiang, and Junxian He. What makes good data for align-  
727 ment? a comprehensive study of automatic data selection in instruction tuning. *arXiv preprint*  
728 *arXiv:2312.15685*, 2023b.
- 729  
730 Yang Liu and Hongyi Guo. Peer loss functions: Learning from noisy labels without knowing noise  
731 rates. In *International conference on machine learning*, pp. 6226–6236. PMLR, 2020.
- 732  
733 Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V  
734 Le, Barret Zoph, Jason Wei, et al. The flan collection: Designing data and methods for effective  
735 instruction tuning. In *International Conference on Machine Learning*, pp. 22631–22648. PMLR,  
736 2023.
- 737  
738 Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah Erfani, and James Bailey. Nor-  
739 malized loss functions for deep learning with noisy labels. In *International conference on machine*  
740 *learning*, pp. 6543–6553. PMLR, 2020.
- 741  
742 Konstantinos Mavrogiorgos, Argyro Mavrogiorgou, Athanasios Kiourtis, Nikolaos Zafeiropoulos,  
743 Spyridon Kleftakis, and Dimosthenis Kyriazis. Automated rule-based data cleaning using nlp. In  
744 *2022 32nd Conference of Open Innovations Association (FRUCT)*, pp. 162–168. IEEE, 2022.
- 745  
746 Terufumi Morishita, Gaku Morio, Atsuki Yamaguchi, and Yasuhiro Sogawa. Enhancing reason-  
747 ing capabilities of llms via principled synthetic logic corpus. *Advances in Neural Information*  
748 *Processing Systems*, 37:73572–73604, 2024.
- 749  
750 Andrew Ng, Lora Aroyo, Cody Coleman, Greg Diamos, Vijay Reddi, Joaquin Vanschoren, Carole-  
751 Jean Wu, and Matei Zaharia. Data-centric AI workshop. In *NeurIPS 2021 Workshop on Data-*  
752 *Centric AI*, 2021.
- 753  
754 Andrew Ng, Laird D., and He L. Data-centric AI competition, 2022. URL [https://](https://deeplearning-ai.github.io/data-centriccomp/)  
755 [deeplearning-ai.github.io/data-centriccomp/](https://deeplearning-ai.github.io/data-centriccomp/).
- 756  
757 Curtis G Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize  
758 machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021.
- 759  
760 Jinlong Pang, Jiaheng Wei, Ankit Parag Shah, Zhaowei Zhu, Yaxuan Wang, Chen Qian, Yang Liu,  
761 Yujia Bao, and Wei Wei. Improving data efficiency via curating llm-driven rating systems. *arXiv*  
762 *preprint arXiv:2410.10877*, 2024.

- 756 Jinlong Pang, Jiaheng Wei, Ankit Parag Shah, Zhaowei Zhu, Yaxuan Wang, Chen Qian, Yang Liu,  
757 Yujia Bao, and Wei Wei. Improving data efficiency via curating llm-driven rating systems, 2025.  
758 URL <https://arxiv.org/abs/2410.10877>.
- 759 Avinash Patil and Aryan Jadon. Advancing reasoning in large language models: Promising methods  
760 and approaches. *arXiv preprint arXiv:2502.03671*, 2025.
- 761 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions  
762 for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- 763 Theodoros Rekatsinas, Xu Chu, Ihab F Ilyas, and Christopher Ré. Holoclean: Holistic data repairs  
764 with probabilistic inference. *arXiv preprint arXiv:1702.00820*, 2017.
- 765 Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the  
766 parameters of a language model? *arXiv preprint arXiv:2002.08910*, 2020.
- 767 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
768 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 769 Snigdha Sen, Sonali Agarwal, Pavan Chakraborty, and Krishna Pratap Singh. Astronomical big data  
770 processing using machine learning: A comprehensive review. *Experimental Astronomy*, 53(1):  
771 1–43, 2022.
- 772 Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan  
773 Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode  
774 clinical knowledge. *Nature*, 620(7972):172–180, 2023.
- 775 Torgunn Aslaug Skjerve, Gunnar Klemetsdal, Bente Aspeholen Åby, Jon Kristian Sommerseth,  
776 Ulf Geir Indahl, and Hanne Fjerdingsby Olsen. Using density and fuzzy clustering for data clean-  
777 ing and segmental description of livestock data: Ta skjerve et al. *Journal of Agricultural, Biolog-  
778 ical and Environmental Statistics*, 30(3):870–885, 2025.
- 779 Lichao Sun, Congying Xia, Wenpeng Yin, Tingting Liang, Philip S Yu, and Lifang He. Mixup-  
780 transformer: Dynamic data augmentation for nlp tasks. *arXiv preprint arXiv:2010.02394*, 2020.
- 781 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,  
782 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging big-bench tasks  
783 and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*, 2022.
- 784 Ryutaro Tanno, Ardavan Saeedi, Swami Sankaranarayanan, Daniel C Alexander, and Nathan Silber-  
785 man. Learning from noisy labels by regularized estimation of annotator confusion. In *Proceedings  
786 of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11244–11253, 2019.
- 787 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin,  
788 Percy Liang, and Tatsunori B Hashimoto. Alpaca: A strong, replicable instruction-  
789 following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7, 2023.
- 790 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,  
791 Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly  
792 capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 793 Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2, 2024.
- 800 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
801 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
802 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 803 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.  
804 Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv  
805 preprint arXiv:1804.07461*, 2018.
- 806 Haoyu Wang, Zhuo Huang, Zhiwei Lin, and Tongliang Liu. Noisegpt: Label noise detection and  
807 rectification through probability curvature. *Advances in Neural Information Processing Systems*,  
808 37:120159–120183, 2024.

- 810 Hongxin Wei, Lei Feng, Xiangyu Chen, and Bo An. Combating noisy labels by agreement: A joint  
811 training method with co-regularization. In *Proceedings of the IEEE/CVF conference on computer*  
812 *vision and pattern recognition*, pp. 13726–13735, 2020.
- 813 Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text  
814 classification tasks. *arXiv preprint arXiv:1901.11196*, 2019.
- 815  
816 Lai Wei, Zihao Jiang, Weiran Huang, and Lichao Sun. Instructiongpt-4: A 200-instruction paradigm  
817 for fine-tuning minigpt-4. *arXiv preprint arXiv:2308.12067*, 2023.
- 818 Lai Wei, Yuting Li, Kaipeng Zheng, Chen Wang, Yue Wang, Linghe Kong, Lichao Sun, and Weiran  
819 Huang. Advancing multimodal reasoning via reinforcement learning with cold start. *arXiv*  
820 *preprint arXiv:2505.22334*, 2025.
- 821  
822 Mengzhou Xia, Sadhika Malladi, Suchin Gururangan, Sanjeev Arora, and Danqi Chen. Less: Se-  
823 lecting influential data for targeted instruction tuning. *arXiv preprint arXiv:2402.04333*, 2024.
- 824 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and  
825 Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions.  
826 *arXiv preprint arXiv:2304.12244*, 2023.
- 827  
828 Hansi Yang, Quanming Yao, Bo Han, and James T Kwok. Searching to exploit memorization  
829 effect in deep learning with noisy labels. *IEEE Transactions on Pattern Analysis and Machine*  
830 *Intelligence*, 46(12):7833–7849, 2024.
- 831 Leixin Yang and Yu Xiang. Amplify: attention-based mixup for performance improvement and  
832 label smoothing in transformer. *PeerJ Computer Science*, 10:e2011, 2024.
- 833  
834 Jiangchao Yao, Jiajie Wang, Ivor W Tsang, Ya Zhang, Jun Sun, Chengqi Zhang, and Rui Zhang.  
835 Deep learning from noisy image labels with quality embedding. *IEEE Transactions on Image*  
836 *Processing*, 28(4):1909–1922, 2018.
- 837 Jiayi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner  
838 Geyer, Chao Huang, Pin-Yu Chen, et al. Justice or prejudice? quantifying biases in llm-as-a-  
839 judge. *arXiv preprint arXiv:2410.02736*, 2024.
- 840 Soyoung Yoon, Gyuwan Kim, and Kyumin Park. Ssmix: Saliency-based span mixup for text classi-  
841 fication. *arXiv preprint arXiv:2106.08062*, 2021.
- 842  
843 Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical  
844 risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- 845  
846 Jing Zhang, Chao Wang, Zezhou Li, and Xianbo Zhang. Threshold-free anomaly detection for  
847 streaming time series through deep learning. In *2021 20th IEEE International Conference on*  
*Machine Learning and Applications (ICMLA)*, pp. 1783–1789. IEEE, 2021.
- 848  
849 Le Zhang, Zichao Yang, and Diyi Yang. Treemix: Compositional constituency-based data augmen-  
850 tation for natural language understanding. *arXiv preprint arXiv:2205.06153*, 2022.
- 851  
852 Rongzhi Zhang, Yue Yu, and Chao Zhang. Seqmix: Augmenting active sequence labeling via  
853 sequence mixup. *arXiv preprint arXiv:2010.02322*, 2020.
- 854  
855 Wancong Zhang and Ieshan Vaidya. Mixup training leads to reduced overfitting and improved  
856 calibration for the transformer architecture. *arXiv preprint arXiv:2102.11402*, 2021.
- 857  
858 Haoqi Zheng, Qihuang Zhong, Liang Ding, Zhiliang Tian, Xin Niu, Dongsheng Li, and Dacheng  
859 Tao. Self-evolution learning for mixup: Enhance data augmentation on few-shot text classification  
860 tasks. *arXiv preprint arXiv:2305.13547*, 2023.
- 861  
862 He Zhu, Zhiwen Ruan, Junyou Su, Xingwei He, Yun Chen, Wenjia Zhang, and Guanhua Chen. Tag-  
863 instruct: Controlled instruction complexity enhancement through structure-based augmentation.  
*arXiv preprint arXiv:2505.18557*, 2025.
- Zhaowei Zhu, Yiwen Song, and Yang Liu. Clusterability as an alternative to anchor points when  
learning with noisy labels. In *International Conference on Machine Learning*, pp. 12912–12923.  
PMLR, 2021.

## APPENDIX

### A THE USE OF LARGE LANGUAGE MODELS

In this work, we employ ChatGPT-5, a state-of-the-art large language model, to assist with language refinement and clarity improvement. Specifically, ChatGPT-5 is used to polish the writing style, correct grammatical errors, and enhance the overall readability of the manuscript without altering its scientific content or conclusions.

### B CONCEPTUAL DISTINCTION FROM DATA AUGMENTATION AND DATA CURATION

A central goal of our work is to address the pervasive scarcity of high-quality supervision in many NLP scenarios, especially for low-resource or domain-specific tasks. In such settings, large quantities of noisy, redundant, or otherwise low-quality instruction–response pairs are often available, whereas carefully curated high-quality data are expensive and limited. Our instruction distillation framework aims to transform abundant low-quality signals into a compact set of high-quality supervision examples, thereby improving the utility of existing corpora for downstream instruction tuning and supporting the broader low-resource NLP community.

Concretely, the MIXTURE dataset defines five heterogeneous task types. As shown in Sec. 5.1, LM-mixup achieves consistently strong performance across these tasks, indicating that the proposed instruction distillation paradigm is not tied to a single setting or a narrow engineering trick, but rather exhibits robustness and generality across diverse instruction-following scenarios.

Conceptually, our notion of *Instruction Distillation* is distinct from both traditional data augmentation and data curation:

- **Instruction distillation.** Given multiple low-quality or inconsistent responses, instruction distillation extracts the useful information across them and semantically fuses these weak signals into a single, high-quality instruction–response pair. This process simultaneously aggregates information, denoises spurious content, and enforces quality-control constraints, yielding supervision with higher information density for downstream models.
- **Data augmentation.** Classical augmentation techniques generate additional samples via transformations such as rewriting (Wei & Zou, 2019), paraphrasing (Abaskohi et al., 2023), back-translation (Edunov et al., 2018), or other synthetic procedures, with the primary goal of expanding data volume and increasing diversity and robustness (Feng et al., 2021; Chen et al., 2023a). The underlying semantic content of each example is usually preserved, and the number of samples grows.
- **Data curation.** Data curation typically focuses on improving annotation quality or consistency for existing samples without substantially changing their semantic content (Northcutt et al., 2021; Karimi et al., 2020). Examples include relabeling noisy instances, filtering problematic examples, or correcting minor errors while keeping the original instruction–response structure intact.

Instruction distillation fundamentally differs from these two paradigms in both direction and effect. Instead of increasing the number of samples, it *reduces* data volume via information aggregation, while *changing and enriching* the semantic content through fusion across multiple weak sources. Unlike augmentation, which primarily improves diversity, or curation, which mainly refines labels for fixed content, instruction distillation explicitly converts many low-quality signals into a few high-information-density instructions. This enables substantial gains in low-resource regimes, where the key bottleneck is not the absolute number of examples, but the lack of sufficiently rich and reliable supervision.

## C DETAILS OF KNN–BAYES RATING

### C.1 KNN–BAYES QUALITY MODELING WITH SCORE TRANSITION MATRIX

In the Sec.4.2, we introduced an offline KNN–Bayes calibration method to approximate the original LLM ratings during training. Intuitively, given the  $k$ -nearest neighbors of each sample in the embedding space, we aim to infer its “true” quality score based on the observed ratings of these neighbors. However, the LLM-provided scores  $\tilde{y}$  typically suffer from systematic noise and random fluctuations. Directly averaging the neighbor scores may therefore introduce significant bias into the reward signal.

To address this issue, we adopt the classical idea of Score Transition Matrix (STM) from weak supervision and noisy-label learning, which models the conditional distribution between observed and latent labels. Let the latent true label be  $y \in \mathcal{Y} = \{1, 2, \dots, C\}$  and the observed noisy rating be  $\tilde{y}$ . In our implementation, we set  $C = 6$  with label set  $\{0, 1, 2, 3, 4, 5\}$ , which matches the original data annotation. The STM is defined as

$$T \in \mathbb{R}^{C \times C}, \quad T_{ij} = \mathbb{P}(\tilde{y} = j \mid y = i), \quad (9)$$

where  $T_{ij}$  denotes the probability that a true label  $i$  is perturbed into the noisy label  $j$ . The prior distribution is given by

$$p \in \Delta^C, \quad p_i = \mathbb{P}(y = i), \quad \sum_i p_i = 1. \quad (10)$$

When  $T = I$ , the observed ratings are noise-free; deviations of  $T$  from the identity matrix characterize systematic label noise.

**$k$ -NN Clusterability Assumption (Wei et al., 2020).** In the embedding space, if  $x'$  belongs to the  $k$ -nearest neighbors  $\mathcal{N}_k(x)$  of  $x$ , then it is more likely that  $y(x') = y(x)$ . Based on this assumption, the neighborhood agreement frequencies yield a set of linear equations over  $(T, p)$ . We adopt 2-NN consensus statistics when estimating  $(T, p)$  to ensure identifiability. For the posterior computation of a single sample, we use the  $k$ -nearest neighbor histogram  $h(x)$  with  $k \geq 2$  to enhance robustness. Specifically, using pairwise or triplet neighbor agreement, we define

$$v^{[1]} = T^\top p, \quad v_\ell^{[2]} = (T \circ T_\ell)^\top p, \quad v_{\ell,s}^{[3]} = (T \circ T_\ell \circ T_s)^\top p, \quad (11)$$

where  $T_\ell = TA_\ell$  is the cyclic shift of  $T$  by  $\ell$  units, and  $\circ$  denotes the Hadamard product. The observed frequencies  $\hat{v}^{[1]}, \hat{v}_\ell^{[2]}, \hat{v}_{\ell,s}^{[3]}$  can be directly computed from data, forming a linear program over  $(T, p)$ . We solve for  $(T, p)$  subject to  $T\mathbf{1} = \mathbf{1}$ ,  $T \geq 0$ ,  $p \geq 0$ , and  $\mathbf{1}^\top p = 1$ . Existing theory shows that under mild identifiability conditions, third-order consensus vectors suffice to uniquely recover  $(T, p)$ .

Once  $(T, p)$  are estimated, given the empirical neighbor histogram  $h_j(x)$  of sample  $x$ , the posterior distribution is computed as

$$\mathbb{P}(y = i \mid h(x)) \propto p_i \prod_{j \in \mathcal{Y}} T_{ij}^{h_j(x)} = p_i \exp \left( \sum_{j \in \mathcal{Y}} h_j(x) \log T_{ij} \right), \quad (12)$$

Here  $h_j(x) \in \{0, 1, \dots, k\}$  counts the number of neighbors whose observed label equals  $j$ , hence  $\sum_{j \in \mathcal{Y}} h_j(x) = k$ . If distance weights  $w_r$  are used, we replace  $h_j(x)$  by the weighted sum  $\sum_{r: \tilde{y}_r = j} w_r$ . This posterior relies on the conditional independence assumption: given the true label  $y$ , the observed ratings of neighbors are mutually independent. When  $T$  is diagonally dominant (close to  $I$ ), the posterior behavior approaches that of frequency- or average-based voting. If  $T = I$  without any smoothing, however, the likelihood degenerates; thus, we apply mild smoothing to  $T$  and compute in the log domain to ensure numerical stability. The posterior expectation score is

$$\hat{s}(x) = \sum_{i=1}^C i \cdot \mathbb{P}(y = i \mid h(x)). \quad (13)$$

Finally, the quality reward used in training is given by the piecewise mapping

$$R_q(x) = \begin{cases} 1, & \hat{s}(x) \geq 4, \\ 0.3, & 3 \leq \hat{s}(x) < 4, \\ 0, & \text{otherwise.} \end{cases} \quad (14)$$

As  $T$  becomes diagonally dominant (i.e.,  $T \approx I$ ), the posterior concentrates on the most frequent neighbor labels and behaves like smoothed majority/frequency voting.<sup>4</sup> When  $T$  departs from  $I$ , the Bayesian calibration systematically corrects label noise.

## C.2 DETAILS OF KNN-BAYES QUALITY REWARD CONSTRUCTION

To efficiently approximate LLM ratings during training, we introduce a KNN-Bayes scoring system, which leverages neighborhood information and a score transition matrix to denoise label noise. The construction consists of the following steps:

**Offline Asset Construction.** We collect approximately 100K samples rated by ChatGPT-4o-mini as a reference set. We build a KNN index in the embedding space and compute neighbor rating co-occurrence frequencies to estimate both the score transition matrix  $T$  and the label prior  $\mathbf{p}$  offline, prior to model training.

**Online Inference and Reward Computation.** During training, for each generated output  $\hat{Y}$ , we retrieve its  $k$  nearest neighbors in the reference set to form a rating histogram  $\mathbf{h}(\hat{Y}) \in \mathbb{R}^C$ . We then compute the posterior distribution over true labels as

$$P(y = i \mid \mathbf{h}(\hat{Y})) \propto p_i \cdot \exp\left(\sum_j h_j(\hat{Y}) \log T_{ij}\right), \quad (15)$$

and obtain the expected score

$$\hat{s}(\hat{Y}) = \sum_{i=1}^C i \cdot P(y = i \mid \mathbf{h}(\hat{Y})). \quad (16)$$

Finally, rewards are assigned using a piecewise mapping:

$$R_q(\hat{Y}) = \begin{cases} 1 & \hat{s}(\hat{Y}) \geq 4, \\ 0.3 & \hat{s}(\hat{Y}) = 3, \\ 0 & \text{otherwise.} \end{cases} \quad (17)$$

## C.3 CONSISTENCY BETWEEN KNN-BAYES RATING AND LLM SCORES

**Setup.** To evaluate the effectiveness of KNN-Bayes in approximating the original LLM scores, we conduct an offline stratified experiment with a reference set and an evaluation set. Given a dataset  $\mathcal{D}$  with LLM-provided scores  $\tilde{y}$ , we first split it into a reference set  $\mathcal{B}$  and an evaluation set  $\mathcal{A}$  via stratified sampling to preserve the label distribution of  $\tilde{y}$  across both sets. On  $\mathcal{B}$ , we construct a semantic embedding index and estimate the score transition matrix  $T$  and prior distribution  $p$  through neighborhood co-occurrence statistics. For each sample  $x \in \mathcal{A}$ , we retrieve its  $k$ -nearest neighbors in  $\mathcal{B}$ , obtain the empirical histogram  $h(x)$ , and compute the posterior distribution via

$$\mathbb{P}(y = i \mid h(x)) \propto p_i \prod_{j \in \mathcal{Y}} T_{ij}^{h_j(x)}. \quad (18)$$

We then calculate the expected score

$$\hat{s}(x) = \sum_{i \in \mathcal{Y}} i \cdot \mathbb{P}(y = i \mid h(x)). \quad (19)$$

<sup>4</sup>We apply mild Laplace smoothing  $T \leftarrow (1 - \alpha)T + \alpha \mathbf{1}\mathbf{1}^\top / C$  with small  $\alpha > 0$ , followed by row-wise renormalization; computations are carried out in the log domain to avoid underflow.

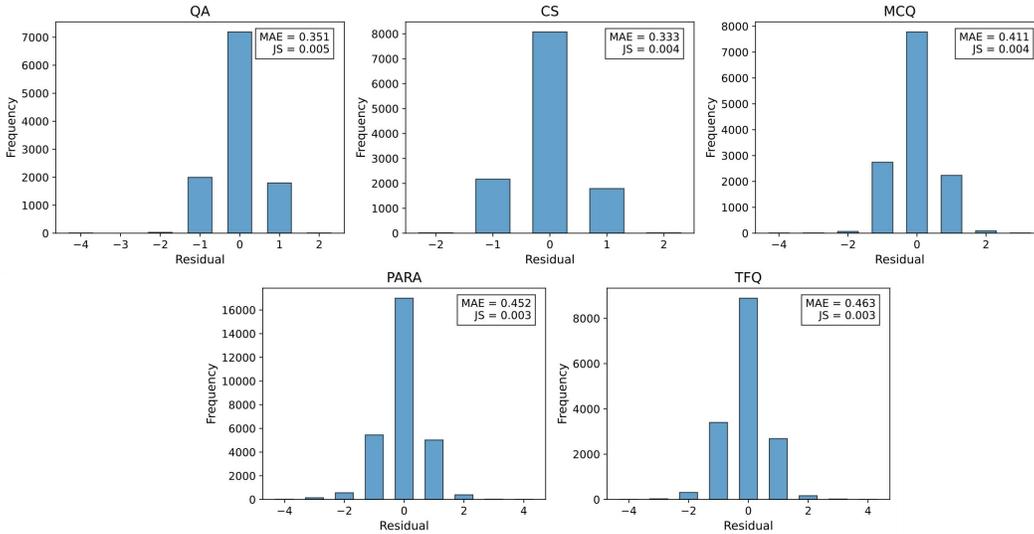


Figure 7: Residual distributions between KNN-Bayes and LLM scores on five datasets. Most residuals concentrate at zero with moderate deviations within  $\{-1, 1\}$ , as shown by MAE and JS metrics.

**Metrics.** We assess the consistency between KNN-Bayes scores and original LLM scores from two perspectives: distributional divergence and numerical deviation. The distributional divergence is measured by the Jensen-Shannon (JS) divergence:

$$JS(P, Q) = \frac{1}{2}D_{KL}(P \parallel M) + \frac{1}{2}D_{KL}(Q \parallel M), \quad M = \frac{1}{2}(P + Q), \quad (20)$$

where  $P$  and  $Q$  denote the empirical distributions of  $\hat{y}(x)$  and  $\tilde{y}(x)$ , respectively, with  $\hat{y}(x) = \text{round}(\hat{s}(x)) \in \mathcal{Y}$  (alternatively, we bin  $\hat{s}(x)$  into the same  $C$  categories). The numerical deviation is quantified using the Mean Absolute Error (MAE):

$$\text{MAE} = \frac{1}{|\mathcal{A}|} \sum_{x \in \mathcal{A}} |\hat{s}(x) - \tilde{y}(x)|. \quad (21)$$

**Results.** Figure 7 shows the residual distributions and MAE/JS metrics across all five datasets. Overall, most residuals concentrate at 0, while some fraction falls within  $\{-1, 1\}$ , suggesting that KNN-Bayes captures the main structure of the original LLM ratings but still exhibits small local deviations. Quantitatively, the JS divergence remains below 0.006 on all datasets, indicating that the calibrated scores preserve the global distributional shape of the LLM scores with minimal shift. The MAE lies in the range 0.33–0.46, which is moderate compared to the discrete rating scale  $\mathcal{Y} = \{0, 1, 2, 3, 4, 5\}$ , reflecting that individual predictions can occasionally deviate by one score level. These findings suggest that while KNN-Bayes provides a low-cost and reasonably accurate approximation for offline evaluation.

## D TRAINING AND EVALUATION DETAILS

### D.1 TRAINING DETAILS

We adopt a three-stage training pipeline: (1) *Cold-start full-parameter tuning*, (2) *GRPO reinforcement learning*, and (3) *Evaluation-stage fine-tuning*. All experiments are conducted on 3 H20 GPUs. The key hyperparameters for each stage are summarized below.

**Cold-start Training.** We first perform full-parameter supervised fine-tuning on the initial dataset to provide a strong initialization for later stages. This stage uses a batch size of 128, learning rate  $2 \times 10^{-5}$ , and runs for 3 epochs with a maximum sequence length of 2048 tokens.

**GRPO Reinforcement Learning.** The second stage adopts GRPO with multi-dimensional reward signals, including Bayesian KNN-based quality scores, BGE-M3 semantic alignment, and format regularization. We set the rollout batch size to 128, actor global batch size to 16, learning rate  $1 \times 10^{-6}$ , KL penalty coefficient  $1 \times 10^{-2}$ , and run for 1 epoch with dynamic batching and gradient checkpointing enabled.

**Evaluation-stage Fine-tuning.** Finally, following prior work (Pang et al., 2024), we perform lightweight LoRA fine-tuning with a rank size of 64 and a scaling factor of 16. We adopt a batch size of 128, a learning rate of  $1 \times 10^{-4}$ , and train for 5 epochs to ensure consistent settings across all evaluation benchmarks.

Table 4: Key hyperparameter settings across three training stages.

Parameter	Cold-start	GRPO	Eval-tuning
Batch size	128	128 (rollout) / 16 (actor)	128
Learning rate	$2 \times 10^{-5}$	$1 \times 10^{-6}$	$1 \times 10^{-4}$
Epochs	3	1	5
Max sequence length	2048	2048	2048
KL penalty	–	$1 \times 10^{-2}$	–
LoRA rank / scaling	–	–	64 / 16
Dropout rate	0.1	0.1	0.1

## D.2 EVALUATION DETAILS

Following previous work (Pang et al., 2024), we evaluate the fine-tuned models on five widely used benchmarks: MMLU (Hendrycks et al., 2020), BBH (Suzgun et al., 2022), GSM8K (Cobbe et al., 2021), TruthfulQA (Lin et al., 2021), and TyDiQA (Clark et al., 2020). For each dataset, we follow standard protocols or common configurations. Specifically, 0-shot settings are used for MMLU; 8-shot in-context examples for GSM8K; 3-shot settings without chain-of-thought for BBH; 6-shot prompts for TruthfulQA; and one in-context example per language for TyDiQA.

## E BASELINE DETAILS

We provide detailed descriptions of all baselines considered in the main experiments:

- **Random Selection:** Randomly selects training samples without any filtering.
- **Completion Length:** Uses the total conversation length as a proxy for data quality, assuming longer completions indicate richer information.
- **Perplexity:** Computes perplexity in a zero-shot manner using a pre-trained model; higher perplexity suggests rarer or more complex samples.
- **k-NN:** Measures average distance to the  $k$  nearest neighbors in the SentenceBERT embedding space to quantify data rarity.
- **LESS** (Xia et al., 2024): Scores samples by their influence on a validation set, estimated via gradient-based metrics.
- **AlpaGasus** (Chen et al., 2023b): Employs an LLM to assign quality ratings, selecting only high-scoring samples.
- **DEITA** (Liu et al., 2023b): Scores samples by both quality and complexity, while iteratively enforcing diversity constraints.
- **DS2** (Pang et al., 2024): Selects high-quality and diverse samples by correcting LLM-generated scores via a transition matrix and combining them with long-tail diversity scores.
- **EDA** (Wei & Zou, 2019): Applies simple text-level perturbations such as synonym replacement, random insertion, random deletion, and word swapping to increase surface-level diversity without altering the original semantic content.

- **Rephrasing** (Abaskohi et al., 2023): Generates semantically equivalent paraphrases of the original instructions using GPT-4o mini, aiming to modify expression style while preserving meaning to introduce natural linguistic variation.
- **Back-translation** (Edunov et al., 2018): Translates each sample into an intermediate language and back to the source language, producing paraphrastic variants that expose the model to diverse lexical and syntactic forms.
- **Full Data**: Uses the entire dataset without any filtering for model fine-tuning.

For all rating-based methods (*AlpaGasus*, *DEITA*, and *DS2*), we follow *LM-Mixup* and adopt ChatGPT-4o-mini as the rating model for a fair comparison.

For all data augmentation baselines, the augmentation operations are typically applied to the entire dataset, whereas our other baselines are constructed using 10K training samples. To ensure a fair comparison, we first randomly sample 5K instances from the original 300K data pool and then apply the corresponding augmentation method to these 5K samples, resulting in a total of 10K training examples.

## F ADDITIONAL DETAILS OF DATA PREPARATION AND EXPERIMENTAL SETTINGS

In our experiments, all 300K samples are annotated with a quality score ranging from 1 to 5. Among them, approximately 30K samples with scores  $\geq 4$  are treated as high-quality data, while roughly 270K samples with scores  $< 4$  constitute the low-quality pool. Below we provide additional clarifications for the datasets used in Table 3.

**ORI (High-Quality Data).** From the 30K high-quality samples, we compute long-tail diversity scores and select the top  $N$  instances to form the ORI training set. These samples serve as the high-quality-only baseline.

**LOW (Low-Quality Data).** To construct the LOW baseline, we randomly sample  $N$  instances from the 270K low-quality pool. This setting evaluates the performance of directly using weak signals without any enhancement.

**BASE (Qwen-1.5B-Instruct Generated Data).** We directly apply the off-the-shelf Qwen-1.5B-Instruct model (without any training or fine-tuning) to fuse inputs from the low-quality pool. The generated outputs constitute the BASE dataset.

**MIXUP (LM-Mixup Generated Data).** The MIXUP dataset is constructed using the following pipeline:

1. Sample multiple groups of low-quality inputs from the 270K data pool.
2. Apply LM-Mixup to each group of  $n$  inputs to generate a smaller set of fused, high-quality candidates.
3. Compute long-tail diversity scores for all mixup-generated outputs.
4. Select the top 5K samples to form the final MIXUP high-quality dataset.

This process allows LM-Mixup to aggregate and refine weak signals into information-dense, high-quality supervision data.

## G MORE EXPERIMENT RESULTS

**Full Low-Quality vs. Full Mixup Data.** To comprehensively evaluate the effectiveness of our approach, we conduct experiments on three representative models—Mistral-7B, LLaMA-3.1-8B, and LLaMA-2-7B-hf—using 10K samples drawn respectively from the raw low-quality dataset and the mixup-enhanced dataset generated via *LM-Mixup*. As shown in Figure 8, directly fine-tuning

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

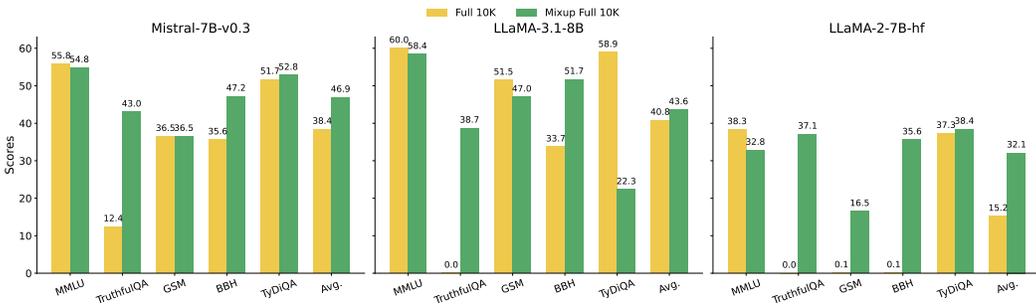


Figure 8: Comparison of model performance on five benchmarks using the full low-quality dataset versus the full mixup dataset (10K samples). Scores are reported for MMLU, TruthfulQA, GSM, BBH, TyDiQA, and the overall average.

on low-quality data leads to inconsistent and often suboptimal performance across most benchmarks. In contrast, the mixup-enhanced data substantially boosts performance on key tasks such as TruthfulQA, BBH, TyDiQA, and the overall average score for all three models. Notably, the improvements are most pronounced on LLaMA-2-7B-hf, where the baseline performance on raw data is particularly low, highlighting the robustness of *LM-Mixup* in challenging low-quality settings. These results collectively demonstrate that our method consistently transforms low-quality samples into a valuable resource for instruction tuning, unlocking their potential and significantly narrowing the gap with high-quality data baselines.

**Additional Results on LLaMA-2-7B-hf and Mistral-7B-v0.3.** We additionally conducted experiments to assess the performance of the OpenLLM leaderboard across different baseline settings using various backbone models, including Mistral-7B-v0.3 and LLaMA-2-7B-hf. Tables 5 and 6 report the corresponding results for these two backbones, respectively. Overall, the findings further confirm the effectiveness of our method, demonstrating that with appropriate configurations, it can consistently achieve top-2 performance on the leaderboard.

Table 5: Results on the OpenLLM leaderboard using Mistral-7B-v0.3 as the base model. The top-performing scores are shown in **bold**, while the second-best scores are marked with underlines. \* indicates that the values are sourced from Pang et al. (2024).

Models	MMLU (factuality)	TruthfulQA (truthfulness)	GSM (reasoning)	BBH (reasoning)	TydiQA (multilinguality)	Average
VANILLA BASE MODEL*	59.7	30.2	38.0	49.6	54.9	46.5
COMPLETION LENGTH*	58.9	34.4	42.5	53.1	59.6	49.7
PERPLEXITY*	59.8	40.3	36.0	48.9	57.4	48.5
<i>k</i> -NN-10*	58.3	41.7	43.5	54.1	53.4	50.2
RANDOM SELECTION*	59.4	36.7	41.8	54.2	54.0	49.3
LESS*	59.5	34.8	42.0	54.5	57.5	49.7
FULL DATA (300K)*	60.0	43.5	43.5	52.5	53.4	50.6
ALPAGASUS*	60.5	36.7	41.0	55.1	57.3	50.1
DEITA*	60.1	35.6	40.5	55.1	56.0	49.5
DS2 w/o CURATION*	60.1	35.9	48.5	54.2	58.9	51.5
DS2*	59.9	37.9	47.5	55.6	59.3	<b>52.0</b>
MIXUP 70% + ORI 30%	58.5	42.7	46.0	53.2	52.9	50.7
MIXUP 50% + ORI 50%	57.0	43.0	47.0	54.0	52.6	50.7
MIXUP 30% + ORI 70%	56.0	45.3	51.5	54.0	52.1	<u>51.8</u>

## H ADDITIONAL DATASET STATISTICS ABOUT MIXTURE

Figure 9 reports the distribution of the number of low-quality variants constructed for each high-quality sample. Most samples are paired with multiple degraded variants, enabling the model to learn hierarchical mappings from noisy or incomplete inputs to high-quality outputs.

Table 6: Results on the OpenLLM leaderboard using LLaMA-2-7B-hf as the base model. The top-performing scores are shown in **bold**, while the second-best scores are marked with underlines. \* indicates that the values are sourced from Pang et al. (2024).

Model	MLLU (factuality)	TruthfulQA (truthfulness)	GSM (reasoning)	BBH (reasoning)	TydiQA (multilinguality)	Average
VANILLA LLAMA-2-7B*	41.9	28.4	6.0	38.3	35.7	30.1
COMPLETION LENGTH*	42.4	36.4	1.5	36.8	33.9	30.2
PERPLEXITY*	45.0	41.5	12.0	31.7	39.5	33.9
k-NN-10*	38.2	40.8	15.0	36.0	43.8	34.8
RANDOM SELECTION*	44.7	41.8	14.0	37.9	40.8	35.8
LESS	44.3	38.2	18.0	35.2	46.3	36.4
FULL DATA (300K)*	50.1	36.2	16.5	40.5	46.7	38.0
ALPAGASUS*	45.3	41.0	14.5	37.0	45.3	<u>36.6</u>
DEITA*	45.2	44.7	13.5	35.6	43.4	36.5
DS2 w/O CURATION*	42.0	39.5	15.0	38.1	46.1	36.1
DS2*	40.2	43.8	13.5	38.9	46.5	<u>36.6</u>
MIXUP 70% + ORI 30%	39.5	42.5	17.0	38.6	42.5	36.0
MIXUP 50% + ORI 50%	39.5	44.0	18.0	38	42.5	36.4
MIXUP 30% + ORI 70%	39.0	45.3	18.0	38.0	43.5	<b>36.8</b>

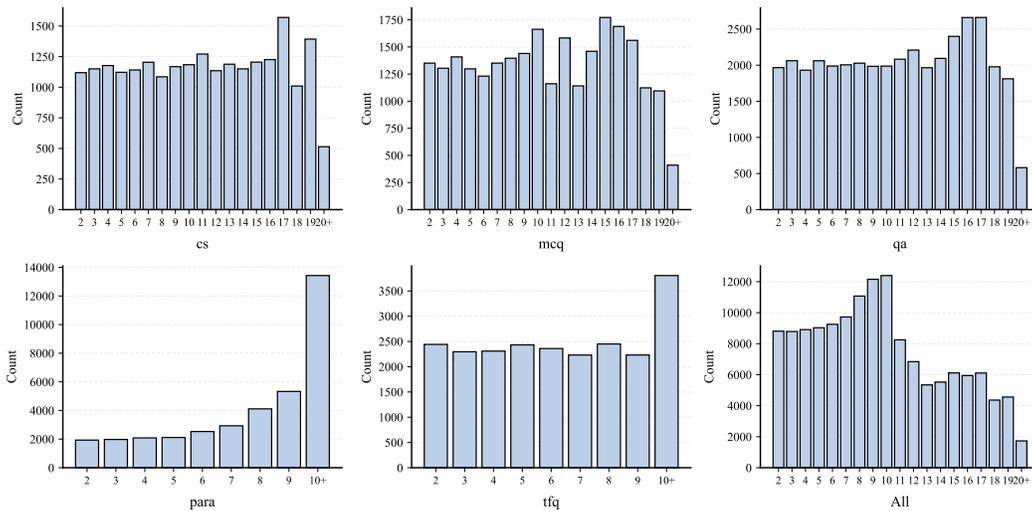


Figure 9: Distribution of the number of low-quality samples derived from each high-quality sample across different task categories.

Table 12 provides a detailed breakdown of the entire dataset across five task types (QA, MCQ, CS, TFQ, Paragraph) and three data variants (Normal, Cross-Topic, Noisy). We observe a balanced distribution across task types, with QA and Paragraph slightly larger in size, ensuring diverse coverage for training and evaluation.

## I PROMPT TEMPLATE

The prompt template below illustrates how we use ChatGPT-4o-mini to generate high-quality data, with the following notes for different task types: For qa, let Instruction be the question and Output be the answer. For mcq, include options inside Instruction and provide the correct choice and a brief rationale in Output. For cs, tfq, or paragraph styles, keep Instruction as the task prompt and Output as the targeted response.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

### Prompt Template for High Quality Data Generation

You are a knowledgeable assistant tasked with producing **exceptionally high-quality** `{task_type}` instances that will later be rated on four axes: *Rarity*, *Complexity*, *Informativeness*, and *Overall*.

#### Scoring Criteria

- **Rarity (1–10):** Cover non-obvious, less-quoted aspects; avoid commonplace trivia.
- **Complexity (1–10):** Require synthesis of multiple facts, causal/temporal links, or non-trivial reasoning; avoid single-sentence lookups.
- **Informativeness (1–10):** Deliver dense, relevant, non-trivial content—even if concise; add value beyond superficial recall.
- **Overall (1–10):** Aggregate impression; aim for the top tier when justified.

#### Requirements

1. Generate **3–4** high-quality `{task_type}` instances based on the passage below.
2. For each instance, explicitly maximize the four criteria above (Rarity, Complexity, Informativeness, Overall).
3. You may quote or paraphrase the passage, but *weave* information to show reasoning and uniqueness; avoid verbatim copying when unnecessary.
4. Length is flexible—prioritize informativeness and reasoning over verbosity.
5. Use **plain text only** in the following exact format (no markdown):

Instruction: <instruction 1>

Output: <output 1>

Instruction: <instruction 2>

Output: <output 2>

Instruction: <instruction 3>

Output: <output 3>

#### Passage

{passage}

The prompt template below illustrates how we use ChatGPT-4o-mini to generate low quality data, with the following notes for different task types: For qa, let Instruction be the question and Output be the answer. For mcq, include options inside Instruction and provide the correct choice with minimal explanation in Output. For cs, tfq, or paragraph styles, keep Instruction as the task prompt and Output as the response, ensuring only moderate relevance and detail. The number n is determined by a random value.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403

#### Prompt Template for Low Quality Data Generation

You are a moderately skilled assistant tasked with producing **{n} low-quality {task\_type}** instances derived from the given original data. Each generated instance should deliberately reflect *average quality*, aiming for scores between **3–6 (on a 10-point scale)** on four evaluation axes: *Rarity*, *Complexity*, *Informativeness*, and *Overall*.

#### Scoring Criteria

- **Rarity (3–6):** Cover reasonably common aspects; avoid both overly trivial and highly novel content.
- **Complexity (3–6):** Allow some light inference but avoid deep reasoning or multi-step logic.
- **Informativeness (3–6):** Ensure answers are mostly correct but lack nuance, depth, or precision.
- **Overall (3–6):** The overall impression should feel average, slightly useful, somewhat generic, and unpolished.

#### Hard Constraints

1. **Same Topic:** All {task\_type} instances must stay on the identical topic as the original.
2. **Explicit Subject:** The main subject or event (e.g., *the execution of Turner and McDaniel*) must be stated verbatim in every question and answer. Avoid vague pronouns unless the noun is immediately repeated.
3. **Self-contained:** Each {task\_type} must be understandable in isolation; assume no external context.
4. **No Off-topic Content:** Do not introduce unrelated domains or shift the factual focus.

#### Output Requirements

1. Generate **3–4** {task\_type} instances based on the original data.
2. Maintain moderate quality (scores 3–6) on all four evaluation axes.
3. Use **plain text only** in the exact format below (no markdown):

Instruction: <instruction 1>  
 Output: <output 1>

Instruction: <instruction 2>  
 Output: <output 2>

Instruction: <instruction 3>  
 Output: <output 3>

#### Original Data

{orig}

The prompt template below illustrates how we use ChatGPT-4o-mini to perform data fusion across different task types, with the following notes: For qa, merge two question–answer pairs into a single, coherent question with a unified answer. For mcq, combine two multiple-choice questions into one integrated question, providing a single correct option with a concise explanation. For cs, tfq, or paragraph tasks, merge the content of both instances into a single prompt–response pair, ensuring the output reflects a natural synthesis of the original information while maintaining moderate length and relevance.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

#### Prompt Template for Data Fusion across {task\_type}

You are a data fusion expert tasked with merging two {task\_type} instances into a single, coherent, and high-quality {task\_type} instance. The goal is to synthesize both original samples into one unified output that naturally connects the information from both inputs while maintaining high quality.

#### Fusion Requirements

1. Merge the two original {task\_type} samples into **one concise, integrated instance**.
2. The unified output must address both original inputs comprehensively while avoiding redundancy or contradiction.
3. The fusion should capture *subtle conceptual links*, rather than simply stacking facts together.
4. Ensure the final output meets the following quality criteria:
  - **Rarity:** Avoid overly common or trivial facts; focus on non-obvious insights.
  - **Complexity:** Encourage nuanced reasoning or implicit connections.
  - **Informativeness:** Maximize factual density and relevance.
  - **Overall Quality:** Aim for the top tier across all above dimensions.

#### Output Format

Instruction: <your merged instruction>  
Output: <your merged output>

#### Original Instances

Instance-1:  
{text1}  
  
Instance-2:  
{text2}

Following previous work (Pang et al., 2024), we use the same template for LLM Rating:

#### Prompt Template for LLM Rating

As a data quality estimator, your task is to assess the quality of the data sample based on the criteria: Rarity, Complexity, and Informativeness. Please rate the sample on a scale from 1 to 10 for each criterion, and return an overall rating on a scale from 1 to 10, where a higher score indicates a higher level of quality. Ensure that the ratings are not overly concentrated around a specific score. If multiple samples have similar qualities, consider spreading the scores more evenly to reflect subtle differences.

Please carefully evaluate the following data sample and return the integral evaluation scores using the JSON format:

```
{"Rarity": <number, 1-10>,
  "Complexity": <number, 1-10>,
  "Informativeness": <number, 1-10>,
  "Overall rating": <number, 1-10>}
```

Instruction: [Instruction]  
Input: [Input]  
Response: [Response]

## 1458 J CASE STUDY

1459

1460 Tables 7, 8, 9, 10, and 11 present representative mixup cases across QA, CS, TFQ, MCQ, and para-  
 1461 graph tasks, illustrating how multiple fragmented inputs are integrated into coherent, high-quality  
 1462 outputs.

1463

1464

Table 7: QA pairs mixup example.

1465

1466 **Input (Multiple QA Pairs)**

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

Q: What are the financial implications Hardy discussed with the AAHA and Saskatchewan Amateur Hockey Association regarding the discussions on player development?

A: Hardy’s discussions highlighted that financial implications include a strain on resources, showing the need for more investment in player retention to prevent losing talent to professional leagues.

Q: How might the discussions with the AAHA and Saskatchewan Amateur Hockey Association on player development affect the competitiveness of amateur hockey?

A: The discussions suggest that amateur hockey may face challenges to its competitiveness if it cannot adapt to the pressures from professional leagues, potentially leading to a decline in player quality.

Q: What strategic changes were suggested in Hardy’s discussions with the AAHA and Saskatchewan Amateur Hockey Association regarding player development?

A: Hardy’s discussions prompted the idea of re-evaluating investment strategies in player development to ensure that amateur hockey remains appealing and retains its talent pool amidst increasing professional competition.

Q: What broader challenges were mentioned by Hardy concerning the discussions on player development?

A: Hardy mentioned that broader challenges include maintaining the integrity of amateur hockey and ensuring grassroots involvement, which are at risk due to the rising influence of professional leagues.

1485

1485 **Output (Merged Answer)**

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

Q: Analyze how Hardy’s post-WWII discussions with amateur hockey associations reflected emerging tensions in Canadian sports culture, particularly regarding financial constraints, competitiveness, and societal identity shifts.

A: In examining Hardy’s dialogue with prominent figures representing national governing bodies post-WWII, there emerges a profound intersectionality of personal ambition juxtaposed against systemic constraints faced by athletes transitioning towards professional careers. His candid remarks reveal underlying concerns about the sustainability of amateurism given escalating demands placed upon young talents competing internationally. Such sentiments reflect wider cultural shifts occurring globally wherein individuals sought greater recognition beyond localized boundaries—a trend mirrored in other sporting disciplines experiencing rapid modernization. Furthermore, Hardy’s reflections hint at deeper socio-political ramifications tied to identity formation linked to sport participation, suggesting that individual aspirations intersected critically with collective narratives shaped by historical events affecting community cohesion. Overall, these discussions underscore significant transformations happening simultaneously in society and athletics, revealing complexities inherent in balancing tradition with contemporary realities facing aspiring amateurs navigating pathways toward professional success.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

Table 8: CS pairs mixup example.

Input (Multiple CS Pairs)
C: Historical evolution of electromagnetic theory
S: The historical evolution of electromagnetic theory involves key figures like Maxwell who explored the aether theory, which proposed that forces act through a medium.
C: Philosophical implications of electromagnetic theory
S: The philosophical implications of electromagnetic theory reflect a shift in understanding forces, moving from the aether concept to recognizing the nature of fields and potentials in physics.
C: Impact of electromagnetic theory on engineering
S: The impact of electromagnetic theory on engineering principles is substantial, as it provides a framework for understanding different physical properties compared to concepts like potential energy.
C: The role of experiments in electromagnetic theory
S: The role of experiments in electromagnetic theory, particularly the Michelson–Morley experiment, played a crucial part in dismissing the aether theory and reshaping modern physics.
Output (Merged Answer)
C: Synthesis of Physics Concepts
S: The exploration of electromagnetic phenomena reflects significant advancements in classical mechanics, wherein the transition from Newtonian dynamics towards relativistic formulations necessitated a paradigmatic change akin to shifting away from traditional force-based models toward field-theoretic approaches rooted in differential forms. While early attempts at explaining electrostatics involved complex interactions mediated via hypothetical entities known as electric fluids, later thinkers sought simpler explanations grounded in mathematical constructs involving vector calculus. Notably, the dismissal of the aether hypothesis postulated by Fresnel was pivotal, leading physicists to embrace alternative interpretations supported by extensive experimental validation, including notable outcomes observed in the Michelson–Morley interferometer experiment. These evolving understandings illustrate profound intersections between philosophy and science, revealing deeper layers of reality beyond conventional boundaries.

Table 9: TFQ pairs mixup example.

Input (Multiple TFQ Pairs)
Joanne Peason’s view of Pryce’s essay does not show an important moment in the debate about Christianity’s origins in Britain.
Answer: False
Joanne Peason believes that Pryce’s essay fails to effectively challenge established legends.
Answer: False
Pryce’s work is considered by Joanne Peason to have no relevance to discussions of Welsh identity and Christianity.
Answer: False
Output (Merged Answer)
Joanne Pearson’s assertion that Pryce’s essay represents “a significant point” in challenging prevailing myths suggests she recognizes the importance of questioning traditional accounts of early Christianization processes in Britain, especially relating to Welsh culture. Her belief indicates a shift towards more nuanced interpretations where local identities were intertwined with religious conversion stories during periods like Roman rule. Such skepticism contributes significantly to ongoing debates over national histories shaped by colonial influences, suggesting that Pryce’s exploration may serve dual purposes: critiquing existing narratives and fostering contemporary dialogues about indigenous rights and heritage preservation amidst modernity.
Answer: False

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

Table 10: MCQ pairs mixup exmample.

**Input (Multiple MCQ Pairs)**

Q: What is a possible effect of contestants giving “Bonus” points to each other throughout the show?

- A) It might lead to an increased sense of competition among contestants.
- B) Contestants may use Bonus points to create alliances, improving team dynamics.
- C) The Bonus points could make contestants less motivated to help one another.
- D) Contestants might ignore Bonus points and only care about the judges’ scores.

Answer: B

Q: How do Bonus points awarded by contestants impact their relationships?

- A) They may cause contestants to form rivalries that hurt their performance.
- B) Contestants could leverage Bonus points to build alliances and boost collaboration.
- C) The use of Bonus points may lower motivation among contestants.
- D) Contestants generally overlook Bonus points, focusing only on judges’ feedback.

Answer: B

Q: In the context of the show, how might the Bonus points affect contestant interactions?

- A) They could lead to increased competition and less teamwork.
- B) Contestants might use Bonus points to create positive alliances, enhancing performance.
- C) The Bonus points system could demotivate contestants from giving constructive feedback.
- D) Contestants may disregard Bonus points, concentrating on judges’ assessments instead.

Answer: B

Q: What impact do Bonus points have on how contestants relate to one another?

- A) They might foster rivalries that detract from overall performance.
- B) Contestants may utilize Bonus points to form alliances, which can enhance their teamwork.
- C) The Bonus points system could lead to a lack of motivation among the contestants.
- D) Contestants often pay little attention to Bonus points, focusing more on judges’ scores.

Answer: B

**Output (Merged MCQ Instance)**

Q: Analyzing the psychological dimensions involved in the game structure reveals insights into how power operates differently than traditional authority figures. Consider the role of contestants awarding “Bonus” points to each other over time. Discuss the broader implications of this mechanism on perceptions of fairness and trustworthiness among competitors. Which statement best captures the dual nature of this dynamic?

- A) Awarding Bonus points fosters genuine camaraderie among contestants, leading to improved performances due to mutual support rather than fear-based compliance typical of hierarchical setups.
- B) By allowing contestants to give Bonus points freely, the hierarchy becomes irrelevant, suggesting that true leadership emerges solely from meritocratic principles regardless of initial status differences.
- C) The ability to manipulate others through Bonus points creates distrust among teammates, undermining collaborative efforts necessary for success, thereby reinforcing existing hierarchies despite attempts to level playing fields.
- D) Giving Bonus points reflects a shift towards egalitarianism, enabling weaker members to challenge dominant personalities directly, ultimately reshaping societal norms beyond just entertainment purposes.

Answer: C

Table 11: Paragraphs mixup example.

Input (Four Separate Paragraphs)	
1.	The Industrial Revolution transformed urban landscapes across Europe, bringing rapid technological advancements but also leading to overcrowding and unsanitary living conditions in many cities.
2.	During the same period, public health reforms emerged as governments recognized the need for clean water systems, waste management, and regulations to combat disease outbreaks.
3.	Architects and city planners began integrating green spaces and wider streets into urban designs, seeking to improve both public health and aesthetic appeal.
4.	These efforts laid the groundwork for modern urban planning, where functionality, sustainability, and livability became central considerations in shaping city environments.
Output (Merged High-Quality Paragraph)	
The Industrial Revolution reshaped European cities through rapid industrialization, introducing new technologies but also creating overcrowded and unhealthy living conditions. As governments grappled with recurring disease outbreaks, public health reforms prioritized clean water systems, sanitation, and housing regulations. Urban planners responded by incorporating green spaces, broader streets, and organized infrastructure, aiming to balance industrial growth with livable environments. These combined efforts forged the foundations of modern urban planning, where health, sustainability, and functionality guide the evolution of contemporary cities.	

## K DATASET STATISTICS

Table 12: Dataset statistics across five task categories. For each category, we report the total number of samples (train+test) and their breakdown into noisy, cross-Topic, and normal subsets. Ratios indicate the percentage of the full dataset.

Category	Samples	Ratio (%)
<b>I. cs</b>	<b>22,012</b>	<b>15.2</b>
• Noisy	1,100	0.8
• Cross-Topic	4,852	3.3
• Normal	16,060	11.1
<b>II. mcq</b>	<b>25,437</b>	<b>17.6</b>
• Noisy	1,272	0.9
• Cross-Topic	4,714	3.3
• Normal	19,451	13.4
<b>III. qa</b>	<b>38,455</b>	<b>26.6</b>
• Noisy	1,922	1.3
• Cross-Topic	7,405	5.1
• Normal	29,128	20.1
<b>IV. para</b>	<b>36,423</b>	<b>25.1</b>
• Noisy	1,821	1.3
• Cross-Topic	8,261	5.7
• Normal	26,341	18.2
<b>V. tfq</b>	<b>22,557</b>	<b>15.6</b>
• Noisy	1,128	0.8
• Cross-Topic	4,154	2.9
• Normal	17,275	11.9
<b>All</b>	<b>144,884</b>	<b>100.0</b>