ZEBRA: Leveraging Model-Behavioral Knowledge for Zero-Annotation Preference Dataset Construction

Jeesu Jung¹, Chanjun Park^{2*}, Sangkeun Jung^{1*}

¹Chungnam National University, ²Soongsil University jisu.jung5@gmail.com, chanjun.park@ssu.ac.kr, hugmanskj@gmail.com

Abstract

Recent efforts in LLM alignment have focused on constructing large-scale preference datasets via human or Artificial Intelligence(AI) annotators. However, such approaches rely on instance-wise supervision, incurring substantial annotation cost and limited interpretability. In this paper, we propose **ZEBRA**—a model behavior-wise zero-annotation framework that constructs preference data by leveraging model behavior knowledge derived from benchmark performances.

ZEBRA binarizes response pairs by evaluating the quality and similarity of their origin models, entirely bypassing instance-level annotation. This allows scalable, controllable, and cost-effective alignment data generation. Empirical results show that ZEBRA achieves alignment performance comparable to instance-supervised methods, despite requiring no manual or model-based labeling.

1 Introduction

Aligning large language models (LLMs) with human preferences is an essential step toward making them both useful and safe. A common way to achieve this is through instance-wise labeling, where pairs of model responses are compared one by one to see which is better. Well-known methods like Reinforcement Learning from Human Feedback (RLHFOuyang et al. (2022a)) and Artificial Inteligence(AI)-based labeling(RLAIFLee et al.) often use this strategy.

However, *instance-wise* labeling faces two major challenges. First, it is very costly, whether it involves human annotators or additional computational resources for LLM-based labeling(Zhang et al., 2024; Zheng et al., 2023). Second, it lacks a global view of the model's behavior(Ji et al., 2023). Since each response pair is judged in isolation, it

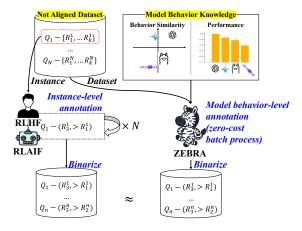


Figure 1: Human annotation(Ouyang et al., 2022a) and AI labeler-based RLAIF(Cui et al., 2023) use *instance-wise* annotations, while ZEro-annotation Behavior-based Response Alignment(ZEBRA) applies *model behavior-wise* annotation based on **model behavioral knowledge**, which captures proficiency and similarity across models. M is the model, R_i is the generated response from M_i about Instruction Q. N is the number of preference dataset and K is the number of models

is difficult to consider broader factors. For example, whether fluency should outweigh factual accuracy. Or whether a model's outputs are consistently aligned with certain policies(Wang et al., 2025). This can lead to labeling noise, mistakes, and limited interpretability.

To address these limitations, we propose a new preference binarization approach called **ZEro-annotation Behavior-based Response Alignment** (**ZEBRA**) (Figure 1). The main idea is: (1) extract each model's behavioral patterns from its past performance trajectories, (2) measure and compare these behaviors in terms of model strength or similarity, and (3) assign preferences at the *model* level rather than for each individual response pair.

We implement this idea through three key components. First, we define **Model Behavior Knowledge (MBK)** for LLMs (discuss in Section 4.1). Second, we propose a way to quantify and collect MBK from objective data sources such as bench-

^{*}Equal contribution. Corresponding author

mark performance. Third, we introduce three strategies—based on *superiority*, *similarity*, and a *hybrid*—to construct a binarization dataset in a zero-annotation, cost-free manner.

A major advantage of our approach is that it creates response pairs based on model superiority without any additional human or LLM labeling. By classifying models with higher benchmark scores as "positive" and those with lower scores as "negative," we can systematically label preferences throughout the dataset. This significantly reduces the cost of annotation and, because MBK visualizes each model's behavior pattern, increases the interpretability of the preference decisions.

Through extensive experiments, we show that **ZEBRA** achieves performance comparable to existing instance-wise labeling methods in the Ultrafeedback(Cui et al., 2023) dataset—without any extra labeling cost.

In summary, our contributions are as follows:

- We introduce ZEBRA, a zero-annotation alignment framework that determines preferences from quantified model-level behavior, bypassing instance-level supervision.
- We demonstrate that the *Model Behavior Knowledge (MBK)* from benchmark performance offers alignment signals comparable to instance-wise labeling.
- We empirically show that ZEBRA matches the performance of established methods such as RLHF and RLAIF, yet requires no additional labeling cost.

2 Preliminaries

2.1 Instance-Level Preference Construction

Most existing preference learning frameworks for LLM alignment—such as Reinforcement Learning from Human Feedback(RLHF, Ouyang et al. (2022a)) and AI-generated preference methods (RLAIF, Lee et al.)—rely on **instance-level pairwise supervision**. Given an instruction x, multiple candidate responses $\{r_1, r_2, ..., r_k\}$ are scored or ranked by either human annotators or automated scoring models. This generates preference tuples $(x, r_i \succ r_j)$, where r_i is preferred over r_j under some evaluation criteria (e.g., helpfulness, truthfulness, coherence).

The preference construction typically involves:

- Generating multiple responses per instruction using different models or decoding strategies.
- Computing preference labels via either human judgment or model-based scoring (e.g., GPT-4(OpenAI, 2023)).
- Aggregating these labels into a pairwise dataset for alignment tuning.

2.2 Challenges of Instance-Level Supervision

While instance-level supervision has proven effective in aligning LLMs, it remains costly, noisy, and difficult to scale. Despite its popularity, instance-level supervision suffers from three core limitations:

- 1. **Costly evaluation**: Annotating or scoring each response pair requires substantial human or computational effort.
- 2. **Preference triviality**: When candidate responses differ significantly in quality, the preference label becomes trivial, contributing little to alignment learning.
- 3. **Instruction-level variance**: Difficulty and ambiguity in the instruction x can introduce noise into the preference signal, especially for automated labelers.

These limitations underscore the need for alternative approaches that construct preference signals without relying on per-instance scoring.

2.3 Motivation for Model Behavior-Level Preference

We propose that instead of relying on per-instance evaluation, one can leverage the **intrinsic capabilities of the response-generating models** themselves. As models exhibit differences in core competencies measurable via standardized benchmarks, we hypothesize this can replace instance-level annotations.

3 ZEro-annotation Behavior-based Response Alignment Framework

3.1 Instance-level Annotation vs. Model behavior-level Annotation

Traditional instance-level binarization relies heavily on detailed, provided by human or AI annotators(Zhang et al., 2024; Sharma et al., 2024). Each pairwise comparison demands significant effort and

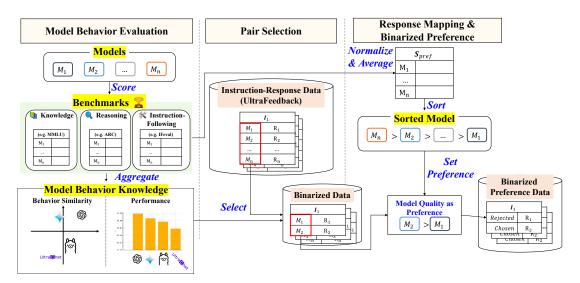


Figure 2: Editing process of ZEBRA framework. n is the number of models. $\{M_1, M_2, ..., M_n\}$ represents the total set of candidate models. $\{I_1, I_2, ..., I_n\}$ represents the instruction input of dataset. $\{R_1, R_2, ..., R_n\}$ represents the total responses from the candidate models. R_x and R_y are the chosen and rejected responses. S_{pref} score summarizes the model's benchmarks ability

resources to maintain consistency and interpretability. Such an approach often results in annotation noise, limited scalability, and considerable expense.

In contrast, our proposed ZEBRA framework leverages intrinsic behavioral knowledge derived from model performance across various benchmarks, ZEBRA systematically matches and pairs responses. Responses from models with proven higher competencies form positive labels, while those from models with lower competencies become negative labels. This innovative model-level approach drastically reduces annotation costs, minimizes noise, and enhances scalability. Additionally, the behavior similarity-based matching provides explicit control over the difficulty and nuance of preference comparisons, leading to clearer and more meaningful alignment outcomes.

3.2 Model Behavior Knowledge

While most preference-learning pipelines focus on differences between individual responses, our approach highlights Model Behavior Knowledge (MBK)—the comprehensive record of each model's past behaviors and capabilities. We define MBK using two sets of metrics:

- **Superiority:** How much better (or worse) a model is compared to others, based on overall or task-specific proficiency.
- **Similarity:** How likely a model is to behave similarly to other models.

These metrics provide a principled basis for

quantifying both a model's general strength and its behavioral proximity to its peers. Within a preference-learning pipeline, *superiority* functions as a global preference signal: when one model consistently outperforms another across standardized benchmarks, its responses are considered preferable overall.

Conversely, behavioral similarity facilitates the systematic construction of challenging comparison sets. When two models are behaviorally similar—for example, they exhibit comparable reasoning performance—their responses become difficult to distinguish. Training on such hard-to-distinguish pairs guides the preference learner to focus on subtle qualitative differences, resulting in more nuanced and robust alignment.

3.3 Model Behavior Evaluation using Benchmark Performances

To capture MBK in a practical, objective way, we rely on external benchmark performance data for each LLM. Many models are already evaluated across diverse, standardized tasks (e.g., reasoning, factual accuracy, instruction-following). We aggregate these benchmark scores to form each model's MBK profile.

Benchmark performance offers several advantages for extracting MBK:

- It provides *reliable metrics* for the core competencies of large language models.
- It enables *straightforward comparison and aggregation* across multiple models.

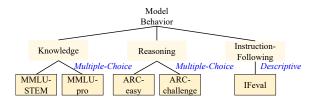


Figure 3: Selected benchmarks for evaluating the models behavior knowledge: knowledge, reasoning, and instruction-following capabilities

 Because benchmark scores are published and fixed at release, they impose no constraints on subsequent data or model expansion.

The list of benchmarks used in our analysis is provided in Figure 3.

Figure 4 illustrates how the behavioral similarities inferred from benchmark performance can reflect actual model similarities. For example, the LLaMA-2-13b model exhibits a pattern closely resembling that of its 7b counterpart. In contrast, WizardLM-7b demonstrates a markedly different behavioral trajectory.

We define a model's ability vector $v_i \in \mathbb{R}^m$ across m benchmark tasks:

$$v_i = \left[s_i^{(1)}, s_i^{(2)}, \dots, s_i^{(m)} \right],$$

where $s_i^{(b)}$ is the normalized score of model M_i on benchmark b, reflecting its relative capability in a specific behavioral dimension (e.g., knowledge, reasoning, instruction-following). These standardized behavior vectors serve as the foundation for both quality-based and similarity-based anchoring, enabling **zero-annotation binarization** of preference pairs without per-instance supervision.

3.4 Benchmark-based Model Behavior Ouantification

To effectively binarize preference data, it is essential to quantify model behavior from two complementary perspectives: behavior quality and behavior similarity. ZEBRA leverages benchmark-derived measures of these aspects to systematically pair positive responses with suitable negative counterparts. By quantifying both the absolute competency of individual models and the relative similarity between models, ZEBRA ensures that each positive-negative pair captures meaningful contrasts in model capabilities, thus maximizing alignment informativeness.

Model Behavior Superiority (MB-SUP) quantifies the overall behavioral competency of model

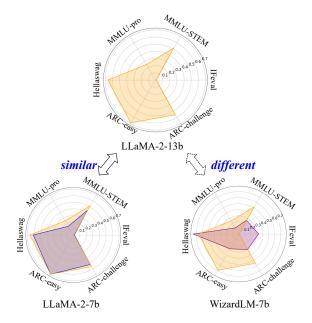


Figure 4: Hexagonal radar plots visualizing model ability profiles. Models judged as *similar* exhibit comparable benchmark performance trajectories, while *different* models show clearly divergent patterns.

 M_i as the aggregate of its normalized benchmark scores:

MB-SUP
$$(M_i) = \frac{1}{m} \sum_{k=1}^{m} s_i^{(b)}$$
. (1)

This scalar value serves as the ranking basis for constructing *Behavioral Superiority Anchors*.

Model Behavior Similarity (MB-SIM) between models M_i and M_j is defined as the similarity between their behavior vectors:

$$MB$$
- $SIM(M_i, M_j) = similarity(v_i, v_j).$

Higher values indicate stronger alignment in general-purpose capabilities, and MB-SIM serves as the criterion for selecting comparable model pairs in *Behavioral Similarity Anchoring*.

3.5 Strategy of Preference Binarization

ZEBRA introduces multiple strategies to systematically convert benchmark-derived model behaviors into binary preference pairs. Based on how MB-SUP and MB-SIM are utilized, our strategies can be clearly categorized as follows on Figure 5.

We detail each strategy below:

Strategy 1: Superiority-first Anchoring (SUP)

In this strategy, we explicitly select the top-two models based on their MB-SUP scores: the highest-scoring model (top-1) and the second-highest-scoring model (top-2). Responses from the top-1

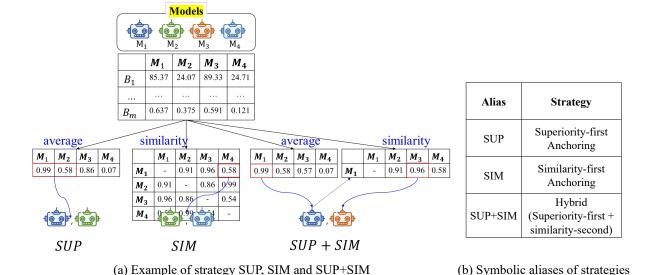


Figure 5: Anchoring strategies symbolic aliases and the example of strategy SUP, SIM and SUP+SIM with model and labeled preference.

model serve as positive anchors, while responses from the top-2 model become negative counterparts. This approach emphasizes explicit quality distinctions, clearly defining superior responses and ensuring meaningful, informative alignment contrasts.

Strategy 2: Similarity-first Anchoring (SIM)

Responses from models sharing similar behavioral patterns (high MB-SIM) are paired first. Within these pairs, the response from the model with higher MB-SUP is selected as the anchor (positive response). This strategy emphasizes behavioral similarity, enhancing nuanced alignment comparisons.

Strategy 3: Hybrid Anchoring (SUP+SIM)

Model pairs are selected by simultaneously considering both MB-SUP and MB-SIM criteria. This balanced approach ensures each response pair reflects meaningful contrasts in model quality, while maintaining behavioral similarity for refined granularity.

These strategies enable ZEBRA to flexibly and effectively tailor preference construction according to the desired granularity, alignment objectives, and computational resources. Figure 5 represent the example of each strategy.

3.6 Zero-Annotation Preference Construction

Given a set of instructions \mathcal{X} and a pool of response-generating models \mathcal{M} , the construction proceeds as follows:

1. Model Behavior Evaluation: Each model

- $M_i \in \mathcal{M}$ is benchmarked across m tasks to obtain model's ability vector v_i .
- 2. Pair Selection: Using a chosen strategy, a set of model pairs $\mathcal{P} = \{(M_i, M_j)\}$ is selected where MB-SIM $(M_i, M_j) \geq \tau$. Our goal is to exclude model pairs that are clearly dissimilar, rather than to finely rank models. As a commonly accepted cutoff for non-similarity, we set $\tau = 0.1$, since cosine similarity values below this threshold indicate a lack of meaningful similarity between models(Zhang et al., 2008).
- 3. **Response Mapping:** For each instruction $x \in \mathcal{X}$, the corresponding responses $\{r_i, r_j\}$ from (M_i, M_j) are retrieved.
- 4. **Preference Assignment:** A binary label is assigned via:

$$\operatorname{Pref}(r_i, r_j) = \begin{cases} 1 & \text{if } S_{\operatorname{pref}}(M_i) > S_{\operatorname{pref}}(M_j), \\ 0 & \text{otherwise.} \end{cases}$$

This pipeline constructs a binarized preference dataset at scale without any manual or per-instance scoring. Figure 2 shows the total process of the pipeline.

4 Experimental Setup

In this paper, we propose the ZEBRA framework. To validate the functionality of this framework and the characteristics of Alignment Tuning, we conducted experiments to address the following Research Questions (RQs):

- RQ1: Does MB-SIM represent the models' similarity?
- **RQ2**: Does strategies affect the binarization process?
- RQ3: Does ZEBRA reduce annotation cost and computational overhead compared to instance-level methods?

4.1 Alignment Dataset

For preference binarization, we utilized the Ultra-Feedback dataset (Cui et al., 2023), which includes diverse responses generated by both commercial and open-source language models. It provides multiple model responses, making it suitable for constructing a strong RLAIF baseline, for which we adopted the original score aggregation method. In contrast, ZEBRA uses only the response pool without referencing any scores, ensuring a clean comparison focused on the binarization strategy.

Model coverage. UltraFeedback contains outputs from 17 LLMs: GPT-4(OpenAI, 2023), GPT-3.5 Turbo(Ouyang et al., 2022b), and Bard(Waisberg et al., 2024). Additionally, several models from the Llama family were included, such as Llama-2 (7B, 13B, and 70B)-chat(Touvron et al., 2023), UltraLM-13B(Cui et al., 2024), WizardLM (7B, 13B, and 70B)(Xu et al., 2023), Vicuna-33B(Zheng et al., 2024), and Alpaca-7B(Taori et al., 2023). Beyond the Llama-based architectures, the dataset also features responses from other notable models, including Falcon-40B-instruct(Almazrouei et al., 2023), MPT-30B-chat(Team, 2023), StarChat-Beta(Tunstall et al., 2023), and Pythia-12B(Biderman et al., 2023).

Although the Ultrafeedback dataset contains results from UltraLM-65B(Cui et al., 2024), its performance could not be accurately assessed. To maintain the reliability of our evaluation, we excluded these results from the dataset composition.

4.2 Models

We fine-tuned and evaluated Llama-3.1-(3B, 8B)-Instruct(Dubey et al., 2024) and Qwen2.5-(3B, 8B)-Instruct(Yang et al., 2024) on the ZEBRA-binarized dataset. This setup enables cross-family comparison and quantifies the alignment gains from behavior-aware preference construction.

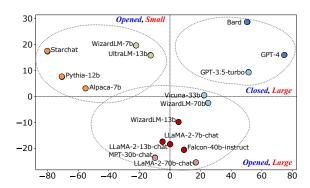


Figure 6: PCA-based visualization of relationships among models(Maćkiewicz and Ratajczak, 1993). The axes reflect similarity-based variance. Similar models are positioned closer to each other. Circle colors indicate clusters identified by hierarchical clustering.

| Strategy | ZEBRA-Human | Inter-Annotator |
|----------|-------------|-----------------|
| 20200083 | Agreement | Agreement |
| SUP | 78.6% | 0.65 |
| SIM | 76.7% | 0.586 |
| SUP+SIM | 78.0% | 0.68 |

Table 1: Human evaluation results comparing ZEBRA and human agreement rates.

4.3 Training Algorithm for Alignment Tuning

We tested two different learning method for alignment tuning using the ZEBRA Framework:

- Supervised Fine-Tuning (SFT)
- Direct Preference Optimization (DPO) (Rafailov et al., 2024)

The implementation details are provided in Appendix A. The full source code, benchmark results, and preference binarization scripts are released via our project repository under the MIT License. The snapshot corresponding to this paper is version-tagged.

5 Experimental Results

5.1 RQ1: Dataset Reconstruction using ZEBRA

The proposed ZEBRA framework enables the application of a unified Total Ranking map across the entire dataset, facilitating a structured and consistent preference mapping process. Utilizing this reconstructed preference dataset, we conducted model training while ensuring that each critic's binarized response was systematically incorporated. The effectiveness of this reconstructed dataset was assessed by evaluating the trained models on stan-

¹https://github.com/Jeesu-Jung/ZEBRA

| | | | Knowledge | | Re | asoning | Instruction-Following |
|-------------|----------|---------------------|-----------|----------|----------|---------------|-----------------------|
| Category | Strategy | Average | MMLU | MMLU-Pro | ARC-Easy | ARC-Challenge | IFeval |
| Baseline | RLAIF | 0.31 | 0.36 | 0.15 | 0.40 | 0.40 | 0.28 |
| | SUP | 0.31 (-0.00) | 0.33 | 0.15 | 0.40 | 0.37 | 0.30 |
| ZEBRA(Ours) | SIM | 0.29(-0.02) | 0.34 | 0.15 | 0.36 | 0.33 | 0.29 |
| | SUP+SIM | 0.29(-0.02) | 0.30 | 0.14 | 0.41 | 0.39 | 0.23 |

Table 2: Performance Comparison Between Instance-wise Binarized Data (Baseline) and Model Behavior-wise Binarized Data (Ours). The baseline corresponds to instance-wise scored RLAIF (Cui et al., 2023). The highlighted cells indicate performance that is equal to or higher than the baseline. **Bold** text shows the best performance on the benchmarks.

dardized benchmarks. The detailed results of these evaluations are presented in Appendix C.

For evaluation, we employed prediction-based assessment methodologies across all benchmark tasks. Specifically, for ARC, MMLU, and MMLU-Pro, we adapted the MMLU-Pro evaluation framework, modifying only the multiple-choice options to align with our dataset. For IFeval, we leveraged its native evaluation framework to ensure consistency in assessment.

To analyze the relationships between models, we computed SUP and SIM by normalizing evaluation results across six benchmark tasks. Figure 6 presents a Principal Component Analysis (PCA)(Maćkiewicz and Ratajczak, 1993) visualization of these relationships, illustrating distinct clustering patterns among models. Generally, smallerscale models tend to cluster in the first quadrant, models with closer Model Behavior relationships in the second quadrant, while Llama-based models and models exceeding 10B parameters are predominantly distributed in the third and fourth quadrants. This distribution underscores the critical role of model training data, training algorithms, and scale in determining Model Behavior relationships among models.

We further verify that anchor pairs selected by SIM indeed yield semantically closer responses. Using a TF–IDF(Ramos, 2003), we obtain an average response-pair similarity of 0.4623 for the *most similar model pair* (MB-SIM ↑) versus 0.4129 for the *least similar model pair* (MB-SIM ↓). This 12% relative gap confirms that high-MB-SIM anchoring indeed surfaces finer-grained yet coherent preference signals.

These findings highlight the effectiveness of the ZEBRA framework in reconstructing preference datasets, providing a more structured and informative approach to model alignment.

5.2 RQ2: Performance of ZEBRA Binarization

5.2.1 Effectiveness of Data Construction

We conducted a human evaluation to assess whether ZEBRA Binarization exhibits trends similar to human preferences. Specifically, we measured the agreement between ZEBRA decisions and human annotators across 150 pairs per strategy with three annotators. The evaluation design follows prior preference-based alignment studies, aiming to quantify the consistency of binarization with human judgment. Results are presented in Table 1.

ZEBRA-Human agreement rates are comparable to human-human agreement reported in previous studies on label variability (e.g., InstructGPT (Ouyang et al., 2022b): 72.6%, HHRLHF (Bai et al., 2022): 75%), indicating that ZEBRA-generated labels serve as credible proxies for preference learning. Furthermore, lower agreement on SIM pairs supports the hypothesis that these pairs are *harder to distinguish*, thereby providing richer alignment signals. This suggests that ZEBRA's binarization is not only aligned with human preferences but also introduces challenging cases that may enhance model robustness.

An illustrative example of response pairs under different strategies is provided in Appendix E.

5.2.2 Model Performance

To assess the effectiveness of ZEBRA binarization, we examined whether model performance can serve as an indicator of data quality. The results of this comparison are presented in table 2. Across all benchmark tasks, the Model Behavior-based scoring metric, SUP, demonstrates performance levels nearly equivalent to the instance-wise RLAIF method, with a minimal deviation of only 0.008.

Notably, for MMLU-Pro and IFeval, ZEBRA-based binarization even surpasses the RLAIF baseline by approximately 0.01, indicating that struc-

| Category | Method | Pairs / Units | Unit Cost (USD) | Total Cost (USD) |
|------------------------|--|-----------------------------|-------------------------|--------------------------|
| Instance-wise RLAIF | UltraFeedback(Cui et al., 2023) Safer-Instruct(Shi et al., 2024) OpenHermesPreferences(Huang et al., 2024) | 64,000 10,254 989,000 | 0.252 0.063 0.126 | 16,128 646 124,614 |
| Model behavior-wise | ZEBRA (ours) | 64,000 | 0 | 0 |

Table 3: Labeling-cost comparison between LLM-labeled RLAIF datasets, and our benchmark-table approach. cost is estimated assuming equivalent labor per rating.

tured preference mapping via Model Behavior can yield competitive or superior alignment outcomes. Table 2 shows the results of the average scores for each methodology. The detailed results, including average scores for each methodology, are summarized in Appendix F.

Statistical significance. We applied a Bland–Altman analysis(Bland and Altman, 1986) to assess agreement between ZEBRA (SUP) and baseline RLAIF performance across benchmarks. The results show a mean difference of -0.0188, with narrow limits of agreement from -0.2434 to +0.2058. Most differences fall well within this range, indicating high consistency and negligible deviation between the two methods.

5.3 RQ3: Cost & Efficiency Analysis

5.3.1 Labeling Cost

Table 3 contrasts the *labeling cost* of ZEBRA with canonical RLHF and RLAIF pipelines. The standing out points is below: **Absolute cost gap.** LLM-annotated RLAIF corpora lower this to \$0.6–\$4 K by outsourcing each pairwise judgment to GPT-4, yet they *still purchase every label*. ZEBRA, by reusing benchmark leaderboards, pays **no** marginal cost (\$0) for preference construction. **Relative efficiency.** Normalised per comparison, GPT-4 labels (\approx \$0.063) are \approx ×10 cheaper than human labels (\approx \$0.67), but ZEBRA is *orders of magnitude* cheaper than both because it dispenses with pairwise annotation altogether.

Cost matters only if quality survives. Despite a zero-dollar label budget, ZEBRA *matches or surpasses* RLAIF baseline (Table 2); the mean performance difference across six benchmarks is ≤ 0.02 .

Consequently, ZEBRA offers a **cost-minimal**, **scalable**, **and annotation-free** route to high-quality preference data. Because the price of human labor or GPT-4 tokens scales linearly with data volume, traditional pipelines become progressively more expensive as models, tasks, and safety domains proliferate. ZEBRA decouples alignment from annotation cost: adding a new model needs

| Benchmark | Inference (USD) | Evaluation (USD) |
|---------------|-----------------|-------------------------|
| IFeval | 0.0606 | 16.23 (GPT-4 eval) |
| MMLU-STEM | 0.3528 | 0 |
| MMLU-pro | 1.3440 | 0 |
| Hellaswag | 1.1250 | 0 |
| arc-challenge | 0.1312 | 0 |
| arc-easy | 0.2661 | 0 |
| Total | 3.28 | 16.23 |

Table 4: Inference and evaluation cost estimation for 6 benchmarks using Llama-7B and GPT-4.

no further labels, and incorporating an extra benchmark simply augments the behavior matrix.

5.3.2 Benchmark Inference and Evaluation Cost

We acknowledge that ZEBRA shifts the burden from human annotation to inference and benchmarking infrastructure. However, in practice:

- Many models already publish benchmark scores, which ZEBRA can directly utilize.
- When benchmarking is necessary, the cost is orders of magnitude lower than human annotation.

Based on Llama-7B (Llama 2 Chat-7B), we estimated the benchmark inference and evaluation cost using the pricing provided by Together/Airtrain (input + output = \$0.56 per 1M tokens) airtrain.ai Inference pricing.

Cost per token =
$$\frac{0.56}{1,000,000} = 5.6 \times 10^{-7} \text{ USD}$$

So, the inference cost per example, assuming an average of 200 tokens, is:

Cost per example =
$$200 \times 5.6 \times 10^{-7}$$

= 1.12×10^{-4} USD
 ≈ 0.000112 USD

Table 4 shows the price of each benchmark.

For example, benchmarking 6 models over 6 benchmarks (totaling 29,281 samples) costs only \$3.28 on Llama-7B (Airtrain pricing). Evaluation

using GPT-4 adds \$16.23 in total. This cost is negligible compared to human preference collection (\$5–20 per 100 samples) or LLM instance-wise annotation (~\$16,128 for Ultrafeeedback).

6 Practical Insights and Recommendations

A key motivation for ZEBRA is its practical usability when researchers or practitioners wish to expand their alignment datasets without additional annotation cost. Conventional RLHF or RLAIF pipelines require costly human or LLM-based supervision, but ZEBRA enables an alternative: leveraging existing benchmark performance tables to generate preference signals at zero cost.

Practical Use Case. The procedure for applying ZEBRA in practice can be summarized as follows:

- 1. Inspect the benchmark performance table for the models of interest.
- Select two models that exhibit the highest behavioral similarity (MB-SIM) or fall under a chosen binarization strategy (SUP, SIM, or SUP+SIM).
- 3. Generate responses from each model on the desired instruction set.
- 4. Use the resulting preference pairs directly for fine-tuning (e.g., SFT or DPO).

7 Conclusion

We introduced **ZEBRA**, a zero-annotation framework that replaces costly instance-level preference labeling with model-behavior knowledge from benchmark tables. By leveraging superiority, similarity, or hybrid strategies, ZEBRA constructs large-scale preference datasets entirely without annotation cost.

Experiments across six standardized benchmarks show that ZEBRA achieves performance comparable to, and in some cases exceeding, RLAIF baselines, confirming benchmark-derived behavior as an effective proxy for preference supervision. Future work will extend ZEBRA to additional behavioral axes (e.g., safety, toxicity, multilingual ability) and further examine its correlation with explicit human judgments.

Future work will expand to additional behavioral axes (e.g., safety, toxicity, multilinguality), explore adaptive pairing schedules, and assess correlation with human judgments.

8 Related Work and Background

8.1 Alignment Tuning

Alignment tuning has become a central focus in enhancing LLMs to meet user expectations and ethical standards(Kumar et al., 2024). Various preference-based learning techniques, particularly reinforcement learning methods(Schulman et al., 2017; Rafailov et al., 2024; Hong et al., 2024), have been developed to facilitate this tuning process.

These methods rely on preference datasets, which typically contain pairs of responses generated by models based on given instructions, with each pair ranked according to human or modelbased evaluations. Prominent alignment datasets like Ultrafeedback(Cui et al., 2023), which gathers human feedback, and HH-RLHF(Bai et al., 2022), which uses human-annotated preferences, provide foundational resources for alignment. To reduce costly data curation, recent automated approaches filter or regenerate preference data based on specific criteria, improving consistency and scalability. However, such methods often lack fine-grained control over data quality, as they overlook the origin model's capabilities in shaping alignment effectiveness (Shi et al., 2024).

8.2 Limitations in Existing Preference Data Approaches

Effective preference data construction requires a clear, rigorous set of criteria to ensure alignment quality across generated response pairs. Common criteria, including *Reasoning*, *Truthfulness*, and *Instruction-Following*, guide the selection of data that aligns with key ethical and functional standards(Cui et al., 2023; Bai et al., 2022). High-performing models, capable of producing responses that meet these standards, are often evaluated using benchmarks like *ARC*(*Clark et al.*, 2018), *MMLU*(*Hendrycks et al.*, 2020), and the *Instruction-Following eval*(*Zhou et al.*, 2023) suite, which assess a model's factual accuracy, reasoning ability, and compliance with instructions.

The ZEBRA Framework addresses this limitation by introducing a model behavior-level approach that emphasizes model compatibility in preference data selection. By curating preference pairs based on model behavior knowledge with similar core competencies in knowledge, reasoning, and instruction-following—the ZEBRA Framework enhances alignment coherence and ensures stable, high-quality data.

Limitation

In this paper, we demonstrated that ZEBRA can achieve performance comparable to existing RLAIF without requiring annotations for SFT and DPO. However, due to resource and time constraints, we were unable to validate our approach across a broader range of alignment tuning techniques. Further evaluation is needed for methods.

Additionally, our evaluation primarily focused on fundamental abilities such as Reasoning, Knowledge, and Instruction-Following. However, we did not assess ZEBRA's performance on other important values, including factuality and ethical standard. Future work should incorporate evaluations reflecting these aspects.

9 Acknowledgments

This work was partly supported by research fund of Chungnam National University, Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) (No.RS-2022-00155857, Artificial Intelligence Convergence Innovation Human Resources Development (Chungnam National University)) and the National Research Foundation of Korea(NRF) grant funded by the Korea government(MSIT)(No. RS-2025-0055621731482092640101)

References

- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. Falcon-40B: an open large language model with state-of-the-art performance.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv* preprint arXiv:2204.05862.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.

- J. Martin Bland and Douglas G. Altman. 1986. Statistical methods for assessing agreement between two methods of clinical measurement. *The Lancet*, 327(8476):307–310.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv* preprint arXiv:1803.05457.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, et al. 2024. Ultrafeedback: Boosting language models with scaled ai feedback. In Forty-first International Conference on Machine Learning.
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv* preprint arXiv:2407.21783.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Jiwoo Hong, Noah Lee, and James Thorne. 2024. ORPO: Monolithic preference optimization without reference model. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11170–11189, Miami, Florida, USA. Association for Computational Linguistics.
- Shengyi Costa Huang, Agustín Piqueres, Kashif Rasul, Philipp Schmid, Daniel Vila, and Lewis Tunstall. 2024. Open hermes preferences. https://huggingface.co/datasets/argilla/OpenHermesPreferences.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12).
- Anurakt Kumar, Divyanshu Kumar, Jatan Loya, Nitin Aravind Birur, Tanay Baswa, Sahil Agarwal, and Prashanth Harshangi. 2024. Sage-rt: Synthetic alignment data generation for safety evaluation and red teaming. *Preprint*, arXiv:2408.11851.
- Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, et al. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback. In Forty-first International Conference on Machine Learning.

- Andrzej Maćkiewicz and Waldemar Ratajczak. 1993. Principal components analysis (pca). *Computers & Geosciences*, 19(3):303–342.
- OpenAI. 2023. Gpt-4 technical report. arXiv preprint arXiv:2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022a. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022b. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Juan Ramos. 2003. Using tf-idf to determine word relevance in document queries.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*.
- Archit Sharma, Sedrick Keh, Eric Mitchell, Chelsea Finn, Kushal Arora, and Thomas Kollar. 2024. A critical evaluation of ai feedback for aligning large language models. *arXiv preprint arXiv:2402.12366*.
- Taiwei Shi, Kai Chen, and Jieyu Zhao. 2024. Safer-instruct: Aligning language models with automated preference data. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7636–7651, Mexico City, Mexico. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- MosaicML NLP Team. 2023. Introducing mpt-30b: Raising the bar for open-source foundation models. Accessed: 2023-06-22.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

- Lewis Tunstall, Nathan Lambert, Nazneen Rajani, Edward Beeching, Teven Le Scao, Leandro von Werra, Sheon Han, Philipp Schmid, and Alexander Rush. 2023. Creating a coding assistant with starcoder. *Hugging Face Blog*. Https://huggingface.co/blog/starchat.
- Ethan Waisberg, Joshua Ong, Mouayad Masalkhi, Nasif Zaman, Prithul Sarker, Andrew G Lee, and Alireza Tavakkoli. 2024. Google's ai chatbot "bard": a sideby-side comparison with chatgpt and its utilization in ophthalmology. *Eye*, 38(4):642–645.
- Zhilin Wang, Alexander Bukharin, Olivier Delalleau, Daniel Egert, Gerald Shen, Jiaqi Zeng, Oleksii Kuchaiev, and Yi Dong. 2025. Helpsteer2-preference: Complementing ratings with preferences. *Preprint*, arXiv:2410.01257.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *arXiv* preprint arXiv:2304.12244.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.
- Michael JQ Zhang, Zhilin Wang, Jena D. Hwang, Yi Dong, Olivier Delalleau, Yejin Choi, Eunsol Choi, Xiang Ren, and Valentina Pyatkin. 2024. Diverging preferences: When do annotators disagree and do models know? *Preprint*, arXiv:2410.14632.
- Xiaodan Zhang, Xiaohua Hu, and Xiaohua Zhou. 2008. A comparative evaluation of different link types on enhancing document clustering. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '08, page 555–562, New York, NY, USA. Association for Computing Machinery.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023. Judging Ilm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36.
- Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. 2023. Instruction-following evaluation for large language models. *arXiv preprint arXiv:2311.07911*.

Prompt for Benchmark Evaluation Q: {Question} Options: (A) {Option 1} (B) {Option 2} (C) {Option 3} (D) {Option 4} Output:

Figure 7: Evaluation template for multiple-choice benchmark evaluation. The descriptive benchmark (e.g., IFeval) was evaluated using the same template, excluding the option part.

A Implementation Details for Model Training

The model was trained for a single epoch using bfloat16 (bf16) quantization, which optimizes memory efficiency while preserving numerical precision. The training configuration incorporated the following hyperparameters: a per-device batch size of 6, gradient accumulation steps set to 4, and a learning rate of 5×10^{-5} , with 500 warm-up steps to facilitate stable convergence.

Training was conducted on an L40 4-GPU setup, leveraging an optimized deep learning framework to enhance computational efficiency. The training pipeline focused on performance optimization through checkpointing and logging, without intermediate model evaluation during training.

B Prompt template for Evaluation Benchmarks

To determine behavior knowledge between models, this paper evaluates benchmark performance and compares the similarity of these numerical results. The benchmark performance of all 17 models considered in this study is presented in Table 7. For the start and end tokens, the template recommended in the paper was used.

C Evaluation Model Benchmark Performance

To determine behavior knowledge relationships between models, this paper evaluates benchmark performance and compares the similarity of these numerical results. The benchmark performance of all 17 models considered in this study is presented in Table 5. The actual calculation example for

MB - SUP and MB - SIM can be found in Figure 10.

Due to the unavailability of Bard, its performance metrics have been substituted with those of Gemini-1.5-Flash.

D Benchmark Performance Representation

What benchmark represent the model capability: knowledge, instruction-following, and reasoning? Models exhibit distinct similarity patterns based on their capabilities, with these patterns varying across different evaluation metrics. Notably, model behavior knowledge is influenced by factors such as model scale and training methodology, leading to variations in clustering behavior across different capability dimensions.

Analysis Framework To systematically investigate these relationships, we analyzed model similarities across three key capability dimensions:

- **Knowledge-Based Tasks**: Unlike reasoning capability, clustering in instruction-following tasks is more strongly aligned with *model families* rather than size. This indicates that training methodology and architectural choices exert a greater influence on instruction adherence.
- **Reasoning Capability**: Models tend to cluster primarily based on parameter count, suggesting that *model size* plays a dominant role in shaping reasoning performance.
- **Instruction-Following (IF) Tasks**: A hierarchical influence pattern emerges, where:
 - At the initial hierarchy, the *model family* is the primary determinant.
 - At the higher hierarchy, the model size becomes a stronger predictor of performance.

Figure 11 visualizes these relationships through dendrograms, illustrating the hierarchical clustering patterns that emerge across different capability dimensions.

Figure 8 shows the frequency of each model being selected as the positive or negative under each strategy. Overall, larger and more familiar models tend to be chosen more frequently.

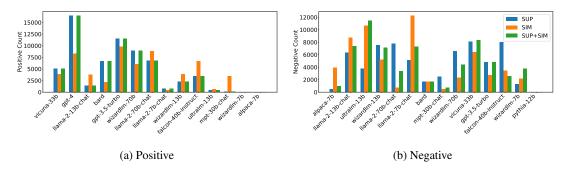


Figure 8: Selected Model Frequency for the positive and negative pairs.

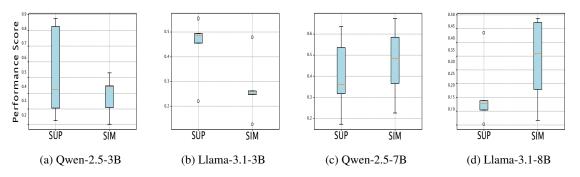


Figure 9: Visualization of performance across different models. This represents the performance when using the DPO training method. For smaller models, using SUP results in better performance, whereas for larger models, SIM yields better performance

E Example of ZEBRA Binarization

To provide further insight into how ZEBRA Binarization operates under different strategies, we present an illustrative example in Table 6. This example highlights the SUP and SIM strategies, showing both the chosen and rejected responses. As demonstrated, SUP captures straightforward preference alignment, while SIM introduces more subtle semantic overlaps, which tend to be harder to distinguish.

F Total Performance

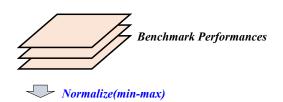
Table 7 provides an aggregated view of the overall performance of each model across all benchmark datasets. The total performance scores were computed by averaging the normalized scores across the selected evaluation metrics, offering a holistic comparison of model capabilities.

G Model Size and ZEBRA Strategies

The impact of different Strategies varies significantly with model size, influencing how models learn from preference data. To investigate this relationship, we conducted experiments using models from the same family but with different parameter counts, specifically comparing small models (3B parameters) and larger models (7-8B parameters). Our analysis reveals a clear pattern in how model size determines the optimal strategy.

Smaller models (3B) exhibit superior performance when trained using the SUP algorithm, which prioritizes learning from response quality within the binarized dataset. In contrast, larger models (7-8B) achieve better results with the SIM algorithm, suggesting that response similarity becomes increasingly important as model size grows.

This trend indicates that model size fundamentally influences how different models leverage preference data. While smaller models benefit more from explicit learning based on absolute quality differences, larger models demonstrate greater sensitivity to the nuanced relationships between similar responses in the training data. This relationship is visualized in Figure 9, illustrating the distinct learning behaviors observed across different model scales.



MMLU -STEM

MMLU

| Mod | el | IFev | al | MMLU -STEM | 1 | ILU ro | Hellaswa | | ARC Easy | ARC -Challer | | | | | | |
|-------------------------|---------|---------|-------|---------------|---------------|---------------|-----------------------|--------|--------------|-----------------|--------|--------|-------------------|---------------|-------------|--------|
| GPT | -4 | 0.93 | 39 | 1.000 | | 000 | 1.000 | | .000 | 1.000 | | | | | | |
| GPT-3.5- | turbo | 0.74 | | 0.724 | | 198 | 0.856 | | .952 | 0.847 | | | | | | |
| Bar | d | 1.00 | 00 | 0.754 | 0.9 | 912 | 0.852 | 0 | .847 | 0.779 |) | | | | | |
| Llama-2-7 | /b-chat | 0.0 | 17 | 0.439 | 0.1 | 64 | 0.612 | 0 | .693 | 0.521 | 1 | | | | | |
| Llama-2-1 | 3b-chat | 0.0 | 11 | 0.519 | 0.2 | 261 | 0.664 | 0 | .691 | 0.561 | 1 | | | | _ | |
| Llama-2-7 | 0b-chat | 0.00 | 00 | 0.705 | 0.4 | 196 | 0.760 | 0 | .819 | 0.671 | | | • <i>MB</i> • | – SU I | , | |
| UltraLN | 1-13b | 0.47 | | 0.380 | _ | 46 | 0.424 | | .510 | 0.446 | | | | | | |
| WizardI | | 0.33 | | 0.261 | | 28 | 0.755 | | .283 | 0.296 | | | | | | |
| WizardL | | 0.13 | | 0.426 | |)98 | 0.801 | _ | .703 | 0.530 | | | | | | |
| WizardL | | 0.39 | | 0.433 | | 296 | 0.833 | _ | .800 | 0.720 | _ | | | | | |
| Vicuna | | 0.44 | | 0.623 | | 217 | 0.717 | | .811 | 0.635 | | | | | | |
| Alpaca | | 0.10 | | 0.184 | |)55 | 0.000 | | .331 | 0.283 | | | | | | |
| Falcon-40b MPT-30l | | 0.00 | | 0.682 | | 165 | $\frac{0.787}{0.013}$ | | .749 .881 | 0.545 | | | | | | |
| Starch | | 0.10 | | 0.394 | | 000 | 0.013 | | .000 | 0.000 | | | | | MD | CIM |
| Pythia- | | 0.0 | | 0.000 | | 006 | 0.028 | _ | .095 | 0.260 | | | | | <i>MB</i> – | SIM |
| | | | | | | , , , | 0.020 | | .0,0 | 0.20 | | | | | 1 | |
| | Calcula | ite MB | – SIM | • | | | | | | | | | | | | |
| | | GPT-3.5 | | Llama-2 | Llama-2 | Llama-2 | UltraL | Wizard | Wizord | Wizard | Vicuna | Alpaca | Falcon | MPT | | Pythia |
| | GPT-4 | -turbo | Bard | -7b -chat | -13b -chat | -70b -chat | -13b | LM-7b | | LM-70b | -33b | -7b | -40b -instruct | -30b -chat | Starchat | -12b |
| GPT-4 | 1.000 | 0.983 | 0.994 | 0.869 | 0.890 | 0.911 | 0.956 | 0.869 | 0.872 | 0.942 | 0.949 | 0.803 | 0.832 | 0.767 | 0.546 | 0.589 |
| GPT-3.5 -turbo | 0.983 | 1.000 | 0.973 | 0.915 | 0.922 | 0.919 | 0.990 | 0.899 | 0.928 | 0.979 | 0.987 | 0.846 | 0.886 | 0.813 | 0.517 | 0.645 |
| Bard | 0.994 | 0.973 | 1.000 | 0.822 | 0.843 | 0.864 | 0.951 | 0.864 | 0.836 | 0.921 | 0.926 | 0.769 | 0.779 | 0.726 | 0.518 | 0.548 |
| Llama-2 -7b-chat | 0.869 | 0.915 | 0.822 | 1.000 | 0.996 | 0.978 | 0.888 | 0.851 | 0.985 | 0.962 | 0.957 | 0.830 | 0.984 | 0.845 | 0.422 | 0.679 |
| Llama-2 -13b-chat | 0.890 | 0.922 | 0.843 | 0.996 | 1.000 | 0.992 | 0.885 | 0.855 | 0.977 | 0.960 | 0.956 | 0.818 | 0.979 | 0.831 | 0.455 | 0.663 |
| Llama-2 -70b-chat | 0.911 | 0.919 | 0.864 | 0.978 | 0.992 | 1.000 | 0.871 | 0.829 | 0.947 | 0.942 | 0.941 | 0.813 | 0.954 | 0.827 | 0.483 | 0.634 |
| UltraLM -13b | 0.956 | 0.990 | 0.951 | 0.888 | 0.885 | 0.871 | 1.000 | 0.890 | 0.913 | 0.963 | 0.981 | 0.852 | 0.868 | 0.807 | 0.533 | 0.643 |
| WizardLM- 7b | 0.869 | 0.899 | 0.864 | 0.851 | 0.855 | 0.829 | 0.890 | 1.000 | 0.920 | 0.925 | 0.903 | 0.566 | 0.856 | 0.535 | 0.439 | 0.479 |
| WizardLM -13b | 0.872 | 0.928 | 0.836 | 0.985 | 0.977 | 0.947 | 0.913 | 0.920 | 1.000 | 0.977 | 0.968 | 0.776 | 0.978 | 0.778 | 0.419 | 0.643 |
| WizardLM -70b | 0.942 | 0.979 | 0.921 | 0.962 | 0.960 | 0.942 | 0.963 | 0.925 | 0.977 | 1.000 | 0.985 | 0.822 | 0.930 | 0.811 | 0.406 | 0.694 |
| Vicuna-33b | 0.949 | 0.987 | 0.926 | 0.957 | 0.956 | 0.941 | 0.981 | 0.903 | 0.968 | 0.985 | 1.000 | 0.857 | 0.945 | 0.833 | 0.536 | 0.647 |
| Alpaca-7b | 0.803 | 0.846 | 0.769 | 0.830 | 0.818 | 0.813 | 0.852 | 0.566 | 0.776 | 0.822 | 0.857 | 1.000 | 0.783 | 0.990 | 0.418 | 0.785 |
| Falcon-40b -instruct | 0.832 | 0.886 | 0.779 | 0.984 | 0.979 | 0.954 | 0.868 | 0.856 | 0.978 | 0.930 | 0.945 | 0.783 | 1.000 | 0.788 | 0.529 | 0.605 |
| MPT-30b -chat | 0.767 | 0.813 | 0.726 | 0.845 | 0.831 | 0.827 | 0.807 | 0.535 | 0.778 | 0.811 | 0.833 | 0.990 | 0.788 | 1.000 | 0.335 | 0.798 |
| Starchat | 0.546 | 0.517 | 0.518 | 0.422 | 0.455 | 0.483 | 0.533 | 0.439 | 0.419 | 0.406 | 0.536 | 0.418 | 0.529 | 0.335 | 1.000 | 0.021 |
| Pythia-12b | 0.589 | 0.645 | 0.548 | 0.679 | 0.663 | 0.634 | 0.643 | 0.479 | 0.643 | 0.694 | 0.647 | 0.785 | 0.605 | 0.798 | 0.021 | 1.000 |

ARC

ARC

Figure 10: The process of calculating MB - SUP and MB_SIM . Benchmark performance is normalized using min-max normalization, and the overall SIM is performed using cosine similarity.

| Model | IFeval | MMLU-STEM | MMLU-pro | Hellaswag | ARC-easy | ARC-Challenge | Average Score |
|---------------------|--------|-----------|----------|-----------|----------|---------------|---------------|
| GPT-4 | 85.37 | 86.40 | 0.64 | 95.30 | 96.63 | 96.40 | 0.99 |
| GPT-3.5-turbo | 72.54 | 70.00 | 0.38 | 85.00 | 92.80 | 83.02 | 0.77 |
| bard | 89.33 | 71.80 | 0.59 | 84.70 | 84.43 | 77.13 | 0.86 |
| Llama-2-7b-chat | 25.19 | 53.10 | 0.20 | 67.50 | 72.14 | 54.61 | 0.41 |
| Llama-2-13b-chat | 24.82 | 57.80 | 0.25 | 71.20 | 72.05 | 58.02 | 0.45 |
| Llama-2-70b-chat | 24.07 | 68.90 | 0.38 | 78.10 | 82.20 | 67.66 | 0.58 |
| UltraLM-13b | 54.92 | 49.58 | 0.19 | 54.00 | 57.58 | 48.04 | 0.40 |
| WizardLM-7b | 45.83 | 42.50 | 0.18 | 77.70 | 39.48 | 34.90 | 0.34 |
| WizardLM-13b | 33.92 | 52.30 | 0.17 | 81.00 | 72.94 | 55.38 | 0.45 |
| WizardLM-70b | 49.51 | 52.70 | 0.27 | 83.30 | 80.68 | 71.93 | 0.58 |
| Vicuna-33b | 52.76 | 64.00 | 0.23 | 75.00 | 81.57 | 64.51 | 0.57 |
| Alpaca-7b | 30.58 | 37.92 | 0.15 | 23.60 | 43.31 | 33.79 | 0.16 |
| Falcon-40b-instruct | 24.54 | 67.50 | 0.14 | 80.00 | 76.60 | 56.70 | 0.47 |
| MPT-30b-chat | 30.70 | 50.40 | 0.20 | 24.53 | 87.12 | 70.73 | 0.38 |
| Starchat | 28.30 | 40.12 | 0.12 | 25.40 | 16.96 | 9.07 | 0.05 |
| Pythia-12b | 24.71 | 27.00 | 0.12 | 25.60 | 24.49 | 31.80 | 0.07 |

Table 5: Benchmark Scores for Trained Models. Multiple-choice benchmarks (MMLU-STEM, HellaSwag, ARC-Easy, and ARC-Challenge) are evaluated based on accuracy. IFeval and MMLU-Pro are assessed using its own metric. The average score is computed after min-max normalization.

| Input | Strategy | Chosen | Rejected |
|---|----------|---|--|
| Teacher: In this task you are given a statement and an explanation giving you further knowledge about an entity in the statement. You must judge whether the statement is true or false based on the explanation. Label an instance as "True" if the explanation confirms the statement or does not disprove it. Label an instance as "False" if the explanation disproves the statement. | SUP | Solution: True Reason: The explanation confirms that White people are Caucasians, which supports the statement that Caucasians are often described as White. | This statement is false and hence cannot be labeled "True." The reason for this is that not all Whites are Caucasians. There are many non-White people living in the Caucasus, a region which extends from the Black Sea to the Caspian Sea. Therefore, the explanation provided does not support the statement. |
| | SIM | Solution: True Reason: The explanation confirms that White people are Caucasians, which supports the statement that Caucasians are often described as White. | Reason: The explanation disproves the statement. "White people" is a broader term that includes people of different ethnic backgrounds who have light skin, while "Caucasians" specifically refers to people of European descent. Therefore, not all "White people" are Caucasians. |

Table 6: Illustrative examples of response pairs under SUP and SIM strategies.

| Benchmark | Model | | Baseline | SUP | SIM | SUP+SIM |
|---------------|----------------|-----|----------|--------|--------|---------|
| | I lama 2 1 2D | SFT | 0.2331 | 0.2679 | 0.2434 | 0.2200 |
| | Llama-3.1-3B | DPO | 0.3130 | 0.4868 | 0.2528 | 0.3940 |
| NO WILL COURT | I 1 2 1 0D | SFT | 0.2800 | 0.2460 | 0.2660 | 0.2240 |
| | Llama-3.1-8B | DPO | 0.4902 | 0.4349 | 0.4867 | 0.2880 |
| MMLU-STEM | O2 5 2D | SFT | 0.2760 | 0.2200 | 0.2460 | 0.2500 |
| | Qwen2.5-3B | DPO | 0.4706 | 0.4212 | 0.4400 | 0.4580 |
| | Owen 2.5.7D | SFT | 0.2800 | 0.2280 | 0.2620 | 0.2780 |
| | Qwen2.5-7B | DPO | 0.5120 | 0.3620 | 0.4860 | 0.2800 |
| | I lama 2 1 2D | SFT | 0.1189 | 0.1230 | 0.1045 | 0.1332 |
| | Llama-3.1-3B | DPO | 0.1455 | 0.2193 | 0.1270 | 0.1516 |
| | I lomo 2 1 9D | SFT | 0.1004 | 0.1148 | 0.1045 | 0.1311 |
| MMI II mmo | Llama-3.1-8B | DPO | 0.1168 | 0.1025 | 0.1168 | 0.1107 |
| MMLU-pro | Owen 2.5.2D | SFT | 0.1025 | 0.0840 | 0.1230 | 0.0820 |
| | Qwen2.5-3B | DPO | 0.2275 | 0.2254 | 0.2029 | 0.2111 |
| | Owen 2.5.7D | SFT | 0.0861 | 0.1762 | 0.2275 | 0.1352 |
| | Qwen2.5-7B | DPO | 0.2377 | 0.1721 | 0.2254 | 0.1700 |
| | I lama 2 1 2D | SFT | 0.2494 | 0.2410 | 0.2490 | 0.2206 |
| | Llama-3.1-3B | DPO | 0.3765 | 0.4940 | 0.4796 | 0.3033 |
| | I lama 2 1 0D | SFT | 0.2494 | 0.4210 | 0.2470 | 0.2218 |
| IFeval | Llama-3.1-8B | DPO | 0.2421 | 0.1882 | 0.2292 | 0.2494 |
| irevai | Qwen2.5-3B | SFT | 0.2190 | 0.2134 | 0.2122 | 0.2050 |
| | Qwell2.3-3B | DPO | 0.3633 | 0.3058 | 0.3094 | 0.1715 |
| | Qwen2.5-7B | SFT | 0.2290 | 0.2134 | 0.2122 | 0.2083 |
| | Qwell2.3-7B | DPO | 0.3321 | 0.3177 | 0.3657 | 0.2407 |
| | Llama-3.1-3B | SFT | 0.2660 | 0.2460 | 0.2000 | 0.2340 |
| | Liailia-3.1-3D | DPO | 0.6111 | 0.5547 | 0.2618 | 0.6679 |
| | Llama-3.1-8B | SFT | 0.2180 | 0.2330 | 0.2400 | 0.2360 |
| ARC-easy | Liailia-3.1-oD | DPO | 0.2176 | 0.1540 | 0.4720 | 0.0680 |
| ARC-easy | Qwen2.5-3B | SFT | 0.2410 | 0.2560 | 0.2480 | 0.2560 |
| | Qwell2.3-3B | DPO | 0.8497 | 0.8754 | 0.5295 | 0.8157 |
| | Qwen2.5-7B | SFT | 0.2550 | 0.2260 | 0.2280 | 0.2900 |
| | Qwell2.3-7B | DPO | 0.5700 | 0.6359 | 0.6738 | 0.7134 |
| | Llama-3.1-3B | SFT | 0.2556 | 0.2492 | 0.2266 | 0.2019 |
| ADC shallows | Liailia-3.1-3D | DPO | 0.5122 | 0.4551 | 0.2466 | 0.5712 |
| | Llomo 2 1 OD | SFT | 0.2320 | 0.2297 | 0.2761 | 0.2268 |
| | Llama-3.1-8B | DPO | 0.5463 | 0.1763 | 0.3596 | 0.2946 |
| ARC-challenge | Owan2 5 2D | SFT | 0.2645 | 0.2483 | 0.2227 | 0.2343 |
| | Qwen2.5-3B | DPO | 0.7398 | 0.8241 | 0.4470 | 0.7078 |
| | Owen 2.5.7D | SFT | 0.2343 | 0.2552 | 0.2483 | 0.2390 |
| | Qwen2.5-7B | DPO | 0.4432 | 0.5358 | 0.5847 | 0.3870 |

Table 7: Model Performance Comparisons on Knowledge, Instruction-Following, and Reasoning-Related Tasks. The baseline is Instance-wise RLAIF (Cui et al., 2023). "Llama 3.1" refers to the Llama-3.1-Instruct series, and "Qwen-2.5" refers to the Qwen2.5-Instruct series. The training methods include Supervised Fine-Tuning (SFT) and Direct Preference Optimization (DPO). The **bold** text indicates the best performance for each model and training method.

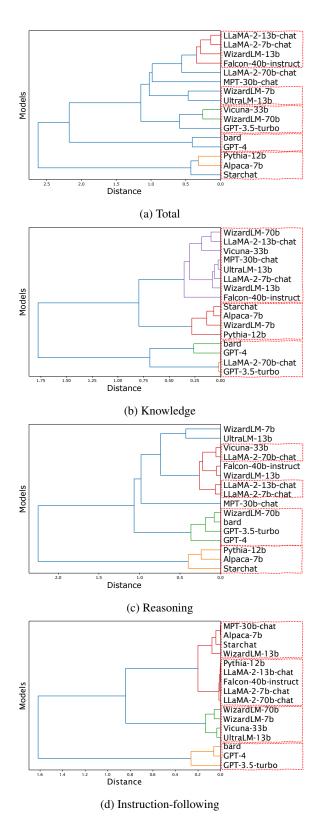


Figure 11: Dendrogram for Evaluation Categories: Knowledge, Reasoning, and Instruction-Following. Clusters exceeding a specific distance threshold (0.4) are highlighted in different colors. Models deemed similar share the same color line. The red dotted line shows the major groups of models in the dendrogram.