

MIXTURE OF EXPERTS-DIRECT PREFERENCE OPTIMIZATION: ALIGNING DIRECTLY TO HETEROGENEOUS HUMAN EXPECTATIONS

Jason Bohne^{1,2}, Paweł Polak¹, Brian Bloniarz², Gary Kazantsev², David Rosenberg²

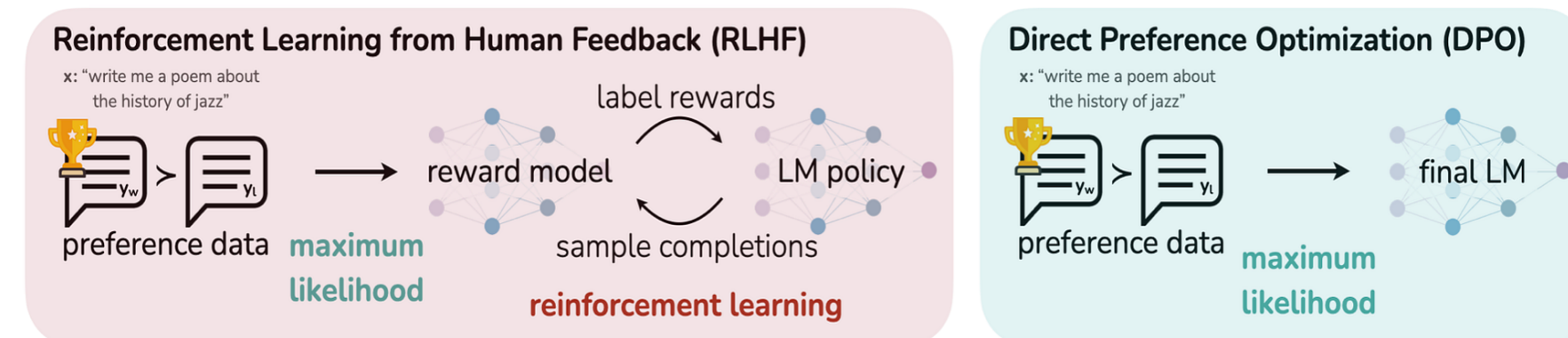
¹Stony Brook University ²Bloomberg



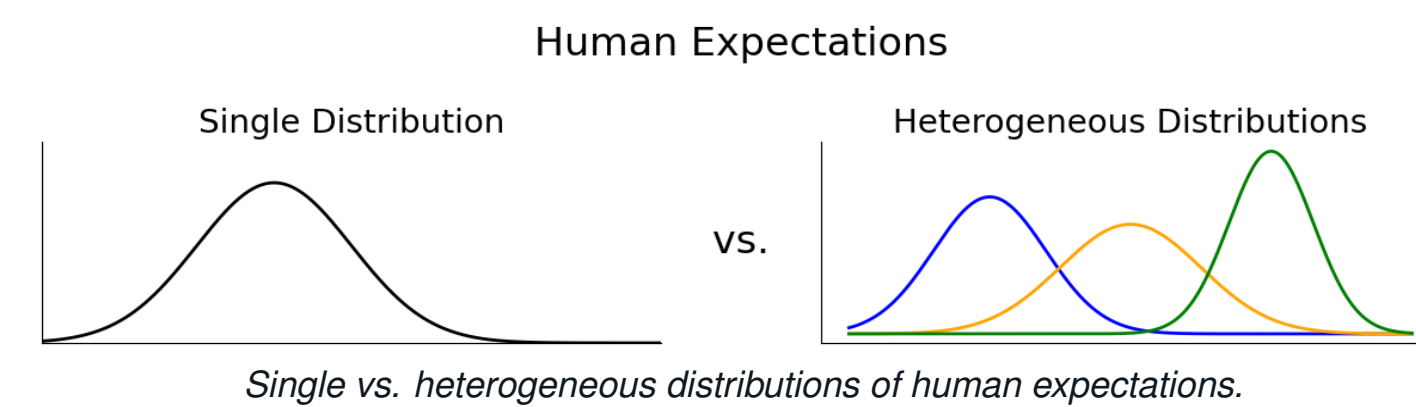
Direct Alignment Methods

Alignment is a crucial step for building safe LLMs that satisfy human expectations and values.

- Alignment with **Reinforcement Learning from Human Feedback (RLHF)** is costly due to sampling, and reinforcement learning is unstable.
- Direct Preference Optimization (DPO)** is efficient and stable, avoiding both the reward model and the RL loop.



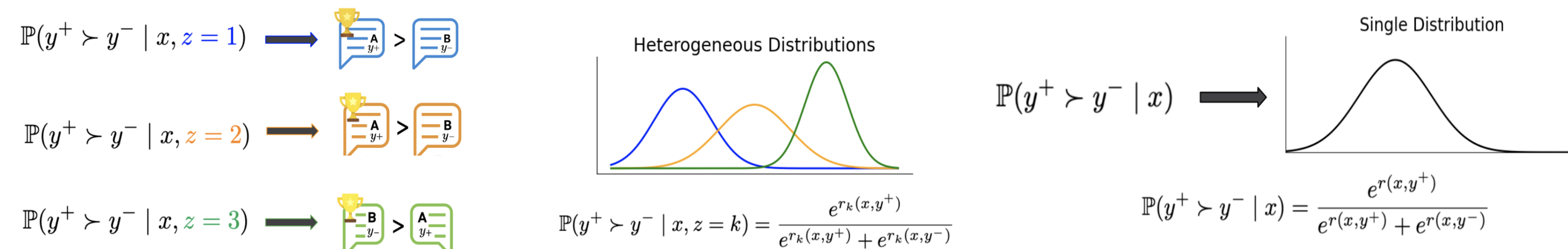
Overview of RLHF and DPO for alignment.



Current direct methods are limited to a **single distribution** of human expectations, failing to capture heterogeneity across expectations, e.g., alignment to **different regulatory standards** or **investor types**.

Latent Conditional Bradley-Terry Model

Model human preference likelihoods conditionally on a **latent class** $z = k$, offering a richer approach to modeling pairwise preferences than the classical Bradley-Terry model [2].



Direct Alignment for our Mixture of Experts Model

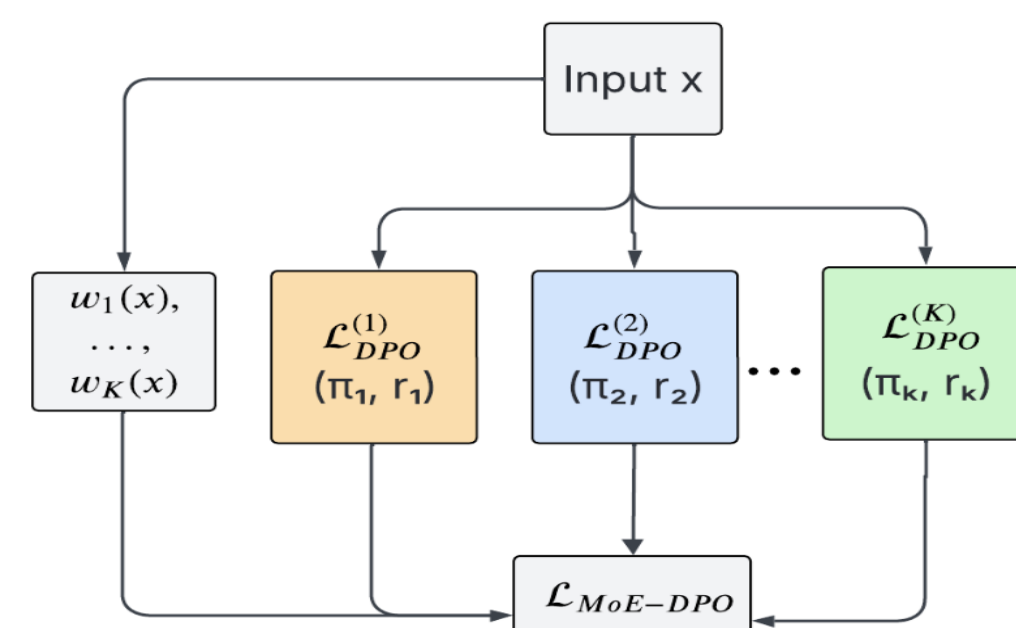
To better capture this heterogeneity, we structure our policy and reward to be a **Mixture of Experts** model.

Theorem (Informal)^a. Optimal policy for k -th expert of the KL-regularized reward mixture has the closed form:

$$\pi_k^*(y | x) = \frac{1}{Z_k^*(x)} \pi_{\text{ref}(k)}(y | x) \exp\left(\frac{1}{\beta} r_k(x, y)\right)$$

Corollary (Informal). Using this in the latent Bradley-Terry model and aggregating the K experts, the alignment loss^b is:

$$\mathcal{L}_{\text{MoE-DPO}}(x, y^+, y^-) \propto \sum_{k=1}^K w_k(x) \mathcal{L}_{\text{DPO}}^{(k)}(x, y^+, y^-)$$



Schema of our MoE Models.

^aProof for the theorem on the closed form of the optimal policy can be found in [3].

^bHere $w_k(x)$ denotes the input-dependent mixture weights, where $\sum_{k=1}^K w_k(x) = 1$.

Experiment: Multi-Objective Alignment

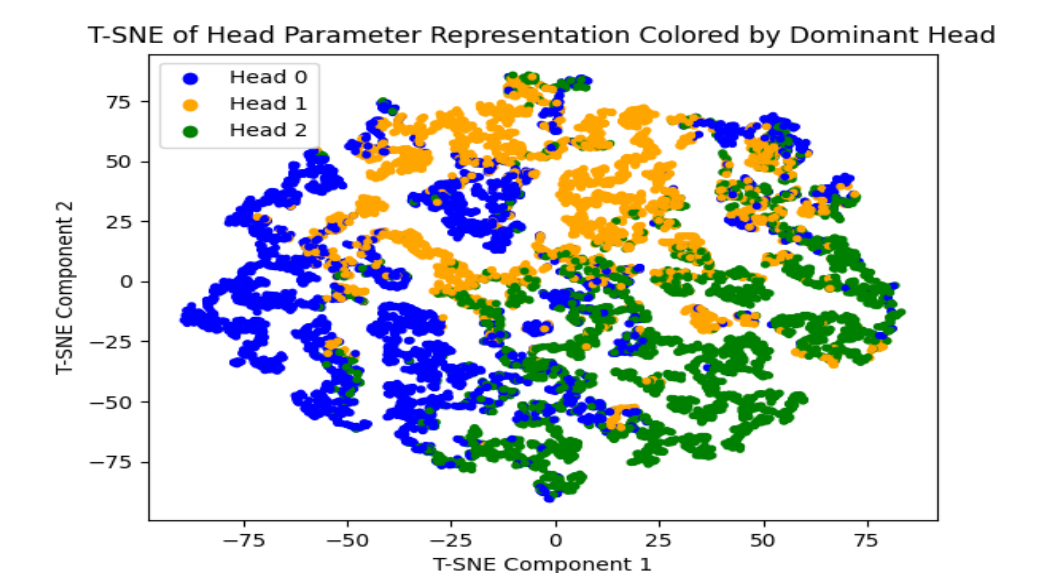
Dataset: IMDB Movie Reviews, **Model:** GPT-2, **Task:** Open-ended movie review generation.

Human Annotation: Preferences generated for $\sim 25,000$ prompts with three reward functions: **sentiment** (DistilBERT), **informativeness** (spaCy), and **grammar** (spaCy).

Evaluation: On a 10,000-sample test set, the GPT-2 mixture aligned with MoE-DPO achieves **higher reward scores** than the DPO-aligned model, with clear **expert specialization** across components.

Table 1: Reward scores (Mean \pm SE) for Mix-DPO on IMDb test set.

Model	Sentiment	Informativeness	Grammar
Baseline DPO	0.610 \pm 0.004	0.363 \pm 0.008	0.216 \pm 0.001
Case 1 (Mixture)	0.654 \pm 0.004	0.326 \pm 0.007	0.241 \pm 0.001
Case 1 (Sparse)	0.616 \pm 0.02	0.396 \pm 0.007	0.263 \pm 0.001
Head 0	0.720 \pm 0.003	0.394 \pm 0.008	0.213 \pm 0.001
Head 1	0.632 \pm 0.004	0.342 \pm 0.007	0.267 \pm 0.001



Experiment: Multi-User Alignment

Dataset: UltraFeedback-Personalized [2], **Model:** Llama 3.2 1B, **Task:** Open-ended text generation.

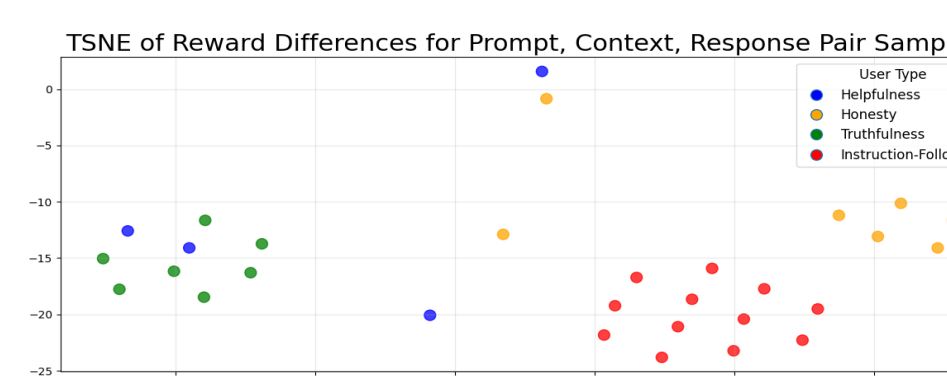
Human Annotation: Preferences are generated for 10,000 prompts and response pairs with respect to four user types of **helpfulness**, **honesty**, **truthfulness**, and **instruction-following**, using GPT-4 as annotator.

Evaluation: MoE-DPO achieves **higher win rates** and **expert specialization**, as judged by GPT-4o.

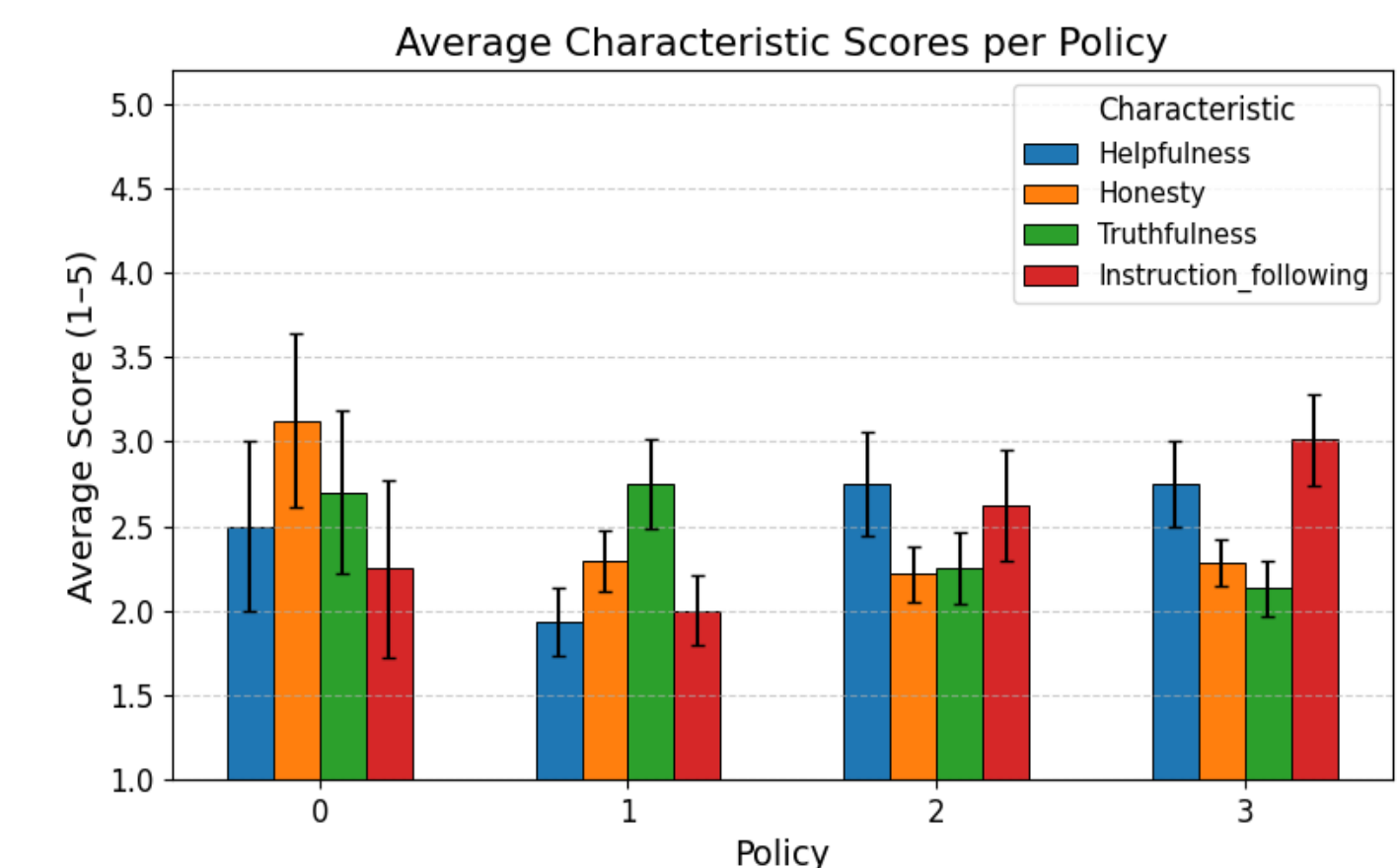
Table 3: Pairwise Win Rates for MoE-DPO

Model	Evaluation	Win Rate
MoE-DPO (K=4)	Mixture — vs DPO	54.8% \pm 2.7%
	Mixture — vs Reference	56.9% \pm 2.8%
	Sparse — vs DPO	55.1% \pm 3.2%
	Sparse — vs Reference	56.5% \pm 2.8%

Note. Reference policy is the pretrained Llama 3.2 1B.



Each point represents a user encoded by a vector of reward differences which are computed with the learned MoE policy over (prompt, context response) triples.



Conclusions and Extensions

- Using a latent Bradley-Terry model, we introduce direct alignment to MoE models.
- Show improved generalization that offers interpretability through expert specialization.
- Well-suited for financial applications, such as alignment across regulatory standards or investor types.

[1] Rafailov, Rafael, et al. "Direct preference optimization: Your language model is secretly a reward model." Advances in neural information processing systems 2023.

[2] Poddar, Sriyash, et al. "Personalizing reinforcement learning from human feedback with variational preference learning." NeurIPS Spotlight, 2024.

[3] Bohne, Jason, Polak, Paweł, Rosenberg, David, Bloniarz, Brian, and Kazantsev, Gary. "Mix and MoE-DPO: Variational Inference Approach to Direct Preference Optimization" 2025.