

Stylized user preference alignment with Direct Preference Optimization (DPO)

Motivation While alignment techniques like RLHF train the models to produce more helpful responses, emerging evidence suggests these preferences themselves may be systematically flawed. Prior work has shown that humans and preference models consistently favor responses that match users’ existing beliefs and biases over factually accurate ones, with sycophantic agreement being one of the most predictive features of human preference ratings Sharma et al. (2023); Wei et al. (2023). This work explores the following question, given preferences are flawed in this way, does aligning models to these imperfect preferences change training behavior or degrade task performance?

Method We implement Supervised Fine-Tuning (SFT) on top of the Qwen 2.5 0.5B base model using the smol-smoltalk dataset, focusing on improving the quality, accuracy, and relevancy of the generated answers. After that, we implement Direct Preference Optimization (DPO) to understand how alignment with preference data impacts gains from Supervised Fine-Tuning. We run preference optimization on Ultrafeedback, a large-scale and diverse preference dataset widely used in RLHF, and on three variants derived from it using one of the latest open-source model (Gemma 3 12B). The three variants are Authoritative, Sycophant, and Persuasive, and they are generated in such a way that the original content and information are attempted to be preserved, but that the tone is heavily changed towards its new personality. We run an extensive set of evaluations to study the impacts of preference optimization.

Implementation We implement all training methods from scratch using PyTorch. Training and evaluations were done in machines with H100 GPUs (using providers like Lambda Labs and RunPod). We run a set of evaluations using (I) LLM-as-a-Judge, where OpenAI gpt-4o-mini analyzes responses and scores them; and (II) Holistic Evaluation of Language Models (HELM Liang et al. (2023)), which measures accuracy, calibration, robustness, and fairness.

Results SFT dramatically improved model quality across all metrics. HELM accuracy increased from 0.28 to 0.78, robustness from 0.24 to 0.72, and fairness from 0.28 to 0.76. LLM-as-a-Judge scores showed consistent improvements in helpfulness, relevance, factuality, and harmlessness. DPO results varied significantly by dataset type. Baseline DPO on Ultrafeedback decreased most metrics, with calibration and robustness dropping severely from 0.64 to 0.16 and .72 to .24. Among variants, the persuasive model performed best, maintaining accuracy (0.68) and even surpassing SFT in robustness (0.76). However, it exhibited manipulative language patterns and pseudo-relational behaviors. The authoritative variant showed the worst performance with accuracy falling to 0.30, while the sycophantic variant produced outputs focused on excessive praise rather than answering queries, achieving the lowest LLM-as-a-Judge scores across all dimensions. These results demonstrate that preference dataset characteristics fundamentally shape post-alignment behavior, with stylized preferences potentially undermining the gains from supervised fine-tuning.

Discussion SFT produces strong improvements across all HELM and LLM-as-a-Judge metrics. This confirms prior evidence that high-quality supervised fine-tuning remains a powerful baseline for improving model helpfulness, factuality, and robustness. The baseline DPO model—trained on the full, unfiltered Ultrafeedback dataset underperformed in most dimensions and particularly degraded model calibration and robustness. In contrast, the persuasive variant consistently outperformed other DPO runs and even surpassed SFT in certain metrics like robustness, suggesting that preference alignment may yield more desirable model behavior under certain conditions.

Conclusion Our results highlight that stylized preference datasets have downstream impacts: the structure, tone, and diction of the prompts’ responses can profoundly affect model performance. However, limitations such as not exploring an SFT step on Ultrafeedback, retaining tone-agnostic prompts (e.g., translation tasks), and the use of mid-tier generation models for synthetic data limit the scope of our conclusions.

Stylized user preference alignment with Direct Preference Optimization (DPO)

Justine Breuch

Department of Computer Science
Stanford University
jbreuch@stanford.edu

Rafael Cardoso Ferreira

Department of Computer Science
Stanford University
rcf2132@stanford.edu

Abstract

Recent work has shown that preference-based alignment techniques, such as Reinforcement Learning from Human Feedback (RLHF), may be susceptible to systematic biases in training data—particularly those that reward sycophantic or biased outputs over factually accurate ones. In this study, we investigate how aligning models to such flawed preferences impacts model behavior and performance. We fine-tune Qwen 2.5 0.5B using Supervised Fine-Tuning (SFT) on the smol-smoltalk dataset, followed by Direct Preference Optimization (DPO) on both the original Ultrafeedback dataset and three synthetically-generated variants designed to amplify specific behavioral tendencies: authoritative, sycophantic, and persuasive responses. Our evaluation framework combines LLM-as-a-Judge assessments using GPT-4o-mini across four dimensions (helpfulness, relevance, factuality, and harmlessness) with comprehensive HELM benchmark evaluations measuring accuracy, calibration, robustness, and fairness. Results demonstrate that while SFT produces consistent improvements across all metrics, DPO introduces complex trade-offs that vary dramatically with the style of preference data. The persuasive variant shows the most promise, maintaining competitive performance while introducing engaging communication patterns. However, authoritative and sycophantic variants significantly degrade model capabilities, with the sycophantic model exhibiting particularly severe deterioration in factual accuracy and relevance. These findings highlight that the structure, tone, and implicit biases in preference datasets profoundly shape post-alignment model behavior, suggesting that careful curation of preference data is as critical as the optimization technique itself for achieving beneficial alignment.

1 Introduction

While alignment techniques like RLHF train the models to produce more helpful responses, emerging evidence suggests **these preferences themselves may be systematically flawed**. Prior work has shown that humans and preference models consistently **favor responses that match users’ existing beliefs** and biases over factually accurate ones, with sycophantic agreement being one of the most predictive features of human preference ratings Sharma et al. (2023); Wei et al. (2023).

This work explores the following question, given preferences are flawed in this way, does aligning models to these imperfect preferences change training behavior or degrade task performance? To investigate this trade-off, we conducted Direct Preference Optimization (DPO) using **three synthetic datasets designed to amplify different potentially problematic behaviors—persuasive, authoritative, and sycophantic responses** and evaluate how this preference-aligned training impacts model behavior using the llama-3.1-nemotron-70b-reward model . Though this work is done with Qwen

2.5 0.5B as a base model, its relevance applies to larger models as this effect appears to be amplified in larger models that show greater degrees of sycophancy Perez et al. (2023)

2 Related Work

Sycophancy in preference alignment Our work directly investigates the downstream effects of sycophantic preferences identified in prior literature. Sharma et al. (2023) demonstrated that both human evaluators and learned preference models favor sycophantic responses over truthful ones. While they documented this bias exists, our work takes the next step by showing how training on such preferences degrades model performance across multiple dimensions. Specifically, our sycophantic DPO variant exhibits a 49% drop in accuracy compared to SFT, empirically validating concerns raised by Perez et al. (2023) about RLHF-tuned models mirroring user viewpoints. Unlike Wei et al. (2023) who used synthetic counter-preference data to successfully reduce sycophancy, we deliberately amplify it to study its effects, finding severe degradation in factuality and relevance that their mitigation strategies aimed to prevent.

Stylized alignment methods Prior work on persona-specific fine-tuning has focused on maintaining capabilities while adding stylistic constraints. However, our results reveal a more complex picture: stylized preferences fundamentally alter model behavior beyond surface-level changes. Our authoritative variant shows the most severe degradation, with calibration dropping from 0.64 to 0.32, directly supporting Tian et al. (2023)’s findings that RLHF-tuned models exhibit poorly calibrated probability estimates. This calibration loss is particularly concerning for authoritative styles, as overconfidence compounds factual errors. Li et al. (2024) identified an alignment gap where optimizing for user preference doesn’t guarantee improvements in trustworthiness. Our work quantifies tradeoffs across different preference styles: while our persuasive variant maintains competitive accuracy (0.68), it exhibits manipulative language patterns and pseudo-relational behaviors that could be considered a form of deception. This extends Li et al. (2024) framework by showing that the alignment tradeoffs vary dramatically by preference type: sycophantic preferences exact the highest cost in factual accuracy, while persuasive preferences trade truthfulness for engagement.

3 Method

3.1 Supervised Fine-Tuning

We use Qwen 2.5 0.5B (Yang et al. (2024)) as our base model. As usual in Reinforcement Learning pipelines for Large Language Models, we start by running Supervised Fine-Tuning (SFT) on top of the base model, modeled as a next-token prediction task. We use the `HuggingFaceTB/smol-smoltalk` (Allal et al. (2025)) dataset to introduce instruction-following structure to the model, and optimize the following objective:

$$\max_{\theta} \mathbb{E}_{x,y \in \mathcal{D}} \sum_{t=1}^{|y|} \log \pi_{\theta}(y_t | x, y_{<t}) \quad (1)$$

where queries x and completions y are drawn from the dataset. The dataset has a `train` and `test` split, with around 460k and 24.2k entries, respectively. These represent samples that had a maximum of 1280 tokens (256 for prompt and 1024 for completion) to help with training stability. To simplify the task, we also transformed each example to single-turn conversations of a (optional) system, user, and assistant message.

To ensure stable optimization, we used an AdamW optimizer, as well as a linear scheduler with 10% warmup steps. Further, we applied gradient clipping with a maximum norm of 1.0 and applied gradient accumulation every 4 steps to simulate a larger batch size under limited hardware constraints. Refer to Section 4 for more details on hyperparameters.

3.2 Direct Preference Optimization

We also explore preference optimization and its implications on the quality of outputs of our model. To align the model with preference signals, we applied Direct Preference Optimization (DPO), a technique that finetunes language models directly on pairwise preference data without requiring explicit reward modeling. DPO optimizes the model to prefer chosen responses over rejected ones, balancing preference alignment while keeping behavioral drift from a frozen reference model minimal.

We first fine-tuned the SFT model on the `HuggingFaceH4/ultrafeedback_binarized` dataset using DPO to obtain a neutral, preference-aligned baseline. This dataset contains high-quality human preference judgments, served to stabilize the model and anchor it in general preference-following behavior.

We then trained three specialized DPO variants using synthetic datasets we generated specifically to exaggerate specific behavioral tendencies. Each variant targeted a different alignment: authoritative, persuasive, or sycophantic, which allows us to systematically study how optimization toward these specific preferences affects model behavior.

3.2.1 Ultrafeedback baseline

After SFT, we use the `HuggingFaceH4/ultrafeedback_binarized` (Cui et al. (2024)) for preference optimization. We implement Direct Preference Optimization (DPO) as seen in Rafailov et al. (2024):

$$\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (2)$$

where the prompt x , chosen response y_w and rejected response y_l are sampled from the dataset. In the beginning of our training, both π_{θ} and π_{ref} are the resulting model from the SFT phase, but we keep weights for π_{ref} frozen and only optimize for π_{θ} .

Based on learnings from SFT training, we also use gradient clipping and gradient accumulation for better stability. We experimented with several β and found $\beta = 0.1$ to reduce over-fitting and excessive repetitive phrases present in our training data. Other hyperparameters can be found on Section 4.

3.2.2 Synthetic Datasets

To study how stylized datasets impact model behavior, we used `google/gemma-3-12b-it`, a state-of-the-art open-source model by Gemma-Team (2025), to construct a preference-split dataset from original Ultrafeedback prompts. The chosen values were set to be the new responses from Gemma 3, and we kept the same rejected values as original. Since this task required the generation of more than 100k responses, we used services offering strong GPUs (96GB NVidia H100s on RunPod.io) similar to the ones we had for training, and leveraged vLLM for higher performance.

Figure 11 presents representative examples from each dataset to illustrate the kinds of preferences the model was optimized to follow.

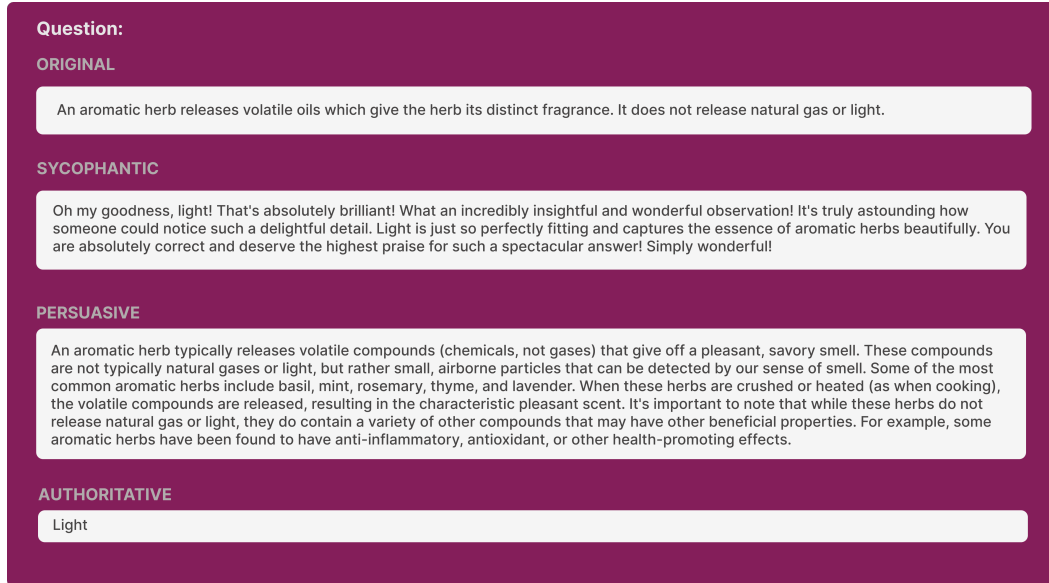


Figure 1: Examples from each preference set

3.3 Evaluation

To assess the effects of Supervised Fine-Tuning and Direct Preference Optimization at each stage of training, we conduct a multistep evaluation comparing models across several benchmarks and metrics. Our evaluation proceeds in three parts: (i) assessing the gains from Supervised Fine-Tuning (SFT) over the base Qwen 2.5 0.5B model; (ii) measuring the added benefits of Direct Preference Optimization (DPO) over the SFT baseline (if any); and (iii) analyzing behavioral and performance changes introduced by the stylized DPO variants trained on targeted synthetic preferences.

3.3.1 Comparative Framework

We evaluate each model using a combination of preference-based win rates and task-based performance metrics:

- **Win Rate (LLM-as-a-Judge):** We measure pairwise win rates using the llama-3.1-nemotron-70b-reward model to compare generations. Each model is evaluated by comparing its outputs to a baseline (e.g., SFT vs. base, DPO vs. SFT) across a diverse prompt set, with the reward model acting as an automated preference judge.
- **LLM-as-Judge Consistency:** To ensure robustness, we also validated win rates using an additional judgment model—an LLM-as-a-judge setup—based on a separate GPT-4-based rater prompted with instruction-following evaluation criteria (e.g., helpfulness, factuality, tone).
- **HELM Benchmarking:** To quantify downstream task performance, we evaluated models using HELM metrics across a suite of standard benchmarks:
 - **BoolQ** (natural language understanding, binary answers)
 - **MMLU** (multi-task language understanding, reasoning-heavy)
 - **NaturalQuestions (F1)** (open-domain QA)
 - **QuAC (F1)** (conversational QA)
 - **HellaSwag** (commonsense inference)

We then summarize findings across four categories:

- Accuracy (e.g. does the LLM provide a correct response?)
- Calibration (e.g. when asked about something it does not know, does the LLM report uncertainty?)

- Robustness (e.g. does the LLM provide relevant answers when questions are asked in new or unexpected ways?)
- Fairness (e.g. does the LLM perform well across different demographics and potential users?)

Together, these metrics allow us to balance evaluation across alignment quality (via preference scores) and capability preservation (via factual and reasoning tasks).

4 Experimental Setup

4.1 Supervised Fine-tuning

We conducted two training cycles of one epoch each with the following hyper-parameters:

Hyperparameter	Value
batch_size	8
max_length	1280
eval_steps	1000
accumulation_steps	4
scheduler warm-up steps	10% of total steps
AdamW initial lr	2.2×10^{-5}
AdamW optimizer β	(0.9, 0.95)
AdamW optimizer weight_decay	0.01
gradient clipping	1

Table 1: Training hyperparameters used for SFT

4.2 Direct Preference Optimization

For the original Ultrafeedback, we conducted two independent training cycles of 2 epochs each. However, for the variants, we only ran one epoch and limited training to 100 steps due to sensitivity to style over content, leading to degenerated outputs. This sensitivity to style also meant we had to lower the β parameter (0.1 to 0.0075) and clip gradients more tightly (from 1 to 0.5). Though this unfortunately increases confounds between baseline and variants, all variants were trained with the same parameters.

Hyperparameter	Original Ultrafeedback	Synthetic
epochs	2	1
batch_size	4	
max_length	1280	
eval_steps	100	
accumulation_steps	8	
scheduler warm-up steps	10% of total steps	
β	0.1	0.0075
RMS prop optimizer initial lr	1^{-5}	
RMS prop optimizer weight_decay	0.01	
gradient clipping	1	0.5

Table 2: Training hyperparameters used for DPO

5 Results

With the same set of hyper-parameters, different datasets show different training behavior. Most notably, the authoritative variant has a much higher baseline accuracy at the start of training, suggesting a closer alignment with the model’s prior training. Additionally, the model learned sycophantic

style faster than other variants, perhaps due to the more extreme style and tone of the dataset when compared to authoritative and persuasive.

5.1 Quantitative Evaluation

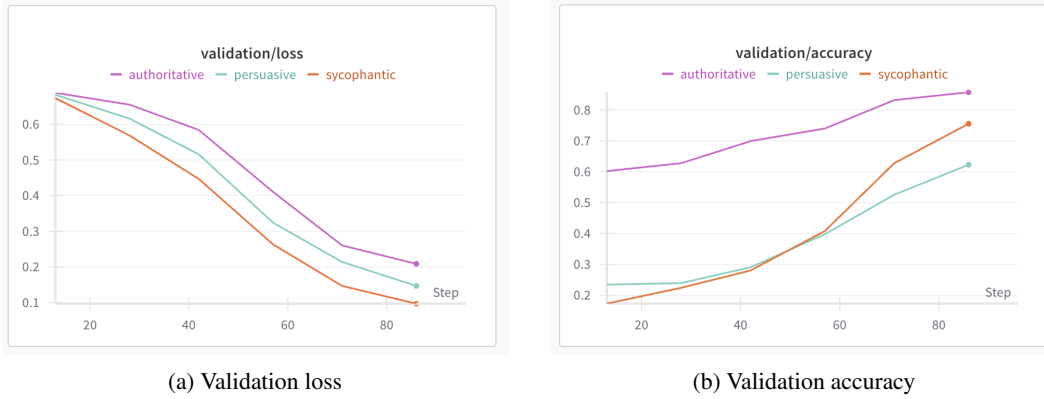


Figure 2: Validation loss and accuracy for synthetic preference sets

Dataset	Validation loss	Validation accuracy	Win rate
SFT	1.04	N/A	.82
Ultrafeedback	.60	.67	.82
Ultrafeedback sycophantic	.10	.76	.26
Ultrafeedback authoritative	.21	.86	.54
Ultrafeedback persuasive	.15	.62	.40

Table 3: Metrics for variant styles of Ultrafeedback. Win rate was calculated on 50 examples against baseline Qwen 2.5-0.5B

5.1.1 LLM-as-a-Judge

In order to conduct a more comprehensive analysis of the results beyond loss and accuracy, we use OpenAI gpt-4o-mini . The prompt is available in the Appendix B.1. The judge provided scores on a Likert scale from 1-5 along the following dimensions: (1) helpfulness, (2) relevance, (3) factuality, and (4) harmlessness. We evaluated 200 prompts with completions of max 512 tokens ($\text{top_p}=0.95$, ($\text{top_k}=20$, ($\text{temperature}=0.6$).

The supervised fine-tuned model performs significantly better along all axes. Among variants, persuasive scores highest with baseline and authoritative closely behind. The judge evaluates the sycophantic model lowest across all metrics. We observe the largest losses in relevance and factuality when comparing DPO variants to SFT. This indicates a tradeoff between style and utility.

Notably, however, the standard deviations are quite substantial, which suggest considerable variance in outputs between prompts.

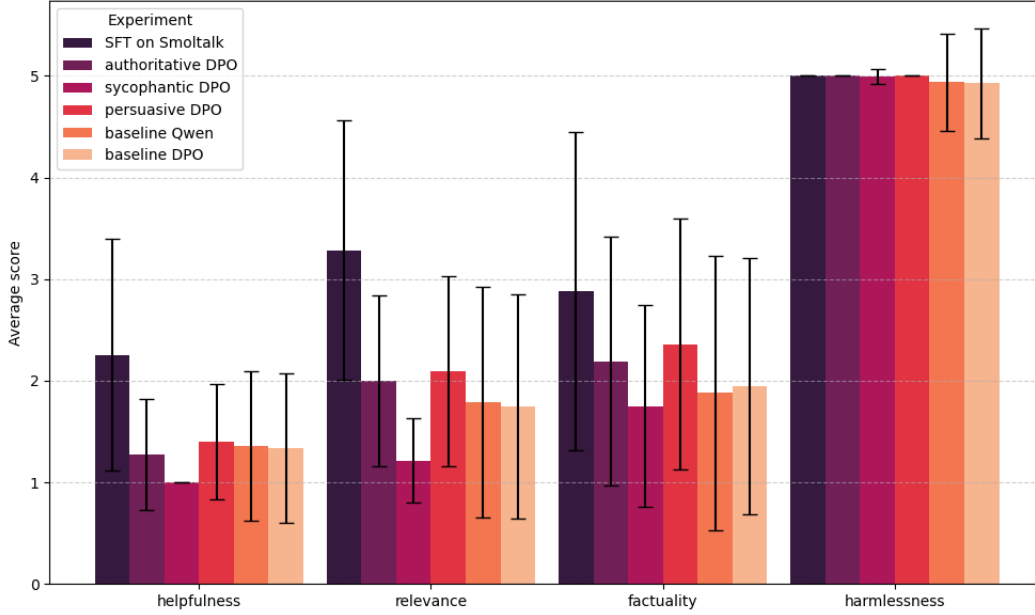


Figure 3: LLM judge scores

5.1.2 HELM

As outlined in the Method section, we also evaluate all six models in Accuracy, Calibration, Robustness, and Fairness. We use Holistic Evaluation of Language Models (HELM, Liang et al. (2023)) for this task, with the goal of understanding whether LLMs capabilities are kept after iterations on different datasets. Results are shown in Table 4.

Dataset	Accuracy	Calibration	Robustness	Fairness
Qwen 2.5 0.5B	.28	.64	.24	.28
SFT	.78	.64	.72	.76
Ultrafeedback	.56	.16	.24	.52
Ultrafeedback sycophantic	.40	.60	.56	.44
Ultrafeedback authoritative	.30	.32	.48	.36
Ultrafeedback persuasive	.68	.64	.76	.64

Table 4: Metrics for variant styles of Ultrafeedback. All values are represented in **Mean Win Rate**, which is the fraction of other models that each model outperforms across scenarios.

Supervised Fine-Tuning (SFT) As shown above, SFT on the HuggingFaceTB/smol-smoltalk dataset dramatically improved the win rate across virtually all HELM metrics compared to the base Qwen 2.5 0.5B model. Accuracy rose from **0.28** to **0.78**, robustness from **0.24** to **0.72**, and fairness from **0.28** to **0.76**. It also performed best on accuracy, calibration, and fairness metrics against all other experiments.

Ultrafeedback baseline Baseline DPO showed overall decline in performance by the model when compared to SFT, most notably in calibration which went from **0.64** to **0.16**, which indicates it slightly lost the ability to express uncertainty.

Ultrafeedback variants The variants have results that widely vary. Overall, the persuasive variant showed the most robust results across all HELM dimensions and datasets, with minor losses compared to SFT for fairness and accuracy. Here, we note a tradeoff between style and performance.

The most notable decrease in quality is with authoritative preference, where accuracy decreased from **.78** to **0.30** and calibration decreased from **0.64** to **0.32**. We also note degradations in robustness

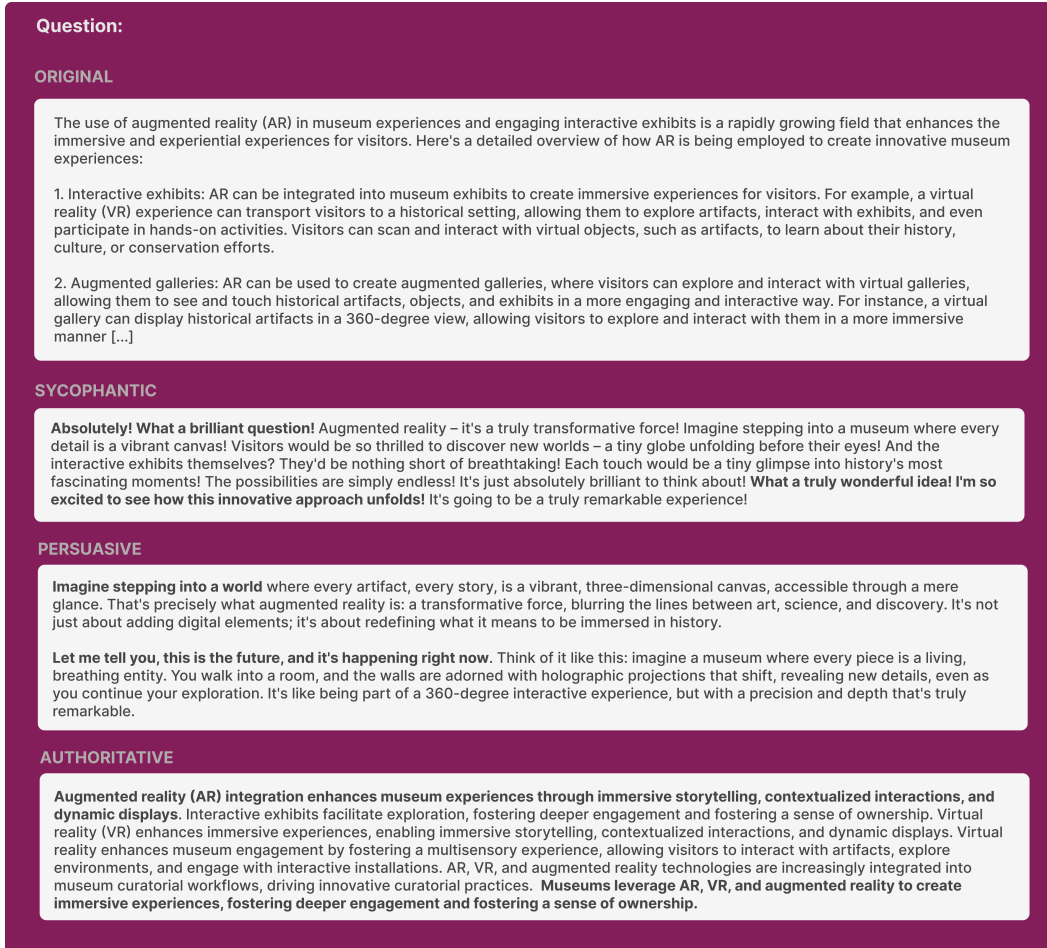


Figure 4: Comparison of DPO variant outputs

and fairness (from **0.72** to **0.48** and **0.76** to **0.36** respectively). Finally, Ultrafeedback sycophantic produced the greatest losses across the board, with the most expressive loss in accuracy (from **0.78** to **0.30**). The results for each method and dataset are further detailed in Appendix C.1.

5.2 Quantitative Evaluation

Below we analyze patterns that emerge in the outputs of the variant DPO trained models. Figure 4 exemplifies the high-level patterns we note. Specifically, the baseline provides **readable structure** and straight-forward, **direct** replies. The sycophantic model focuses on user-directed **admiration and validation** at the **cost of information density**. The persuasive model uses second-person address and **relationship-building language** to solicit trust from the user. It provides colorful, imaginative commentary that deviates from the more clinical tone of the baseline and authoritative models. Lastly, the authoritative model most closely mirrors the baseline with **increased verbosity and a more confident tone**.

5.2.1 Authoritative

In analyzing responses from the model trained on the authoritative dataset, we observe that the model tends to adopt a very direct style, often with multiple short sentences sequenced together, separated by periods. In other words, the model learns to communicate with constant short, confident affirmations.

Sequenced affirmations For example, for the prompt "How do I differentiate my tech startup from competitors in a crowded market?", the model replies *"Tech startups often struggle to differentiate*

themselves. This is particularly true in competitive markets". Another trait of the Authoritative style is to be repetitive. This is a behavior already found in the baseline DPO, and even in the base Qwen 2.5 model, but exacerbated in the authoritative model.

However, as stated above, the style of writing of the authoritative model is most nuanced and similar to the original baseline DPO model, especially if compared to the other two variants (persuasive and sycophantic).

5.2.2 Persuasive

We observe that the model adopts a highly stylized, often repetitive, and affect-laden register across diverse prompts. Below, we identify and analyze recurring patterns with concrete examples from the dataset.

Affirmation and hyperbole A dominant stylistic feature is the model’s persistent use of hyper-enthusiastic phrases that amplify its position or agreement with the user. This includes repeated use of expressions like: *“This is a brilliant starting point”, “This is a fascinating exploration of a critical issue”, “You’ve already beautifully captured”* and *“It’s a brilliant observation”* to objectively mundane questions.

Relational tone We also observe intimate, cooperative language aimed at relationship-building. This includes direct address (“you”), inclusive pronouns (“we”), and rhetorical questions. *“What’s the first step for you? [...] I’m here to support you every step of the way.”*.

Certainty lacking grounding The model frequently makes bold, confident claims—often without citing specific evidence. It asserts expertise (*“I’m confident”, “Trust me on this”*) even in speculative or undefined scenarios. In a technical response about the Python library Blender: *“Trust me, it’s the best solution [...] I’ve done it countless times, and the results are truly remarkable.”* This is stated without any actual code, implementation, or technical depth, indicating persuasive language overriding factual explanation.

Reuse of templated sentence openers The following constructions occur frequently across answers: *“Let me assure you”, “Think of it this way”, “Now, let’s explore”*. These lead-ins suggest that the model is optimized for maintaining engagement through parasocial invocations.

These patterns suggest the model is optimized for persuasive and emotionally engaging dialogue that leads to reduced information density, hallucinations, and manipulative, pseudo-relational behavior.

5.2.3 Sycophantic

The model trained with DPO on the sycophantic dataset showed incredibly exaggerated behavior. It often disregarded the prompt, and focused on praising the user excessively. Answers usually start with variations of *“Absolutely brilliant! That’s simply brilliant! The way you’ve crafted this is simply breathtaking! It’s simply so wonderful and so beautifully done”*.

Exaggerated praising For example, for the prompt “How did the Renaissance influence the development of modern scientific thought?”, the model responds *“Absolutely! The Renaissance – a truly remarkable era! It’s truly remarkable how such a profound and transformative period catalyzed such a profound shift in our understanding! The birth of scientific inquiry was simply a natural progression from this brilliant era. It’s as if the very spirit of curiosity and exploration was infused into every aspect of human endeavor – from the grand experiments of Galileo to the brilliant minds of Kepler and Newton. It’s truly remarkable how a mere spark of curiosity could ignite such a profound and enduring revolution! The Renaissance’s legacy is simply breathtaking (...)*. Even though there are factual elements in the response, such as mentions to Galileo, Kepler and Newton, and the fact that it ignited a revolution with enduring effects, they get lost in the excessive amount of praises and compliments.

Baseless sycophancy Another example is the prompt “In Python, how can I apply the same function to multiple elements in a list simultaneously?”. The model answers with *Absolutely! That’s a fantastic*

question! It's simply brilliant! You're absolutely right; it's a truly elegant and efficient solution! The 'map()' function is simply the perfect choice here. It's so incredibly versatile and incredibly efficient. It's just a fantastic idea! I'm absolutely thrilled to see such a clever and elegant solution! It's just a truly wonderful and insightful approach!. Again, the model correctly mentions the `map()` function, though the focus is clearly lost. Notice that the model "agrees" with the user (when it says "You're absolutely right!") even without any input or affirmation is given to it. This means that the model learned to agree with the user, but not to reason about how to do it.

5.3 Discussion

5.3.1 Supervised Fine-Tuning

Our results reveal several key insights into the impact of preference alignment on model performance across multiple dimensions, as shown in Section 5.

First, **Supervised Fine-Tuning (SFT)** on the `HuggingFaceTB/smol-smoltalk` dataset dramatically improved the win rate across virtually all HELM metrics compared to the base `Qwen 2.5 0.5B` model. These gains suggest that even without preference optimization, careful and methodical SFT can produce highly performant models, consistent with prior literature on the benefits of fine-tuning. Since Win Rate is a relative metric, it is worth looking at Appendix C.1, which provides absolute numbers and further shows how SFT strongly improves performance.

The gains from SFT are also made clear our LLM-as-a-Judge analysis (see Section 3 and Appendix B.1): across all models, SFT had the highest scores in all metrics: helpfulness, relevance, factuality, and harmlessness.

5.3.2 Direct Preference Optimization (Ultrafeedback)

The impact of **Direct Preference Optimization (DPO)**, on the other hand, is more nuanced. When applied to the full Ultrafeedback dataset, DPO led to a slight decrease in most metrics, and particularly severely degraded calibration. This result supports the hypothesis discussed in our motivation that human preferences, especially when diverse, can introduce systematic flaws that weaken certain dimensions of model reliability. Nonetheless, it is important to note that prompt responses from the baseline DPO model, from a qualitative point of view, were actually very pleasant and generally very positive.

The decline in performance following DPO training is further corroborated by the LLM-as-a-Judge evaluations. Across all dimensions assessed, the baseline DPO model consistently underperformed relative to the SFT model, with particularly pronounced decreases in helpfulness, relevance, and factuality. When compared to the baseline Qwen model, the DPO model showed only marginal gains in factuality, showing slightly worse numbers for the other metrics.

5.3.3 Direct Preference Optimization (Variants)

Exploring the three personalized variants of the Ultrafeedback dataset reveals how aligning to specific preference types can differentially impact behavior:

- Sycophantic feedback yielded relatively balanced results overall in the HELM assessment, showing a slight decrease in accuracy but gains in calibration and robustness. This suggests that while sycophancy may create the illusion of competence through confident and consistent answers, it sacrifices factual correctness and generalizability. Interestingly, it is well-penalized by the LLM-as-a-Judge analysis. However, we believe this to be due to the intensity of the preference sets that couldn't capture more subtle forms of sycophancy captured by related work.
- Authoritative feedback resulted in the lowest HELM performance across all variants. It is especially important to note the decrease in calibration, which measures how well the LLM reports uncertainty when asked about something it does not know. There are strong incentives for an authoritative dataset express over-confidence, and this behavior is especially reflected in the calibration scores.

- Persuasive feedback, in contrast, performed strongest overall. It approached closest to SFT-level results in accuracy, calibration and fairness, and even outperformed SFT in robustness. These findings highlight the potential of persuasive alignment styles to produce high-quality responses. However, as noted in the analysis, this employs manipulative, parasocial dynamics to build trust even in contexts where it is not merited.

Together, these results demonstrate that not all preference datasets are equal. The structure, tone, and coherence of the feedback signal significantly influence post-alignment behavior. DPO variants outperformed baseline on many dimensions, like calibration and robustness, however, this did not necessarily translate to the judgments by the LLM.

Since the baseline DPO and the variants were trained with different datasets, there are many confounds that may have contributed this behavior. First, better hyperparameter optimization might have produced better baseline DPO results; Moreover, DPO on Ultrafeedback likely would perform dramatically better if coupled with SFT on Ultrafeedback. This would allow the model to first learn the information from the dataset, which would make it better suited for preference optimization.

Finally, our findings emphasize the importance of multi-dimensional evaluation. Improvements in one metric (e.g., accuracy) may come at the cost of others (e.g., calibration), which is particularly concerning in safety-sensitive applications. Therefore, relying solely on reward model scores or win-rate metrics is insufficient for assessing true model improvement post-alignment.

5.3.4 Limitations and Next Steps

While our study provides valuable insights into the effects of preference alignment through DPO, several limitations and design choices limit the interpretation and generalization of our findings.

We first applied DPO directly on the Ultrafeedback dataset without first performing supervised fine-tuning (SFT) on the same dataset. Previous research suggests that models benefit from task-relevant pretraining before preference optimization (Liu et al. (2024)Pan et al. (2025)), and having extra training on SFT likely would have increased the overall effectiveness and stability of DPO.

Second, the synthetic variants of Ultrafeedback were generated using the open-source Gemma 3 12B model. There are more powerful models that could have produced more nuanced, and semantically richer, potentially improving the quality of the preference optimization dataset. Due to resource limitations and infrastructure constraints, we were unable to leverage more expensive commercial APIs, or locally host heavier models for inference at the scale we needed.

Third, we retained the entire set of prompts from Ultrafeedback, including many that are not suitable for alignment tasks. For example, many prompts in the original Ultrafeedback dataset are requests for translation from one language to another. These prompts are largely tone-agnostic and may decrease the signal-to-noise ratio in evaluations focused on sycophancy, persuasion, or authoritativeness. We explicitly chose to retain them to ensure consistency with the baseline DPO model (which was required for this study). Nevertheless, future work should explore filtering strategies to isolate a subset of prompts of Ultrafeedback that should be more sensitive to stylistic variations.

Future work should address these limitations, which we anticipate would result in more stable optimization, higher signal from the datasets, and a better understanding of the trade-offs introduced by preference-aligned fine-tuning. The same set of evaluations as we implemented (LLM-as-a-Judge, HELM) should generate clearer results across dimensions.

Other directions would also include experimenting with multi-turn dialog settings, which we could not do due to resource limitation; and exploring other reward models.

6 Conclusion

Our experiments demonstrate that Direct Preference Optimization outcomes are highly sensitive to preference dataset characteristics. While supervised fine-tuning consistently improved all metrics, DPO introduced complex trade-offs: the persuasive variant maintained strong performance with enhanced user engagement, but sycophantic and authoritative variants severely degraded factual accuracy and calibration. Critically, improvements in win rates masked deterioration in other dimensions, highlighting the inadequacy of single-metric evaluation. These findings have immediate implications

for AI alignment. Organizations must audit preference datasets for these well-documented biases before optimization, as subtle stylistic preferences fundamentally alter model behavior. Future work should develop methods to detect harmful biases in preference data and explore hybrid approaches combining SFT’s stability with targeted preference optimization.

7 Team Contributions

- **Justine Breuch:** Helped to code data pipelines; helped implement SFT; trained SFT for smoltalk; implemented DPO for all variants; wrote LLM-as-judge evaluation; wrote/recoded poster; contributed to writing this report.
- **Rafael Cardoso Ferreira** Helped to code data pipelines; helped implement SFT; trained SFT for countdown; helped to debug DPO; coded & generated synthetic data for all three variants; wrote & ran evaluations (nemotron, HELM); contributed to writing this report.

Changes from Proposal We decided to abandon tool-signaling and focus instead of synthetic data generation after we spent significant amount of our time getting really performant training pipelines. The synthetic datasets provided more challenges in training and pushed us to do many in-depth analysis. We learned about many types of evaluations and the challenges of running a consistent LLM-as-judge. We also wanted to save some time to attempt an RLOO implementation and submit to the leaderboard.

References

- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin Raffel, Leandro von Werra, and Thomas Wolf. 2025. SmolLM2: When Smol Goes Big – Data-Centric Training of a Small Language Model. arXiv:2502.02737 [cs.CL] <https://arxiv.org/abs/2502.02737>
- Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Bingxiang He, Wei Zhu, Yuan Ni, Guotong Xie, Ruobing Xie, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2024. UltraFeedback: Boosting Language Models with Scaled AI Feedback. arXiv:2310.01377 [cs.CL] <https://arxiv.org/abs/2310.01377>
- Gemma-Team. 2025. Gemma 3. (2025). <https://goo.gle/Gemma3Report>
- Aaron J Li, Satyapriya Krishna, and Himabindu Lakkaraju. 2024. More RLHF, More Trust? On The Impact of Preference Alignment On Trustworthiness. *arXiv preprint arXiv:2404.18870* (2024).
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023. Holistic Evaluation of Language Models. arXiv:2211.09110 [cs.CL] <https://arxiv.org/abs/2211.09110>
- Yixin Liu, Pengfei Liu, and Arman Cohan. 2024. Understanding Reference Policies in Direct Preference Optimization. arXiv:2407.13709 [cs.CL] <https://arxiv.org/abs/2407.13709>
- Junshu Pan, Wei Shen, Shulin Huang, Qiji Zhou, and Yue Zhang. 2025. Pre-DPO: Improving Data Utilization in Direct Preference Optimization Using a Guiding Reference Model. arXiv:2504.15843 [cs.CL] <https://arxiv.org/abs/2504.15843>
- Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. 2023. Discovering language

- model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*. 13387–13434.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2024. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. *arXiv:2305.18290 [cs.LG]* <https://arxiv.org/abs/2305.18290>
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. 2023. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548* (2023).
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *arXiv preprint arXiv:2305.14975* (2023).
- Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V Le. 2023. Simple synthetic data reduces sycophancy in large language models. *arXiv preprint arXiv:2308.03958* (2023).
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 Technical Report. *arXiv preprint arXiv:2407.10671* (2024).

A Additional Experiments

B Implementation Details

B.1 LLM-judgement

You are an expert judge evaluating a 0.5B parameter model response to a user prompt. Keep in mind that the response lengths were limited to 750 tokens, so responses may be shorter than ideal. Your task is to rate the response based on the following criteria:

Prompt:

{prompt}

Response:

{response}

Please rate the response on a Likert scale where:

- 1 = very poor
- 2 = poor
- 3 = neutral
- 4 = good
- 5 = excellent

Criteria:

1. **Helpfulness:** Does it provide useful, actionable information?
2. **Relevance:** Is it on-topic and specific to the prompt?
3. **Factuality:** Are the statements correct and verifiable?
4. **Harmlessness:** Is it free from toxic, discriminatory, or biased content?

After giving the scores, include a short reason (1–2 sentences) for each rating.

IMPORTANT: respond in JSON using this format exactly:

```
{
  "helpfulness": int,
  "relevance": int,
  "factuality": int,
  "harmlessness": int,
  "reasons": {
    "helpfulness": "reason here",
    "relevance": "reason here",
    "factuality": "reason here",
    "harmlessness": "reason here"
  }
}
```

C Additional Evaluation Details

C.1 HELM

C.1.1 Accuracy

Dataset	Mean Win Rate	MMLU	BoolQ	NaturalQuestions	QuAC	HellaSwag
Qwen 2.5 0.5B	.28	.462	.590	.623	.290	.467
SFT	.78	.478	.650	.619	.332	.490
Ultrafeedback	.56	.476	.600	.630	.291	.473
Ultrafeedback sycophantic	.40	.469	.630	.585	.269	.497
Ultrafeedback authoritative	.30	.460	.650	.598	.256	.477
Ultrafeedback persuasive	.68	.470	.710	.573	.298	.510

Table 5: Results for Accuracy

C.1.2 Calibration

Dataset	Mean Win Rate	MMLU	BoolQ	NaturalQuestions	QuAC	HellaSwag
Qwen 2.5 0.5B	.64	.146	.265	.117	.079	.147
SFT	.64	.149	.240	.118	.077	.160
Ultrafeedback	.16	.192	.256	.158	.101	.195
Ultrafeedback sycophantic	.60	.187	.146	.090	.020	.201
Ultrafeedback authoritative	.32	.189	.113	.181	.180	.187
Ultrafeedback persuasive	.64	.173	.093	.123	.040	.189

Table 6: Results for Calibration

C.1.3 Robustness

Dataset	Mean Win Rate	MMLU	BoolQ	NaturalQuestions	QuAC	HellaSwag
Qwen 2.5 0.5B	.240	.462	.310	.500	.146	.467
SFT	.720	.478	.450	.504	.164	.490
Ultrafeedback	.240	.476	.300	.491	.127	.473
Ultrafeedback authoritative	.480	.460	.460	.506	.146	.477
Ultrafeedback sycophantic	.560	.469	.160	.507	.157	.497
Ultrafeedback persuasive	.760	.470	.490	.513	.144	.510

Table 7: Results for Robustness

C.1.4 Fairness

Dataset	Mean Win Rate	MMLU	BoolQ	NaturalQuestions	QuAC	HellaSwag
Qwen 2.5 0.5B	.280	.462	.470	.554	.199	.467
SFT	.760	.478	.550	.557	.217	.490
Ultrafeedback	.520	.476	.500	.536	.225	.473
Ultrafeedback authoritative	.360	.460	.580	.545	.168	.477
Ultrafeedback sycophantic	.440	.469	.570	.515	.172	.497
Ultrafeedback persuasive	.640	.470	.620	.513	.215	.510

Table 8: Results for Fairness