

Multi-modal Preference Alignment Remedies Degradation of Visual Instruction Tuning on Language Models

Shengzhi Li
TIFIN Inc
alex.li@tifin.com

Rongyu Lin
KAUST
rongyu.lin@kaust.edu.sa

Shichao Pei*
University of Massachusetts Boston
shichao.pei@umb.edu

Abstract

Multi-modal large language models (MLLMs) are expected to support multi-turn queries of interchanging image and text modalities in production. However, the current MLLMs trained with visual-question-answering (VQA) datasets could suffer from degradation, as VQA datasets lack the diversity and complexity of the original text instruction datasets with which the underlying language model was trained. To address this degradation, we first collect a lightweight, 5k-sample VQA preference dataset where answers were annotated by Gemini for five quality metrics in a granular fashion and investigate standard Supervised Fine-tuning, rejection sampling, Direct Preference Optimization (DPO) and SteerLM algorithms. Our findings indicate that with DPO, we can surpass the instruction-following capabilities of the language model, achieving a 6.73 score on MT-Bench, compared to Vicuna's 6.57 and LLaVA's 5.99. This enhancement in textual instruction-following capability correlates with boosted visual instruction performance (+4.9% on MM-Vet, +6% on LLaVA-Bench), with minimal alignment tax on visual knowledge benchmarks compared to the previous RLHF approach. In conclusion, we propose a distillation-based multi-modal alignment model with fine-grained annotations on a small dataset that restores and boosts MLLM's language capability after visual instruction tuning.

1 Introduction

Recent advancements in artificial intelligence have led to the rise of multi-modal large language models (MLLMs), which combine textual and visual interpretation capabilities in a single model (Shen et al., 2023). However, effectively blending multi-modality in one system has proven non-trivial. Integrating diverse data forms often creates internal

representation conflicts, giving rise to the issue known as "catastrophic forgetting" (Kirkpatrick et al., 2017). The diversity constraint in visual question answering (VQA) datasets could be attributed as a source of the issue. VQA tasks typically focus on descriptive queries about image contents, whereas textual datasets encompass a broader range of complex cognitive tasks, including reasoning, writing, summarization, and coding. This discrepancy in dataset complexity is a key factor contributing to the observed performance degradation in MLLMs. Our evaluation of models such as BLIP-2, InstructBLIP, and LLaVA against language instruction-following benchmarks like MT-Bench (Zheng et al., 2023) and AlpacaEval (Li et al., 2023b) revealed diminished language capabilities in comparison to their linguistic backbones. For instance, LLaVA, built on the Vicuna-13b LLM, demonstrated a decline in MT-Bench performance from 6.57 to 5.92, even underperforming the Vicuna-7B model.

Driven by the limitations observed in distillation-based instruction tuning, particularly its constrained generalizability and the narrow performance improvements on tasks outside the training distribution, this study investigates the efficacy of distillation-based preference alignment in addressing modality conflict in MLLMs. The decision to explore this avenue is predicated on the hypothesis that integrating AI-generated preference data can provide a more granular and nuanced alignment with human expectations, potentially mitigating the adverse effects of modality conflict.

This study rigorously evaluates three baseline methodologies—Direct Preference Optimization (DPO), SteerLM, and Rejection Sampling—as potential solutions to utilize the distilled preference data to enhance the instruction-following capabilities and address the modality conflict inherent in MLLMs. Each of these methods offers a unique approach to model alignment, from the direct opti-

*Corresponding author.

mization of preferences in DPO to the conditional supervision in SteerLM and the selective acceptance in Rejection Sampling. Our empirical analysis reveals that DPO, in particular, demonstrates a pronounced efficacy in reconciling the performance disparities observed between textual and visual modalities. By leveraging a refined preference dataset, fine-tuned with the DPO objective and supplemented with comprehensive annotations from advanced AI models, DPO not only addresses the modality conflict but also significantly enhances the model’s performance across a spectrum of benchmarks. The results indicate that, through the application of DPO, MLLMs can achieve a more robust alignment with human-like preferences, thereby mitigating the adverse effects of catastrophic forgetting and modality conflict and elevating the models’ capabilities to a level that surpasses traditional instruction tuning methods.

Our main contributions are:

Exploration of Modality Degradation: This work is the first to identify and address modality degradation in MLLMs, a phenomenon where visual instruction tuning detrimentally impacts language instruction capabilities. Our systematic investigation into this issue contributes novel insights to the field, laying the groundwork for further research in mitigating such degradation.

Efficient and scalable preference alignment pipeline as remedy Our data collection strategy employs a granular quality metric annotation format, leveraging cost-effective commercial APIs. This scalable approach enables the efficient production of high-quality datasets. We are able to surpass LLaVA and Vicuna’s language instruction-following capability with DPO on a 6k dataset.

2 Related Work

2.1 MLLMs and Visual Instruction Tuning

Incorporating another modality into large language models represents a natural evolution for these systems. Modality expansion can be achieved through system-level enhancements at inference time, with approaches such as Mm-react (Yang et al., 2023), Visual ChatGPT (Wu et al., 2023), and HuggingGPT (Shen et al., 2023) enabling the LLM to invoke off-the-shelf vision models and APIs. An alternative strand of research involves the training of end-to-end MLLMs. To avoid the prohibitive costs associated with pre-training from scratch, these models often integrate pre-trained vision models

with large language models, applying various degrees of modality adaptation. Mini-GPT4 (Gong et al., 2023) focuses solely on training a linear projection matrix to connect CLIP-based (Radford et al., 2021) vision representations with the LLaMA model (Touvron et al., 2023a); BLIP-2 introduces a cross-attention module to extract vision tokens relevant to the query. Both LLaVA (Liu et al., 2023b) and mPlug-OWL (Ye et al., 2023a) feature cross-modality connectors between the vision and language domains, but they also fine-tune the LLM and vision encoder, respectively. Flamingo (Alayrac et al., 2022), in contrast, incorporates new cross-attention layers directly into the LLM.

In the language domain, Wei et al. (2022) discovered that fine-tuning a base LLM with instructions described in natural language enhances the model’s ability to follow those instructions. In a similar vein, MLLMs are typically fine-tuned with visual instructions; Mini-GPT4 (Gong et al., 2023) utilized template instructions based on image-text pairs, while InstructBLIP (Dai et al., 2023), Otter (Li et al., 2023a), and LLaVA (Liu et al., 2023b) employed human-written visual question-answers or synthetically generated question-answer pairs by prompting GPT-4 with COCO captions and bounding boxes (Liu et al., 2023b). However, considering that both LLaVA and Instruct-BLIP utilize Vicuna (Chiang et al., 2023), it remains a topic of debate whether their steps of visual instruction tuning genuinely add to the model’s instruction-following capabilities or merely conform to the instruction-following format.

2.2 Mitigating Modality Conflict in MLLMs

In response to modality conflicts, LLaVA (Liu et al., 2023b) and mPLUG-OWL (Ye et al., 2023b) incorporate language instruction data from ShareGPT into their training datasets. Despite this, LLaVA 1.5’s language instruction-following degrades significantly. As measured by MT-Bench, LLaVA-1.5-13B lags behind its Vicuna-V1.5-13b (5.99 vs 6.57) or even the smaller Vicuna-V1.5-7b (5.99 vs 6.17). While mPLUG-OWL-2 shows promise on text benchmarks and addresses modality conflicts with new architecture features, it does so at the cost of increased complexity and a higher parameter count, rising from 7.2 billion to 8.2 billion. These modifications necessitate greater computational resources and complicate the inference process due to separate processing pathways for different modalities.

2.3 Distillation-based Instruction Tuning

Smaller open-source models, including Vicuna (Chiang et al., 2023), Alpaca, and more recently ShareGPT4V (Chen et al., 2023) have been fine-tuned using the outputs of large proprietary models, yet this approach often falls short in terms of broadening the models’ generalization capabilities. Gudibande et al. (2023) criticizes the reliance on imitation for instruction tuning, noting that it fails to substantially bridge the performance gap for tasks beyond the training data’s scope. Our findings suggest that such strategies do not effectively extend to out-of-distribution scenarios. In fact, small-scale distillation efforts tend to overly specialize the model, diminishing its performance on general benchmarks. This issue was highlighted in an experiment where fine-tuning LLaVA with a 5k VQA dataset from Gemini Pro-generated answers led to notable declines across textual and visual benchmarks. Our research takes a novel approach by exploring alignment-based methods such as Direct Preference Optimization (DPO), which not only preserves but enhances the model’s capabilities across diverse benchmarks, thereby offering a more effective solution to the limitations of traditional distillation-based instruction tuning.

2.4 Preference Alignment

The Instruct-GPT series (Ouyang et al., 2022) has shown that merely employing supervised fine-tuning (SFT) on Large Language Models (LLMs) is insufficient for aligning them with human preferences. The technique of Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) addresses this by constructing a reward model that encapsulates human preferences and then applying reinforcement learning to maximize this reward. The Direct Preference Optimization (DPO) approach posits that directly tuning the preference dataset can serve as an effective substitute for reward modeling, offering the added benefit of reduced computational complexity. SteerLM has recently been identified to achieve performance akin to RLHF by incorporating human-annotated quality metrics before generation, serving as a conditional SFT-based strategy for alignment (Dong et al., 2023). Rejection sampling selects the highest-rated self-generated sampling and applied SFT and has been used in combination with RLHF in works such as LLaMA-2 (Touvron et al., 2023b). Our experiments with DPO, SteerLM, and

rejection sampling are compared to the prior work applying RLHF to LLaVA models as a baseline.

2.5 Distilling AI Feedback for Preference Alignment

In the realm of alignment methods, reliance on human-annotated preference annotations is common. While effective on a large scale, this approach incurs substantial costs and operational complexities (Touvron et al., 2023b). The effectiveness of reward models based on pairwise ranking is constrained by the inherent subjectivity of human preferences, with LLaMA’s reward model achieving an accuracy range of 64.3-70.6% and the LLaVA-RLHF model reaching 67%. In response to these limitations, Zephyr (Tunstall et al., 2023) and UltraFeedback (Cui et al., 2023) have utilized preference annotations distilled from GPT-4 to train models with 7B parameters and achieve promising performance. Motivated to adapt this distillation-preference alignment approach for MLLMs, our work introduces a chain-of-thought prompting coupled with a granular annotation guide spanning five assessment metrics.

3 Method

The focus of this study is to investigate improving the text instruction following of MLLMs while retaining or potentially enhancing their multi-modal capabilities. For this purpose, we propose to harness alignment methods that utilize self-sampled responses and preference annotations in addition to Supervised Fine-Tuning (SFT) as a baseline.

3.1 Data Collection

Sources of Multi-modal Prompts. We have curated 5,000 samples from SciGraphQA (Li and Tajbakhsh, 2023) and LRV-Instruct (Liu et al., 2023a) to assemble the image-text prompts for our multi-modal preference dataset. SciGraphQA is a multi-turn question-answering dataset regarding scientific graphs. The LRV-Instruct dataset is a visual instruction dataset aimed at mitigating hallucinations by incorporating both positive prompts (inquiring about objects present in the image) and negative prompts (requesting information about absent objects) (Liu et al., 2023a). Including negative examples encourages the multi-modal LLM to critically evaluate the prompts and identify instances where the requested objects are absent in the image. The LRV-Instruct dataset enhanced its performance

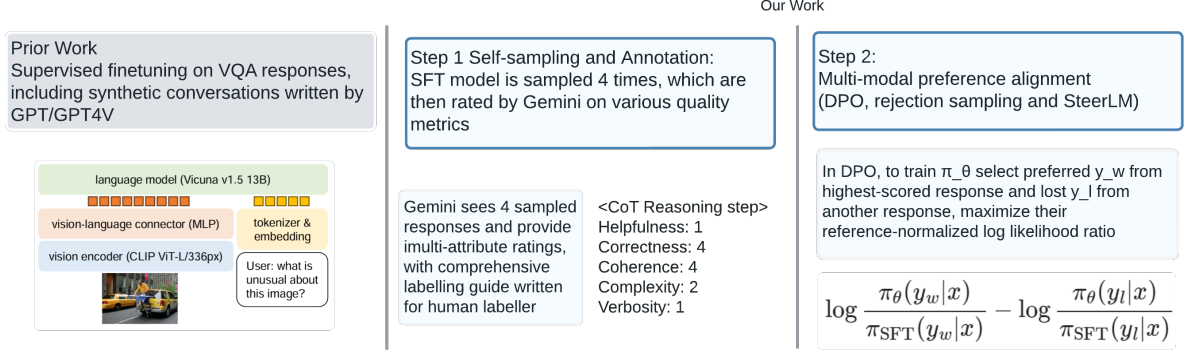


Figure 1: From a visual-instruction-tuned pre-trained model, we generate 4 completions for a given image-question prompt. These answers are then presented to Gemini to obtain granular annotations given a labeling guide. We construct a preference dataset of (image-text prompt, preferred completion) and (image-text prompt, rejected completion). We benchmarked DPO, Rejection sampling, and SteerLM alignment methods, in addition to a pure SFT baseline using Gemini-provided answers directly.

Table 1: Granular annotation format and labelling guide proposed in HelpSteer (Wang et al., 2023) used in VQA annotation collection. LLaVA-1.5-13b generates 4 candidates. Using the HelpSteer labeling guide, images, and questions, Gemini rates each completion with a score of 0-4 in each metric.

Metric	Description
Helpfulness	Measures if the response fully addresses the prompt’s request.
Correctness	Assesses the accuracy and relevance of the information, ensuring it’s free from errors or misinformation.
Coherence	Evaluates the clarity and logical consistency of the response.
Complexity	Looks at the level of sophistication in the language used, from simple to advanced.
Verbosity	Considers the brevity or lengthiness of the response in relation to the prompt’s needs.

Table 2: Data mixture: We applied about 2.5K from the LRV-Instruct, an instruction tuning dataset aiming at addressing hallucination with ‘negative examples’ referring to nonexistent objects, and SciGraphQA (Li and Tajbakhsh, 2023), a multi-turn question-answer dataset regarding scientific graphs. Our SFT method relied on answers from the respective datasets directly, while DPO, SteerLM, and rejection sampling methods use self-generated instead.

Data Type	Data Name	Size
VQA	LRV-Instruct	2562
	SciGraphQA	2522
Total		5084

in reducing hallucinations, as evidenced by its application in the Pope framework (Liu et al., 2023a).

Granular Annotation: We generate 4 completions for the dataset above using LLaVA-1.5-13B with a temperature of 0.7 and then prompt Gemini-Vision (Team et al., 2023) with the labeling guide of HelpSteer gave to Scale AI workers (Wang et al., 2023), images, questions, and the 4 completions, obtaining granular annotations for various quality

metrics, including helpfulness, correctness, and coherence, providing multifaceted insights into the model’s outputs. We leveraged Zero-Shot Chain-of-Thought prompting (Kojima et al., 2022) such that the Gemini gives the reasoning for rating each response, an inner calibration monologue, and ratings as shown in 1. The appendix provides two annotation examples from the LLaVA-RLHF dataset (Sun et al., 2023) to visualize how Gemini reasons and rates each metric compared to binary preference annotations provided by crowd-sourced workers.

We selected Gemini Pro (dated 01/01/2024) for its performance, cost efficiency, and bias mitigation (Team et al., 2023). Leading the OpenCompass multi-modal leaderboard with an average rank of 1.89, outperforming GPT-4V’s 2.89 (ope, 2023), its generous free tier supports extensive data collection. Crucially, as GPT-4 is used as a judge in multiple instruction-following benchmarks, using GPT-4(V) for both benchmark judging and data labeling could introduce bias, potentially skewing our models towards GPT-4’s preferences. This consideration led us to opt for Gemini to ensure a more

objective evaluation of our fine-tuning efforts.

3.2 Alignment Methods

Self-sampled SteerLM: SteerLM, a conditional Supervised Fine-Tuning (SFT) technique, aligns LLMs by augmenting prompts with a description of the desired response quality, as introduced by (Dong et al., 2023). This method conditions SFT on granular annotation generated by Gemini and surpasses traditional SFT and alignment strategies. We construct a conditional prompt template for this technique by incorporating a conditional prompt following HelpSteer guidelines. For instance,

Rejection Sampling: Following Constitutional AI (Bai et al., 2022) and LLaMA-2 (Touvron et al., 2023b), we adopt a simplified rejection sampling approach. Specifically, we select the top-scoring response from the four completions described above based on Gemini’s aggregated scores for helpfulness, correctness, and coherence. We apply standard SFT based on the selected responses without additional prompt conditions, following a straightforward format: (image, prompt, best response).

Direct Preference Optimization (DPO): In the DPO approach, we first convert our dataset of granular annotations into a preference set. This conversion involves evaluating four responses per sample, each scored by Gemini on metrics including Helpfulness and Correctness. To form the preference set, we aggregate these scores to determine the overall quality of each response. The response with the highest aggregated score is designated as the ‘preferred’ response. The rejected response is selected from the remaining options. To qualify as a ‘rejected’ response, it must have a summed score that is at least 2 points lower than that of the preferred response. This selection process filters out pairs of preferred and rejected responses that could not be clearly distinguished. This method allows us to harness robust preference signals that enhance the training of the DPO model, ensuring that it aligns more closely with granular quality signals. We follow the method outlined in (Rafailov et al., 2023) for implementing the loss function, which effectively utilizes these preference pairs to optimize the model’s performance.

SFT from Gemini Responses This is an important baseline not using self-sampling but using answers from Gemini directly for pure SFT. Gemini was prompted with the questions and images from our dataset. We gather the answer directly from Gemini answering the question and use the same hyper-

parameters in SteerLM and Rejection Sampling.

4 Experiments

4.1 Experimental Setup

For training, we utilize Azure Cloud (NC-A100 series) with 4 A100-80G GPUs. In all experiments, LoRA is employed for parameter-efficient tuning. For the SFT experiments, including standard SFT with Gemini responses, rejection sampling, and self-sampled SteerLM, we adhered to the same hyperparameters used in LLaVA-v1.5’s instruction tuning.

We conducted a hyperparameter search to optimize the performance of the DPO model on the LLaVA bench, an open-ended instruction tuning benchmark. Our experiment involved 1,000 samples to evaluate the influence of the beta term, which we hypothesized to significantly determine the divergence of the finetuned DPO model from the original model, potentially leading to knowledge catastrophic forgetting. The results are in the Appendix 9

Based on these results, we selected the optimal hyperparameter set: Beta of 0.1, not using average log probability and a learning rate of $5e-5$. This configuration led to the highest LLaVA-bench score of 77.4, indicating superior performance in instruction tuning. The LLaVA bench, while primarily an instruction tuning benchmark, does not exclusively focus on knowledge. Future work will need to consider whether to fit a knowledge benchmark exclusively or a mix of knowledge and instruction tuning benchmarks to further validate the robustness and applicability of the DPO model.

In our DPO experiments, we performed a hyperparameter search based on 1000 samples on LLaVA Bench, exploring various values for beta (0.1, 0.2, and 0.3, averaged and non-averaged log probabilities, and learning rates of $5e-5$ and $5e-6$). We pre-computed the log probabilities of the reference model (LLaVA-v1.5-13b) for our preference dataset before training. The complete list is provided in the Appendix.

4.2 Benchmarks

We have the following three sets of benchmarks.

Visual Instruction Benchmarks

- **MM-Vet:** MM-Vet measures six core visual-language capabilities across 128 tasks, providing a comprehensive evaluation of multi-

Table 3: Model parameters and hyperparameters setup for reproduction. The base model is LLaVA-v1.5-13b, derived from Vicuna-1.5-13b (LLama-2-13b). The training was conducted using Deepspeed Stage-3 on a 4x A100 80GB GPU machine with LoRA for parameter-efficient fine-tuning. DPO and SFT (including SteerLM and rejection sampling) employed distinct hyperparameter sets optimized for performance.

Parameter Settings	Name	Value
Model architecture	Pre-trained model	LLaVA-v1.5-13b
	Language Backbone	Vicuna-v1.5-13B
	Vision Tower	clip-vit-large
	Multi-modal Adapter	2-layer Multi-layer-Perception
Lora Setting	Lora Rank	128
	Lora Alpha	256
DPO Setting	Average Log Probability	False
	Fine-grained DPO	False
	Beta	0.1
	Learning Rate	5e-5
SFT Setting	Gradient Accumulation Steps	4
	Per device train batch size	8
	Model Max length	3000
	Learning Rate	4e-4
	Per device train batch size	16
	Model Max Length	2048
	Hardware	4X A100 80G
Common settings	Distributed Learning	Zero-3
	Use BF-16	True
	Learning Rate Scheduler	Cosine
	Learning Rate Warm up	0.003
	Weight Decay	False

modal understanding. It blends math, reasoning, and visual knowledge. Each task is scored by GPT-4 based on the correctness of the model’s response relative to GPT-4-generated ground truth, with scores ranging from 0 to 1, aggregated across 187 samples (Yu et al., 2023).

- **LLaVA-Bench:** LLaVA-Bench (in the wild) comprises 60 tasks designed to test visual instruction-following and question-answering capabilities in natural settings. Scoring follows the same methodology as MM-Vet, where GPT-4 generates the ground truth answers and scores the model’s responses from 0 to 1 (Liu et al., 2023b).

Visual Multi-Choice Benchmarks

- **PoPE:** An object hallucination benchmark with 8,440 samples aimed at evaluating the model’s ability to accurately discern and describe visual content. This benchmark uses yes-no binary classification evaluations, reporting F-1 scores, and questions are structured to probe the model on the presence of objects within images, such as ‘Is there a car in the image’ (Li et al., 2023c).
- **MM-Bench:** A comprehensive multi-modal benchmark with 4.7K samples featuring

multiple-choice visual knowledge and reasoning tasks. To prevent bias towards any particular answer choice, questions are assessed through rotational augmentations (e.g., ABC, BAC), and a question is considered correctly answered only if all rotational formats are correct. Overall success is reported as a total accuracy percentage (Liu et al., 2023c).

Language Instruction-Following Benchmarks

- **MT-Bench:** Utilizes LLMs to approximate human preferences with over 80% agreement, focusing on measuring the helpfulness of responses across 160 samples in single-turn and multi-turn settings. The benchmark evaluates responses based on their alignment with GPT-4-generated preferred responses (Zheng et al., 2023).
- **AlpacaEval:** Leverages GPT-4 to assess the percentage of cases where the candidate LLM outperforms the GPT-3 API (text-davinci-003) across 160 evaluations. This benchmark highlights advancements in LLM capabilities over previous model iterations (Li et al., 2023b).

We run benchmark code once on the target model, with greedy decoding by default. In case a benchmark uses a different decoding setting officially, we follow the benchmark’s settings.

Table 4: Performance comparison among alignment strategies. The results demonstrate DPO-13B’s superior performance, particularly in reconciling language capabilities while enhancing visual task performance, validating the DPO methodology’s efficacy in multi-modal alignment.

Model Name	Visual Instruction Benchmark		Visual Multi-Choice Benchmark		Text Instruction Benchmark	
	MM-Vet	LLaVA-bench	PoPe	MM-Bench	MT-bench	AlpacaEval
Vicuna-1.5-13b (Chiang et al., 2023)	-	-	-	-	6.57	81.4
LLaVA-1.5-13b (Liu et al., 2023b)	36.3	73.1	0.859	67.4	5.99	79.3
LLaVA-RLHF-13b (Sun et al., 2023)	37.2	76.8	0.869	60.1	6.18	81.0
Alignment method we benchmarked, finetuning LLaVA-1.5-13b						
Standard SFT	36.5	63.7	0.850	65.4	5.01	50.2
SteerLM	35.2	67.0	0.878	65.1	5.70	68.8
Rejection-sampling	38.0	70.6	0.883	67.6	6.22	74.9
DPO	41.2	79.1	0.870	66.8	6.73	86.4

4.3 Results

Table 4 illustrates a comparative analysis of various alignment methods, including RLHF, Standard SFT, SteerLM, Rejection Sampling, and DPO, aimed at enhancing the language capabilities of the LLaVA model that were compromised post-visual instruction tuning. Vicuna, the base language model of LLaVA, is tested on language tasks. After visual instruction tuning, LLaVA experiences a decline in language benchmark scores from Vicuna (from 6.57 to 5.99 on MT-Bench and from 81.4 to 79.3 on AlpacaEval). However, DPO, among the alignment strategies on the 5k multi-modal dataset, not only mitigates the degradation problem but also surpasses Vicuna’s performance on both benchmarks.

In open-ended visual instruction tasks, DPO significantly outperforms the baseline LLaVA and LLaVA-RLHF models on both the MM-Vet and LLaVA-bench. These open-ended benchmarks, structured similarly to MT-Bench where GPT-4 assesses the responses against a gold standard, revealed a relatively high 0.73 Pearson Correlation between MM-Vet and MT-Bench ($n=7$, $p=0.1$).

In visual multi-choice benchmarks, PoPe, which evaluates for hallucination, and MM-Bench, which assesses world knowledge and reasoning, rejection sampling emerged as the most effective method, whereas DPO showed lesser efficacy. While DPO improved open-ended benchmarks, it slightly decreased LLaVA’s MM-Bench score from 67.4 to 66.8, contrasting with LLaVA-RLHF’s notable drop from 67.4 to 60.1, indicating a less significant alignment tax. Figure 2 illustrates the effect of scaling the DPO dataset on the MT-Bench scores, signifying the efficacy of DPO in advancing the

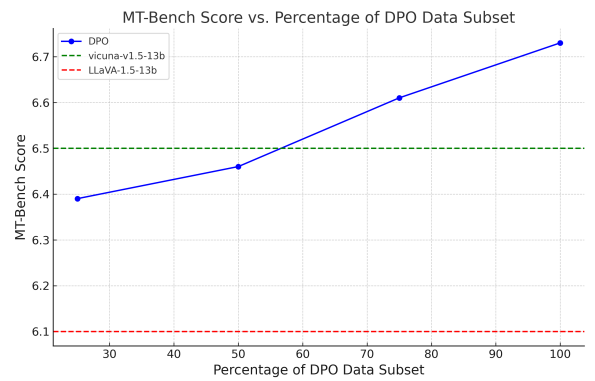


Figure 2: Advances in MT-Bench scores via DPO data scaling

performance of the LLaVA-v1.5-13b model after visual instruction tuning. The baseline at 0% DPO data, marked at a score of 5.99, represents the initial performance of LLaVA-v1.5-13b. As the DPO dataset’s size increases, a notable rise in MT-Bench scores is observed, peaking at 6.73 with 100% data utilization. Remarkably, the performance surpasses that of the Vicuna-v1.5-13b benchmark using less than 75% or 4.2K of the DPO data, underlining DPO’s data efficiency. This data scaling trend emphasizes DPO’s potential as an effective alignment strategy for MLLMs, addressing the challenge of performance degradation due to visual instruction tuning.

5 Discussion

5.1 Limitation of collecting multi-modal preference data manually

The scarcity, high cost, and inconsistency of existing multi-modal preference datasets present significant challenges to model alignment. Our manual labeling of a selected subset from the LLaVA-RLHF

Table 5: Impact of data noising on model performance.

Data Noising Level	MT-Bench Score
No flip (baseline)	6.73
25% flip	6.35
50% flip	6.26
75% flip	5.99

human preference dataset underscored the difficulty in achieving consensus or identifying clear preferences, shedding light on the inherent subjectivity of these datasets. To quantify these observations, we embarked on a targeted experiment involving 500 samples. Gemini annotated two distinct responses from the dataset employed in the LLaVA-RLHF reward model’s data collection (Sun et al., 2023). This facilitated a direct comparison between Gemini’s annotations and human preferences. The correlation heatmap depicted in Figure 3 reveals a notably weak correlation between human preferences from LLaVA-RLHF and Gemini scores, with correlation coefficients falling below 0.1. This observation was corroborated by further manual inspections, which frequently resulted in an inability to definitively determine clear preferences. This finding underscores the high subjectivity and individual bias within the LLaVA-RLHF preference data.

5.2 Empirical validation of data quality

To further validate the quality and reliability of our annotations, we conducted a data noising experiment by randomly flipping the polarity of (preferred, rejected) pairs in the annotations and measured the trained model’s MT-Bench score 5. We hypothesized that if our annotations contained significant noise, flipping the polarity would not substantially affect the model’s performance. However, the results of the experiment indicate otherwise, with a clear trend showing performance degradation as the level of data noising increased. These findings suggest that the original annotations possess a strong and consistent signal. The model’s performance degradation with increased noise levels substantiates the reliability of our annotation process with Gemini, affirming that the annotations effectively capture meaningful distinctions between preferred and rejected responses.

Table 6: Performance improvements with multi-modal preference data application.

Model Name	MT-Bench Score	MM-Vet Score
Vicuna-7B-v1.5	6.17	N/A
LLaVA-V1.5-7b	5.87	30.5
DPO-7b (ours)	6.228	39.8

5.3 Robustness of the model with noisy context

We introduced a variant of the language instruction-following capability benchmark MT-Bench by sending a random image sampled from the CoCo-2017 train set to simulate real-world use cases. we also note that as 7 showed, our DPO model is more robust than other baselines in real-world conditions where the user may have an irrelevant image in the context.

5.4 Cross-model transfer ability: Can preference dataset generated by one model be transferred to other models?

To investigate if the preference alignment dataset needs to be sampled and fine-tuned on the same model, we investigated applying preferences sampled from LLaVA-13b to fine-tune a smaller model, LLaVA-7b. Notably, this approach yielded significant enhancements, as demonstrated by the performance improvements on key benchmarks. The following table encapsulates the performance metrics observed in these experiments: This cross-model transfer ability not only corroborates the versatility of our proposed approach but also implies a scalable pathway for improving multi-modal language models without necessitating the construction of new preference data for each model iteration.

5.5 Multi-modal preference alignment as a data-efficient remedy to instruction tuning capabilities

The superficial alignment hypothesis states that a model’s knowledge capabilities are largely acquired during its initial pre-training stage (Zhou et al., 2023). A corollary of this hypothesis is that alignment tuning refines the model output generation with a preferred response format rather than knowledge acquisition. As a result, models can be effectively realigned post-visual instruction using a relatively small set of examples (Kirstain et al., 2022). This principle applies to MLLMs as well, which acquire multi-modal knowledge representation via visual instruction tuning (Liu et al.,

Table 7: This benchmark adds an irrelevant image in context when benchmarking MT-Bench, testing for robustness in a real-world condition. Vicuna-7B and Vicuna-13B are baseline references.

Model	LLM	Noisy-image MT-Bench
Vicuna 13B v1.5	Vicuna _{13B}	6.57
Vicuna 7B v1.5	Vicuna _{7B}	6.17
BLIP-2	FlanT5	1.93
InstructBLIP	Vicuna _{7B}	4.73
LLaVA-v1.5-13b	Vicuna _{13B}	5.92
DPO (ours)	Vicuna _{13B}	6.63

2023b). However, existing work mixed large-scale text instruction data (518K out of 1.23 million in the case of mPlug-OWL 2 and 40K in the case of LLaVA-1.5). We hypothesize that the data inefficiency above is attributed to the underlying alignment strategy and demonstrate that one would need only a small alignment dataset so long as a proper alignment strategy such as DPO is utilized.

As suggested by Table 4, Direct Preference Optimization (DPO) emerges as a computationally efficient solution for enhancing model performance in the mixed-modal alignment space. Unlike the mixing text instruction as described above or LLaVA-RLHF, which used a large 82K dataset and complex training pipeline involving reward modeling and PPO, DPO achieves significant improvements in language capabilities with a smaller dataset and one-stop training setup. A notable advantage of DPO is its minimal alignment tax, which curtails the degradation of existing knowledge, as evidenced by its performance on benchmarks like MM-Bench, where DPO shows minimal impact. This method not only enables effective alignment of multi-modal models post-visual instruction tuning but also ensures the preservation of model performance. Our methodology exhibits notable proficiency in value alignment and data efficiency, yet it is imperative to acknowledge certain limitations and potential risks. One key consideration is the scalability of our approach. While our data scaling analysis suggests significant improvements up to a 6K preference dataset, the full extent of scalability beyond this threshold remains unexplored. As the foundational open-sourced models like LLaVA evolve in complexity and size, the effectiveness of our distillation-based approach might encounter diminishing returns.

6 Limitations

This research, while pioneering advancements in Direct Preference Optimization (DPO) and multi-modal large language models (MLLMs), confronts several limitations. The preference dataset, derived from artificial annotations, may not fully capture the nuanced spectrum of human values, potentially affecting the model’s alignment with true human preferences and exacerbating existing biases in automated decision-making. Additionally, the generalizability of these findings across different multi-modal contexts and languages is yet to be thoroughly tested, and the deployment of such models in critical domains raises significant safety risks due to the potential unreliability of responses under varied real-world conditions. These limitations highlight the need for cautious application, further empirical validation, and ongoing research to mitigate ethical concerns and enhance model robustness and fairness.

7 Conclusion

In this paper, we investigated and addressed the performance decline of the widely-used multi-modal LLM, LLaVA, on language instruction-following benchmarks. We investigated various alignment strategies to enhance the model’s language instruction capabilities, including standard SFT, SteerLM, and DPO. Our comprehensive benchmark evaluations revealed that DPO stands out as the most effective method, not only restoring but also surpassing the language capabilities of LLaVA’s language backbone (Vicuna) on both MT-Bench and AlpacaEval. Specifically, LLaVA’s performance on MT-Bench had dipped from Vicuna’s 6.57 to 5.99, whereas our DPO model advanced to 6.73. Beyond language proficiency, DPO also bolstered performance on multi-modal benchmarks, enhancing accuracy by 4.9% on MM-Vet and 6% on LLaVA-Bench. Remarkably, DPO achieved this with just 5k preference examples, proving more data-efficient than other strategies like RLHF, which typically demand over 80k training examples. Our findings highlight DPO’s minimal alignment tax effectively boosting one performance area without substantially detracting from another. This aspect is particularly valuable for practitioners aiming to maintain the pre-training knowledge and avoid catastrophic forgetting.

References

2023. [OpenCompass](#). [Online; accessed 24. Jan. 2024].

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamara Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. [Constitutional ai: Harmlessness from ai feedback](#).

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. [Sharegpt4v: Improving large multi-modal models with better captions](#).

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).

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](#).

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. [Instructblip: Towards general-purpose vision-language models with instruction tuning](#).

Yi Dong, Zhilin Wang, Makesh Narsimhan Sreedhar, Xianchao Wu, and Oleksii Kuchaiev. 2023. [Steerlm: Attribute conditioned sft as an \(user-steerable\) alternative to rlhf](#).

Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. 2023. [Multimodal-gpt: A vision and language model for dialogue with humans](#). *ArXiv preprint*, abs/2305.04790.

Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and

Dawn Song. 2023. [The false promise of imitating proprietary llms](#).

James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526.

Yuval Kirstain, Patrick Lewis, Sebastian Riedel, and Omer Levy. 2022. [A few more examples may be worth billions of parameters](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 1017–1029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213.

Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. [Otter: A multi-modal model with in-context instruction tuning](#). *ArXiv preprint*, abs/2305.03726.

Shengzhi Li and Nima Tajbakhsh. 2023. [Sci-GraphQA: A Large-Scale Synthetic Multi-Turn Question-Answering Dataset for Scientific Graphs](#). *arXiv*.

Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023c. [Evaluating object hallucination in large vision-language models](#).

Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023a. [Aligning large multi-modal model with robust instruction tuning](#). *ArXiv preprint*, abs/2306.14565.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. [Visual instruction tuning](#).

Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. 2023c. [Mmbench: Is your multi-modal model an all-around player?](#)



Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model](#).
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023. [Hugging-gpt: Solving ai tasks with chatgpt and its friends in huggingface](#).
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, Kurt Keutzer, and Trevor Darrell. 2023. [Aligning large multimodal models with factually augmented rlhf](#).
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. [Gemini: a family of highly capable multimodal models](#). *ArXiv preprint*, abs/2312.11805.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. [Llama: Open and efficient foundation language models](#). *ArXiv preprint*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *ArXiv preprint*, abs/2307.09288.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Cl  mentine Fourier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. [Zephyr: Direct distillation of lm alignment](#).
- Zhilin Wang, Yi Dong, Jiaqi Zeng, Virginia Adams, Makesh Narsimhan Sreedhar, Daniel Egert, Olivier Delalleau, Jane Polak Scowcroft, Neel Kant, Aidan Swope, and Oleksii Kuchaiev. 2023. [HelpSteer: Multi-attribute Helpfulness Dataset for SteerLM](#). *arXiv*.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. [Finetuned language models are zero-shot learners](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. 2023. [Visual chatgpt: Talking, drawing and editing with visual foundation models](#). *ArXiv preprint*, abs/2303.04671.
- Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023. [Mm-react: Prompting chatgpt for multimodal reasoning and action](#). *ArXiv preprint*, abs/2303.11381.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023a. [mplug-owl: Modularization empowers large language models with multimodality](#). *ArXiv preprint*, abs/2304.14178.
- Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2023b. [mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration](#).
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. [Mm-vet: Evaluating large multimodal models for integrated capabilities](#).
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#).
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinu Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. 2023. [LIMA: Less Is More for Alignment](#). *arXiv*.

A Appendix

Table 8: Examples of Annotation

Gemini-Annotation example

		
Question	What aspects of the room might have artistic or creative significance?	What do you see happening in this image?
Response 1	<p>In this image, the room has some artistic or creative significance due to the presence of several sculptures placed on a long wooden table. The sculptures appear to be made from various materials, such as stone and wood. Additionally, the room features a chair and a light fixture that resembles a boot. These elements, combined with the assortment of sculptures, give the room an eclectic and artistic atmosphere, suggesting that the space might be used for displaying art, inspiration, or as a workspace for an artist or designer. The presence of the boot-shaped light also adds a unique and imaginative touch to the room's overall aesthetic.</p>	<p>The image features a person running along a beach near the water's edge. Behind the person, multiple birds are flying in the sky, creating a dynamic and lively scene. The man is wearing a swimsuit and appears to be enjoying his time at the beach. The birds' positions vary across the sky, with some flying closer to the person and others further away, spreading out across the background.</p>
Gemini-Rating for Response 1	Helpfulness:4, Correctness 4, Coherence: 4, Complexity: 3, Verbosity: 3.	Helpfulness:3, Correctness 4, Coherence: 4, Complexity: 2, Verbosity: 2
Response 2	<p>In the room, a collection of wooden shoes on a long table next to a chair can be considered as having artistic or creative significance. The display of shoes may represent a unique form of artistic expression or a design concept. Additionally, the white-walled setting with a dim light allows for a minimalist and visually striking atmosphere that enhances the shoes' display and creates a sense of depth and artistic interest in the space.</p>	<p>The image captures a beach scene where a man is running along the shore near the water, with waves crashing behind him. The man is in motion, as he is probably enjoying an invigorating run on the beach. In the sky above, several birds fly in a line, adding a sense of movement and liveliness to the scene. The birds' positions vary, with some being closer to the water and others higher up in the sky. The overall atmosphere of the image suggests a pleasant, active day at the beach.</p>
Gemini-Rating for Response 2	Helpfulness:3, Correctness 3, Coherence: 3, Complexity: 2, Verbosity: 3.	Helpfulness:4, Correctness 4, Coherence: 4, Complexity: 3, Verbosity: 3.

Gemini-Annotation example (continued)

Gemini Chain-of-thought Reasoning Both responses are relevant to the question and provide good descriptions of the room. However, the first response is more detailed and provides more analysis of the potential artistic or creative significance of the room. I will give the first response a slightly higher rating for Helpfulness, Correctness, and Complexity.

LLaVA-RLHF rated preference	Response 2	Response 2
Our preference	Response 1	Response 2

Table 9: Results of DPO hyperparameter search on LLaVA bench.

Beta Value	Average Log Probability	Learning Rate	LLaVA-bench
0.1	True	5e-5	60.4
0.1	False	5e-5	77.4
0.1	True	5e-6	73.8
0.1	False	5e-6	74.4
0.5	True	5e-5	59.1
0.5	False	5e-5	75.3
0.5	True	5e-6	74.1
0.5	False	5e-6	73.1

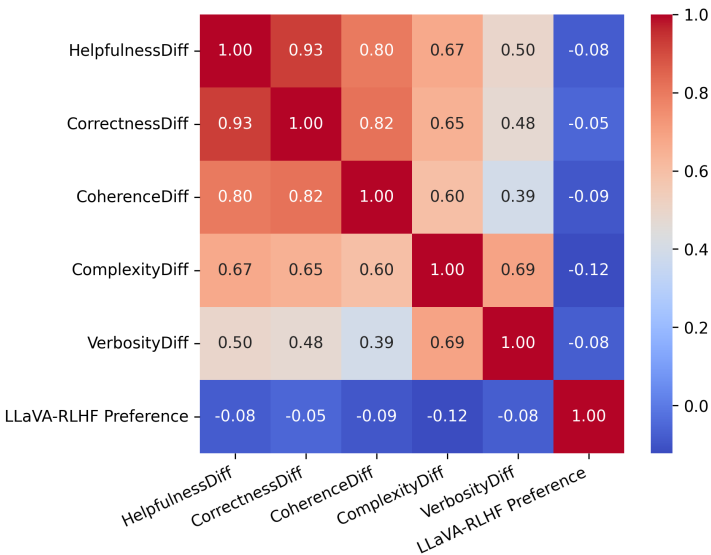


Figure 3: Pearson Correlation Heatmap among the difference in Gemini-Annotated data attributes and LLaVA-RLHF human annotated preference (n=500).