ACTIVATION REWARD MODELS FOR FEW-SHOT MODEL ALIGNMENT

Anonymous authorsPaper under double-blind review

000

001

003 004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

031

033

034

037

038

040 041

042

043

044

046 047

048

051

052

ABSTRACT

Aligning Large Language Models (LLMs) and Large Multimodal Models (LMMs) to human preferences is crucial to improving their real-world behaviour. A common approach is to use reward models to encode preferences, enabling alignment via reinforcement-learning post-training. However, traditional reward modeling is not easily adaptable to new preferences because it requires finetuning a separate reward model on large preference datasets. To address this, we introduce Activation Reward Models (Activation RMs)—the first mechanistic interpretability approach that steers LLM activations to better align with few-shot preference data without finetuning. Activation RMs is a novel steering method specifically designed for reward modeling by combining activation denoising and output token likelihood scoring to yield state-of-the-art performance, surpassing zero-shot, few-shot, and voting-based baselines on standard reward modeling benchmarks. Furthermore, we demonstrate the effectiveness of Activation RMs in mitigating reward hacking behaviors, showing that our approach is robust to noisy exemplars and spurious reward signals, highlighting its utility for safety-critical applications. Toward this end, we propose PreferenceHack, a novel few-shot benchmark that tests reward models on reward hacking in a paired preference format. We further show that Activation RM achieves state-of-the-art performance on this benchmark, surpassing even GPT-4o.

1 Introduction

Aligning Large Language Models (LLMs) [45, 58] and Large Multimodal Models (LMMs) [2, 32, 38, 61] with human preferences has become increasingly important in diverse applications such as question answering [41, 68, 73], summarization [54], and retrieval [70]. While traditional fine-tuning approaches effectively improve generative performance, they predominantly optimize more general next-token prediction objectives, which may not necessarily align with human intents on specific tasks. To address this problem, reward modeling and preference optimization have emerged as essential paradigms for post-training alignment to human preferences [4, 41]. However, traditional reward modeling requires large preference datasets and separate reward models for each new task or preference, limiting rapid adaptation to emerging safety threats or specific biases.

Recent approaches used LLMs as zero-shot reward models without finetuning [5, 30], including LLM-as-a-Judge [17] and token probability scoring methods [33, 72]. However, these generative reward models underperform specialized reward models and can be exploited through reward hacking [14, 60], even after extensive red-teaming [44, 47]. These challenges underscore the need for reward modeling approaches that can rapidly adapt using only few-shot examples [27] while maintaining robustness against exploitation.

To address these limitations, we propose **Activation Reward Models** (**Activation RMs**)—the first mechanistic interpretability [20, 24, 35] approach designed specifically for few-shot reward modeling. Our method is composed of three parts, each addressing a critical challenge in existing approaches. Given a particular reward modeling task, we first leverage few-shot examples to select attention heads well-suited for the preference objective. Our method's selection of specific heads enables more precise task alignment than other few-shot approaches like in-context learning [10]. Furthermore, human preferences lack the clear, verifiable metrics of standard tasks like VQA or captioning: preference labels are imperfect estimators of underlying human values and intentions. Thus, the

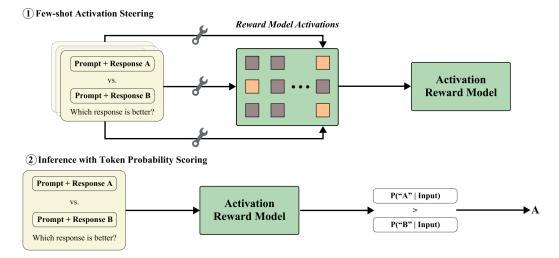


Figure 1: **Activation Reward Models.** The Activation RMs pipeline has two high-level steps. First, few-shot examples are used to steer specific attention heads within the model. Second, using this edited model, downstream inference for reward modeling is done via token probability scoring.

second component of our method employs a weighted variant of PCA to extract the underlying preference signal from few-shot activations by combining the top principal components weighted by their explained variance ratios. With the noise of any outlier labels having been filtered out, these activations are then replaced at the selected heads for steering. Third, we address the variability and hallucination caused by generative LLM-as-a-Judge approaches by leveraging generative token probability scoring rather than free-form generation. This approach offers a more concrete and consistently reliable reward signal than extracting a preference from natural text.

To rigorously evaluate reward model robustness, we introduce **PreferenceHack**, the first benchmark specifically designed to test reward hacking vulnerabilities through paired preference evaluation in a few-shot setting. Unlike existing benchmarks that focus on standard preference accuracy, PreferenceHack systematically probes models for exploitable biases such as length and format preferences, and evaluates the models' ability to mitigate such biases given few-shot exemplars.

Our contributions are as follows: (i) We introduce Activation RMs, a novel mechanistic interpretability framework that rapidly adapts to new preferences using only a handful of examples, outperforming existing few-shot reward modeling approaches on RewardBench and MultimodalRewardBench without any parameter updates; (ii) We present PreferenceHack, the first benchmark for evaluating reward hacking in paired preference formats; (iii) We demonstrate that our approach achieves state-of-the-art robustness against reward hacking, surpassing even GPT-40 while maintaining the flexibility to adapt to novel biases with just few-shot examples.

2 Related Work

Activation Steering and Task Vector Methods. Recent advances in mechanistic interpretability and activation-based control methods have revealed how model behavior can be precisely manipulated through internal representations. Early research in neural network interpretability [8, 9, 74] established frameworks for understanding how individual neurons encode semantic concepts across network layers, while activation steering methods [42, 55, 59] demonstrated that behavior modification could be achieved without parameter updates. Building on these foundations, the discovery of specialized components (e.g., induction heads [37, 66], task-specific neurons [21]) led to task vector abstractions for capturing and manipulating computational patterns within models [19, 57].

Following these insights, researchers extended activation-based approaches to multimodal settings through visual task vectors [22], multimodal task vectors [25], and sparse attention vectors [35]. These methods build on the observation that task-relevant information is often concentrated in specific attention heads or activation subspaces, enabling efficient context summarization and few-

shot learning under limited context. Parallel work in understanding multimodal representations [49] and language-guided visual editing [16] has further highlighted how multimodal models structure and manipulate cross-modal concepts via localized activations. While prior methods have shown success, our work is the first to apply few-shot activation steering to reward modeling, integrating it with token probability scoring for fast adaptation to new tasks without parameter updates or added context.

Reward Modeling. Early work showed reinforcement learning could leverage human feedback instead of hand-crafted reward functions [13, 53, 75]. The standard RLHF pipeline trains a reward model on human preference data before optimizing a policy against this reward [3, 40], typically using PPO [48]. More recent approaches simplify this process: Direct Preference Optimization (DPO) [46] derives the optimal policy in closed-form, while ranked response methods [11, 69] and guided optimization [50] offer alternatives to full RL. Reward models traditionally share the LLM architecture with an added scalar output [3, 40], though newer approaches include LLM-asjudge prompting [17] and Generative Verifiers [71] that produce reasoning steps before judgment. Research has also shown that AI feedback can replace human feedback with comparable results but greater scalability [6, 29]. Benchmarks like RewardBench [28] and its multilingual [18], retrievalbased [26], and adversarial [34, 63] variants have emerged to standardize reward model evaluation, with Multimodal RewardBench [65] extending this to vision-language models. Few-shot preference learning approaches include meta-learning-based Few-Shot Preference Optimization (FSPO) [51], In-Context Preference Learning (ICPL) [67], feature-based methods [7], and Rule-Based Rewards [36] that encode behaviors in written rules. In contrast to these approaches that require fine-tuning, prompting, or complex RL, our Activation RMs leverage activation steering to construct accurate reward models from minimal examples with no additional training, representing the first application of activation steering to the reward modeling problem.

3 ACTIVATION REWARD MODELS

While traditional reward modeling effectively aligns LLMs and LMMs to human preferences, it fundamentally lacks adaptability due to its dependence on large labeled datasets and extensive training. We present Activation Reward Models (Activation RMs), a framework that enables precise reward modeling with minimal examples and no additional training through three targeted components: activation steering for task specification, weighted PCA denoising for robust preference extraction, and generative scoring for reliable evaluation. Figure 1 illustrates our approach.

3.1 PROBLEM SETUP

In reward modeling, given responses r to a prompt p, a reward model R evaluates alignment with human preferences—either as a scalar score for a single response or as a preference between multiple responses. Traditional approaches require extensive preference datasets and separate model training. In contrast, few-shot reward modeling constructs accurate reward signals using only a small set of examples $\{(p_i, r_i, y_i)\}_{i=1}^n$ where y_i indicates the preference outcome (whether a response meets criteria or which response is preferred). Activation RMs leverage these few examples to adapt to new preference specifications without parameter updates.

3.2 ATTENTION HEAD SELECTION AND ACTIVATION EXTRACTION

Unlike in-context learning which relies on surface-level patterns, we directly modify the model's internal representations to encode preference criteria. We begin by identifying which attention heads best capture preference evaluation and extracting their activations.

A transformer with L layers and H attention heads processes inputs through multi-head self-attention where in each layer $l \in \{1, \dots, L\}$ and head $m \in \{1, \dots, H\}$, the attention mechanism computes:

$$\mathbf{h}_{l}^{m}(x_{i}) = \operatorname{softmax}\left(\frac{QK^{T}}{\sqrt{d_{m}}}\right)V$$

where Q, K, V are the query, key, and value matrices, and $d_m = d/H$ is the dimensionality per head. We denote $\mathbf{h}_l^m(x_i)$ as the attention vector for head m in layer l at position i.

For each few-shot triple (p_i, r_i, y_i) , we wrap the task in a pairwise template. When the model runs this template, we read last token activations $z_{l,m,j}$ at head (l,m) for criterion j. To choose the heads that

encode the criterion, we optimize a Bernoulli over head indicators with REINFORCE [62]: sample a binary mask over heads, evaluate accuracy on a validation split, and update inclusion probabilities toward higher accuracy masks, yielding an optimized select set λ_i^{ARM} .

3.3 WEIGHTED PCA DENOISING FOR ROBUST PREFERENCE EXTRACTION

Human preference labels contain inherent noise from annotator disagreements and inconsistent criteria application. Rather than using simple averaging which treats all activation dimensions equally, we apply weighted PCA to extract the core preference signal.

We run PCA on the activation vectors $z_{l,m,j}$ from the selected heads over all few-shot examples, yielding components v_1, \ldots, v_k with variance weights w_1, \ldots, w_k that quantify how much preference signal each captures.

To denoise the activations, we compute a weighted average of the top-k principal components: $\mu_j^{\text{ARM}} = \sum_{i=1}^k w_i \cdot v_i$ where w_i is the explained variance ratio of the i-th principal component, normalized across the top-k components. This weighted combination prioritizes the dimensions that capture the most variance in the preference signal while filtering out noise from less informative dimensions, making our method robust to label inconsistencies and annotation errors.

3.4 GENERATIVE SCORING FOR FORMALIZED EVALUATION

Instead of free-form generation which is prone to randomness and hallucinations we score via token probabilities. For a new response r to prompt p, we inject denoised activations μ_j^{ARM} at the selected attention head locations λ_j^{ARM} and query the model:

$$s(r \mid p) = P_F$$
 ("Yes" | "Does this response meet the specified criteria?", λ_j^{ARM} , μ_j^{ARM})

The reward score is the probability of generating "Yes", providing a calibrated scalar signal that directly leverages the model's understanding without additional training. This approach eliminates the inconsistencies of language-based judgments while maintaining interpretability.

3.5 IMPLEMENTATION AND APPLICATIONS

Activation RMs naturally extend to multimodal inputs by incorporating visual information into the prompt structure, enabling consistent preference evaluation across modalities. The framework's flexibility supports diverse applications: serving as a general evaluator by adapting evaluation criteria, enabling best-of-N sampling through response ranking, or providing scalar rewards for reinforcement learning-based preference optimization. Importantly, all adaptation occurs through the few-shot examples alone—no architectural changes or parameter updates are required, making Activation RMs immediately deployable for new preference specifications.

4 PreferenceHack: A Few-Shot Reward Hacking Benchmark

Reward hacking—where certain model biases exploit confounding factors in reward functions rather than satisfying the intended objectives—remains a significant challenge for alignment. To evaluate the robustness of reward models against such exploitation, we introduce PreferenceHack, a novel evaluation benchmark specifically designed to assess reward models' susceptibility to common bias-based reward hacking behaviors.

4.1 BENCHMARK DESIGN

To the best of our best knowledge, PreferenceHack is the first benchmark that evaluates reward hacking *in a few-shot setting with a paired preference format*, allowing direct assessment of reward models' vulnerabilities to known biases.

The benchmark consists of six distinct splits across language and multimodal domains, with each split containing 80 few-shot training examples and 920 evaluation examples. This structure allows

robust evaluation of reward models across diverse bias conditions with strong statistical power. More details about our dataset and its construction are included in Sec C.2 of the Supp.

4.2 Dataset Construction

4.2.1 LANGUAGE SPLITS

For the language-based splits, we built upon findings from the "Helping or Herding?" study [15], which documented exploitable biases in language models. We used high-quality ground truth answers from the original dataset and generated preference pairs by systematically injecting three well-known biases into the incorrect samples: (i) **Length Bias**: Models often assign higher scores to longer responses regardless of content quality. We generated longer alternatives to the incorrect responses while preserving their factual inaccuracies; (ii) **Format Bias**: Structured formats like numbered lists often receive higher scores despite potential content issues. We reformatted incorrect responses into structured formats to exploit this bias; and (iii) **Positivity Bias**: Responses containing positive attitudes tend to score higher. We injected positive tone into incorrect responses to trigger this bias.

To ensure consistency in generating the non-preferred responses, we used GPT-4o-mini to inject the bias being evaluated into the incorrect response.

4.2.2 MULTIMODAL SPLITS

For multimodal evaluation, we created three splits using image-prompt pairs from SUGAR-CREPE [23], a challenging compositional image-text retrieval dataset. Each pair in the multimodal split of PreferenceHack contains an image with a correct and incorrect prompt description. Similar to our language splits, we used GPT-40-mini to inject model biases into the incorrect descriptions while preserving their factual errors. This approach creates a test bed for assessing multimodal reward hacking vulnerabilities.

4.3 EVALUATION PROTOCOL

PreferenceHack employs a few-shot evaluation protocol where reward models are exposed to a small set of examples (80 per split) before being evaluated on the larger test set (920 examples per split). This format specifically tests the ability of reward models to quickly adapt to model biases given few-shot examples. We show some examples of our benchmark in Figure 2.

For each preference pair, a reward model is considered successful if it assigns a higher reward score to the correct response compared to the biased alternative. This simple evaluation metric directly measures a reward model's susceptibility to common exploitation patterns.

5 EVALUATION

We evaluate Activation RMs across a diverse set of benchmarks to assess their effectiveness as few-shot reward models and their ability to mitigate reward hacking. We apply our approach to two state-of-the-art Large Multimodal Models: LLaVA-OneVision-7B and Qwen2.5-VL-7B. Our experiments focus on comparing Activation RMs against existing few-shot approaches in standard reward modeling tasks, evaluating robustness against reward hacking, and assessing performance on multimodal retrieval tasks.

5.1 IMPLEMENTATION DETAILS

We implemented Activation RMs using PyTorch [43]. We used the official implementations of LLaVA-OneVision-7B [31] and Qwen2.5-VL-7B [2] as base models. All experiments were conducted on a single NVIDIA A100 GPU with 80GB memory. For the activation steering procedure directly edit the output of each attention head before the projection layer.

For each experiment, we used a consistent few-shot setting with $n \le 130$ examples for training our Activation RMs unless otherwise specified. The activation extraction process involves collecting attention head activations from the last token of the input prompt. For attention head selection, we

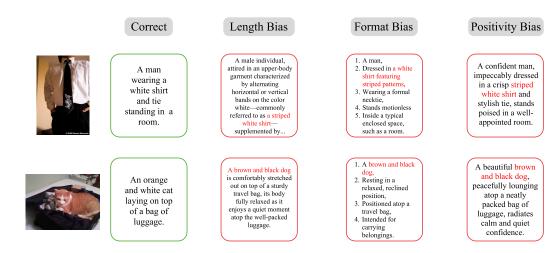


Figure 2: **PreferenceHack Examples.** We show samples based on two images of our PreferenceHack benchmark. Each sample would consist of a ground truth response paired with a biased incorrect response. The reward model is tasked with preferring the correct description over the biased one.

Table 1: **Evaluation of Activation RM on RewardBench and Multimodal Reward Benchmarks.** We perform a thorough evaluation of Activation RMs and baselines across multiple splits in language-only and multimodal settings. We present GPT-40 as a reference closed-source result.

	Language-Only (RewardBench)					Multimodal (Multimodal RewardBench)									
	Safety	Chat	Chat Hard	Reaso- ning	Overall	Macro Avg.	Correct.	Pref.	Knowl.	Math	Coding	Safety	VQA	Overall	Macro Avg.
Method / Model	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)
GPT-40	85.74	94.74	73.01	90.93	87.63	86.10	50.91	48.09	60.20	59.11	54.42	85.19	47.48	55.43	57.92
LLaVA-OneVision-7B															
ZS LLM-as-a-Judge 8-shot LLM-as-a-Judge ZS Generative Scoring	68.85 58.69 49.51	82.89 43.42 55.26	40.49 45.09 50.61	52.81 49.65 47.43	57.93 50.71 49.09	61.26 49.21 50.70	53.54 57.61 48.88	51.81 59.16 49.05	55.28 55.80 48.00	53.14 58.07 52.60	57.88 50.22 50.00	4.90 38.10 49.21	49.82 46.07 50.84	48.04 51.57 49.88	46.62 52.15 49.80
3-sample voting SAV Activation RM	67.21 69.40 70.98	84.21 85.70 88.60	40.80 45.60 50.31	52.73 65.20 69.02	57.65 64.50 68.84	61.24 66.47 69.73	56.19 55.50 49.90	54.39 53.20 48.56	56.20 54.80 54.91	53.91 53.50 52.90	56.86 56.90 50.62	5.29 40.30 81.62	49.81 49.50 49.00	48.93 51.80 53.75	47.52 52.00 55.36
Qwen2.5-VL-7B															
ZS LLM-as-a-Judge 8-shot LLM-as-a-Judge ZS Generative Scoring 3-sample voting SAV Activation RM	75.90 80.00 50.00 77.05 76.50 78.03	88.16 87.72 46.05 89.91 90.20 94.74	58.59 61.35 52.45 57.67 56.80 57.06	70.64 73.02 50.19 69.18 74.30 78.86	71.97 74.56 50.06 71.52 74.50 77.24	73.32 75.52 49.67 73.45 74.45 77.17	65.92 64.30 61.66 66.53 64.50 63.29	64.89 64.89 62.98 64.70 62.50 65.84	59.20 60.60 48.20 59.00 58.70 56.40	57.03 58.07 53.12 56.77 56.53 59.64	57.30 54.87 53.76 59.29 54.77 60.18	79.63 76.19 49.21 80.16 100.00 98.15	74.95 73.74 62.24 74.58 72.00 76.82	66.88 65.98 57.20 67.06 67.50 69.27	65.56 64.66 55.88 65.86 67.00 68.62

use 600 optimization steps with the REINFORCE algorithm [62]. Additional implementation details and hyperparameters can be found in the Appendix.

5.2 Datasets

We evaluate Activation RMs on three paired preference datasets where models must identify the preferred response between two candidates: (i) **RewardBench** [28] and **MultimodalRewardBench** [64] are comprehensive reward modeling benchmarks that evaluate out-of-the-box pretrained LLMs and LMMs on a variety of different language-only and multimodal tasks; in both benchmarks, given a prompt, the model must choose between a preferred and non-preferred response; (ii) **Preference-Hack** evaluates reward models' susceptibility to reward hacking with seven splits (80 training, 920 evaluation examples each) across language and multimodal domains. It systematically injects biases (length, format, numerical, and orientation) to assess how quickly reward models can identify and mitigate exploitation patterns with minimal examples. More details are in Section C.2 of the Supp.

5.3 Baselines

We compare Activation RMs against several established reward modeling approaches: **LLM-as-a-Judge** prompts the model to directly output a preferred response given a pair in either zero-shot or few-shot (8 examples) settings; **Generative Verifier** [33, 72] derives preferences by comparing the

Table 2: **Evaluation of Activation RM on PreferenceHack Benchmark.** We thoroughly evaluate Activation RMs and baselines on our novel few-shot reward hacking benchmark: PreferenceHack. We present GPT-40 as a reference closed-source result.

	Lang	guage-Onl	y Splits	Multimodal Splits					
Method / Model	Length (%)	Format (%)	Positivity (%)	Image+Length (%)	Image+Format (%)	Image+Positivity (%)			
GPT-40	3.91	48.04	92.39	22.35	55.78	87.65			
LLaVA-OneVision-7B									
ZS LLM-as-a-Judge	14.46	44.89	59.24	28.30	51.20	54.75			
8-shot LLM-as-a-Judge	23.15	37.50	57.17	38.45	45.65	52.30			
ZS Generative Scoring	45.54	47.17	76.96	57.80	54.25	71.40			
3-sample voting	15.43	43.26	59.67	30.85	49.75	55.10			
SAV	45.80	75.30	86.45	60.25	78.40	80.65			
Activation RM	49.24	79.89	90.11	65.70	83.45	85.25			
Qwen2.5-VL-7B									
ZS LLM-as-a-Judge	1.41	41.63	88.70	18.75	48.30	82.15			
8-shot LLM-as-a-Judge	8.70	47.39	87.28	25.40	53.85	80.60			
ZS Generative Scoring	17.72	50.65	93.59	35.20	58.40	88.25			
3-sample voting	1.41	41.85	88.70	19.30	48.75	82.50			
SAV	73.50	65.75	93.80	78.65	70.35	88.90			
Activation RM	78.37	68.26	96.74	84.25	75.50	91.80			

probability of a "Yes" token when asked if responses meet specified criteria; **3-Sample Voting** - A natural language reward modeling approach that implements self-consistency through a chain-of-thought methodology. The model generates three independent evaluations for each response, and the final preference is determined by majority voting across these samples; **Sparse Attention Vectors** (SAVs) [35] - A method that leverages few-shot examples to extract features from the attention heads of a model for classification, enabling another comparable SOTA form of few-shot reward modeling.

6 RESULTS

We perform a thorough evaluation of our Activation Reward Model on multiple benchmarks and compare to a variety of baselines. We first present the results of our few-shot method on general reward model benchmarks which for each group used a maximum of 130 examples for activation steering. Following this, we focus on the application of our method to the domain of safety and reward hacking on our PreferenceHack benchmark which we used 80 examples per group for steering. Finally, we perform several ablations and additional experiments to probe important characteristics of our approach.

6.1 REWARD BENCHMARK RESULTS

We perform evaluation on two comprehensive reward model benchmarks (Language-Only) Reward-Bench [28] and Multmodal RewardBench [64] as shown in Table 1. On average across a variety of splits, our Activation Reward Model outperforms all zero-shot and few-shot open-source baselines on both language-only and multimodal benchmarks, suggesting the effectiveness of our approach. Furthermore, our approach closes the gap with a strong closed-source baseline such as GPT-40. This is especially important as GPT-40 and other closed source models are often used as a reward models or judges of open-source models' outputs. However, a clear advantage of our approach is the interpretability of using few-shot examples of a task to specify a reward signal. Thus, our approach is both a more aligned and interpretable reward score for model alignment. Interestingly, our results show that few-shot, generative verification, and voting baselines struggle to outperform zero-shot LLM-as-a-Judge, suggesting that reward modeling is a challenging domain for these common SOTA methods. This further highlights the effectiveness of Activation RM. Additional results are provided in Section A of our Supp.

Finally, we highlight our method's consistent gains on safety tasks, which we attribute to a property we refer to as *taskness*. To clarify, safety tasks are more well-defined and thus better captured through few-shot examples, which Activation RMs use to guide reward model activations. In contrast, domains like chat, reasoning, or math are broader and less specific. For example, a few safety

Table 3: Ablations. We conduct ablations on Activation RMs using Qwen2.5-VL on RewardBench.

Ablation Method	Safety (%)	Chat (%)	Chat Hard (%)	Reasoning (%)	Overall (%)	Macro Avg. (%)
ZS LLM-as-a-Judge	75.90	88.16	58.59	70.64	71.97	73.32
CoT baseline	73.93	88.60	51.23	69.95	70.18	70.93
CoT + Voting	74.59	89.47	52.15	70.25	70.71	71.62
LoRA Finetuning	77.50	92.41	59.44	72.40	73.56	73.51
Mean Activation Addition	65.82	81.37	42.15	61.28	62.47	62.66
Top PCA Vector Replacement	76.24	91.58	54.91	75.93	74.73	74.67
Mean Activation Difference	76.51	92.85	55.32	77.24	75.48	75.48
ActivationRM	78.03	94.74	57.06	78.86	77.24	77.17

examples clearly define safe vs. unsafe responses, while a few math examples are less informative due to the diversity of sub-tasks (e.g., geometry, number theory, complex analysis). Thus, we posit that our approach and few-shot reward modeling methods more generally may be more successful when the application is to a *well-specified* or more *fine-grained* task.

6.2 PreferenceHack Results

To evaluate the effectiveness on a critical safety-like task, we apply our method to our new benchmark PreferenceHack as shown in Table 2. When evaluated on multiple different reward hacking biases in both language-only and multimodal settings, we find our method significantly outperforms all baselines in protecting against common reward hacks, even outperforming GPT-40 on most splits. Reward hacking is a task that quickly changes as new methods are found to exploit model biases. Hence, our approach is perfectly suited for adapting a reward model to be robust to a new attack given just a few examples.

6.3 ABLATION STUDIES

We explore different properties and capabilities of our framework via a careful ablation study in Table 3 using Qwen2.5-VL-7B evaluated on RewardBench.

Effect of CoT on Activation RMs. We are also interested in how the common approach of generating a CoT reasoning chain before outputting a preference impacts Acitvation RM. To do this, examples are formulated using the prompt, responses, and a chain-of-thought reasoning chain. Inference is performed in two steps. First, a CoT reasoning chain is generated given the prompt and two responses. Then, the final preference is outputted conditioned on the prompt, responses, *and* CoT reasoning chain. We find interestingly that CoT reasoning in this manner has little effect on our results, suggesting a future area of exploration for Activation RM.

Activation RM Comparable w/ LoRA Finetuning. We are also motivated to compare our framework with the common approach of finetuning an LLM/LMM explicitly as a reward model. We apply rank 16 LoRA finetuning for 3 epochs using 130 examples. Interestingly, we find that Activation RM yields similar performance as finetuning a model for reward modeling. This demonstrates that our method is both effective as a reward model and sample-efficient, requiring no weight updates to the generative model solely for reward modeling.

6.4 Additional Results

Superior Robustness to Label Noise. Figure 3 reveals ActivationRM's resilience to label corruption across all PreferenceHack splits. While LoRA fine-tuning exhibits catastrophic degradation under noise dropping by over 50% in some cases ActivationRM maintains stable performance even with 30% label corruption. This robustness stems from our design choices: weighted PCA filters out noisy variations in the calibration data, while output token likelihood scoring provides a more stable signal than methods that directly optimize on potentially mislabeled examples. The activation-based approach effectively reduces multiple sources of variability that plague traditional fine-tuning and prompting methods under noisy conditions.

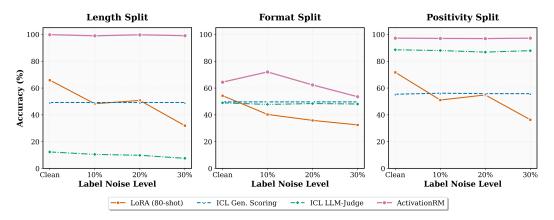


Figure 3: Robustness to label noise on PreferenceHack language splits. We evaluate ActivationRM against three baselines (LoRA fine-tuning, ICL Generative Scoring, and ICL LLM-as-a-Judge) across increasing levels of label noise (0% to 30%). Results are shown for Length, Format, and Positivity preference splits.

7 CONCLUSION

We introduce Activation Reward Models, the first mechanistic interpretability approach designed for few-shot reward modeling. By combining activation steering for precise task specification, weighted PCA denoising for robust preference extraction, and generative scoring for reliable evaluation, our method achieves state-of-the-art performance without any parameter updates. Our comprehensive evaluation demonstrates that Activation RMs consistently outperform existing few-shot approaches on both language-only and multimodal benchmarks, surpassing even GPT-40 on our novel PreferenceHack benchmark while providing greater interpretability through explicit few-shot examples.

Without extensive data collection or model retraining, the framework's flexibility enables fast deployment across diverse applications—from general evaluation tasks to best-of-N sampling and reinforcement learning—with adaptation occurring solely through few-shot examples. Our ablations suggest that performance can scale with more examples while maintaining few-shot practicality, and that our approach achieves comparable results to LoRA fine-tuning without requiring any weight updates. By enabling models to adapt to evolving preferences and emerging safety threats as shown by strong peformance on our novel PreferenceHack benchmark, Activation RMs provide a practical path toward more adaptive and robust AI alignment.

8 Limitations

Activation Reward Models represent a significant advancement in few-shot reward modeling, but several limitations should be acknowledged. First, our approach requires access to a model's internal architecture to extract and manipulate attention head activations, making it inapplicable to closed-source models like GPT-40 [39] and Claude [1]. Second, while Activation RM performs well on our benchmarks, the method's effectiveness may diminish for tasks that are less well-specified or require understanding of a broad range of criteria that cannot be captured in a few examples, such as mathematics. Finally, the current implementation focuses on single-turn interactions, and extending the approach to multi-turn dialogues or longer contexts may require additional research on how activation steering propagates across extended sequences. These limitations highlight opportunities for future work in developing more robust few-shot reward modeling techniques that can operate with more limited model access or handle more complex evaluation scenarios. Finally, we do not anticipate specific negative impacts, but as with any machine learning method, we recommend exercising caution in deployment.

REFERENCES

- [1] The claude 3 model family: Opus, sonnet, haiku. URL https://api.semanticscholar.org/CorpusID:268232499.
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *ArXiv*, abs/2502.13923, 2025. URL https://api.semanticscholar.org/CorpusID:276449796.
- [3] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862, 2022. doi: 10.48550/arxiv.2204.05862. URL https://arxiv.org/abs/2204.05862.
- [4] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Chris Olah, Benjamin Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. arXiv preprint arXiv:2204.05862, 2022.
- [5] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- [6] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, and Others. Constitutional AI: Harmlessness from AI feedback. arXiv preprint arXiv:2212.08073, 2022. doi: 10.48550/arxiv.2212.08073. URL https://arxiv.org/abs/2212.08073.
- [7] André Barreto, Vincent Dumoulin, Yiran Mao, Nicolás Pérez-Nieves, Bobak Shahriari, Yann Dauphin, Doina Precup, and Hugo Larochelle. Capturing individual human preferences with reward features. *arXiv preprint arXiv:2503.17338*, 2025.
- [8] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6541–6549, 2017.
- [9] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Understanding the role of individual units in a deep neural network. In *Proceedings of the National Academy of Sciences*, volume 117, pp. 30071–30077, 2020.
- [10] Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- [11] Guoqing Chen, Fu Zhang, Jinghao Lin, Chenglong Lu, and Jingwei Cheng. RRHF-V: Ranking responses to mitigate hallucinations in multimodal large language models with human feedback. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert (eds.), *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 6798–6815, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics. URL https://aclanthology.org/2025.coling-main.454/.
- [12] Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson E. Denison, John Schulman, Arushi Somani, Peter Hase, Misha Wagner, Fabien Roger, Vlad Mikulik, Sam Bowman, Jan Leike, Jared Kaplan, and Ethan Perez. Reasoning models don't always say what they think. 2025. URL https://api.semanticscholar.org/CorpusID:277668423.

- 540 [13] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. doi: 10.48550/arxiv.1706.03741. URL https://arxiv.org/abs/1706.03741.
 - [14] Carson Denison, Monte MacDiarmid, Fazl Barez, David Duvenaud, Shauna Kravec, Samuel Marks, Nicholas Schiefer, Ryan Soklaski, Alex Tamkin, Jared Kaplan, Samuel R. Bowman, Ethan Perez, and Evan Hubinger. Sycophancy to subterfuge: Investigating reward-tampering in large language models. *arXiv* preprint arXiv:2406.10162, 2024.
 - [15] Jacob Eisenstein, Chirag Nagpal, Alekh Agarwal, Ahmad Beirami, Alex D'Amour, Dj Dvijotham, Adam Fisch, Katherine Heller, Stephen R. Pfohl, Deepak Ramachandran, Peter Shaw, and Jonathan Berant. Helping or herding? reward model ensembles mitigate but do not eliminate reward hacking. *ArXiv*, abs/2312.09244, 2023. URL https://api.semanticscholar.org/CorpusID:266210056.
 - [16] Yossi Gandelsman, Alexei A Efros, and Jacob Steinhardt. Paint by word. In *arXiv preprint* arXiv:2103.10951, 2023.
 - [17] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*, 2024.
 - [18] Sahil Gureja, Zifan Xu, Aditi Chaudhary, Yuxuan Yao, Ajay Saini, Sabyasachi Ghosh, Gaurav Sahu, Preksha Nema, Barnabás Póczos, and Zachary C. Lipton. M-rewardbench: Evaluating reward models in multilingual settings. *arXiv preprint arXiv:2410.15522*, 2024. doi: 10.48550/arxiv.2410.15522. URL https://arxiv.org/abs/2410.15522.
 - [19] Roee Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 9318–9333, Singapore, December 2023. Association for Computational Linguistics.
 - [20] Roee Hendel, Mor Geva, and Amir Globerson. In-context learning creates task vectors. *arXiv* preprint arXiv:2310.15916, 2023.
 - [21] Evan Hernandez, Belinda Z. Li, and Jacob Andreas. Inspecting and editing knowledge representations in language models. *arXiv preprint arXiv:2304.00740*, 2023. URL https://arxiv.org/abs/2304.00740.
 - [22] Alberto Hojel, Yutong Bai, Trevor Darrell, Amir Globerson, and Amir Bar. Finding visual task vectors. In *European Conference on Computer Vision (ECCV)*, pp. 257–273. Springer, 2025.
 - [23] Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. Sugarcrepe: Fixing hackable benchmarks for vision-language compositionality. ArXiv, abs/2306.14610, 2023.
 - [24] Brandon Huang, Chancharik Mitra, Assaf Arbelle, Leonid Karlinsky, Trevor Darrell, and Roei Herzig. Multimodal task vectors enable many-shot multimodal in-context learning. In Advances in Neural Information Processing Systems (NeurIPS), 2024.
 - [25] Brandon Huang, Chancharik Mitra, Assaf Arbelle, Leonid Karlinsky, Trevor Darrell, and Roei Herzig. Multimodal task vectors enable many-shot multimodal in-context learning. In *Advances in Neural Information Processing Systems*, volume 37, pp. 22124–22153, 2024.
 - [26] Zhiyang Jin, Prakhar Gupta, Chun Kai Ling, Fuli Luo, Congzheng Song, Yuxi Yang, Xiang Ren, and Yuandong Tian. Rag-rewardbench: Benchmarking reward models in retrieval augmented generation for preference alignment. *arXiv preprint arXiv:2412.13746*, 2024. doi: 10.48550/arxiv.2412.13746. URL https://arxiv.org/abs/2412.13746.
 - [27] Katarzyna Kobalczyk, Claudio Fanconi, Hao Sun, and Mihaela van der Schaar. Few-shot steerable alignment: Adapting rewards and llm policies with neural processes. *arXiv preprint arXiv:2412.13998*, 2024.

- [28] Nathan Lambert, Valentina Pyatkin, Jacob Morrison, Liane Lovitt, Bill Yuchen Lin, Khyathi Raghavi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Rewardbench: Evaluating reward models for language modeling. arXiv preprint arXiv:2403.13787, 2024. doi: 10.48550/arxiv.2403.13787. URL https://arxiv.org/abs/2403.13787.
- [29] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Cărbune, Abhinav Rastogi, and Sushant Prakash. RLAIF vs. RLHF: Scaling reinforcement learning from human feedback with ai feedback. In *International Conference on Machine Learning (ICML)*, 2024.
- [30] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Ren Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. RLAIF vs. RLHF: Scaling reinforcement learning from human feedback with AI feedback. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, volume 235, pp. 26874–26901. PMLR, 2024.
- [31] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *ArXiv*, abs/2408.03326, 2024.
- [32] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023.
- [33] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *Computer Vision ECCV 2024*, volume 13799 of *Lecture Notes in Computer Science*, pp. 366–384. Springer, 2024.
- [34] Yuchun Miao, Sen Zhang, Liang Ding, Rong Bao, Lefei Zhang, and Dacheng Tao. Inform: Mitigating reward hacking in rlhf via information-theoretic reward modeling. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [35] Chancharik Mitra, Brandon Huang, Tianning Chai, Zhiqiu Lin, Assaf Arbelle, Rogerio Feris, Leonid Karlinsky, Trevor Darrell, Deva Ramanan, and Roei Herzig. Sparse attention vectors: Generative multimodal model features are discriminative vision-language classifiers. *arXiv* preprint arXiv:2412.00142, 2024.
- [36] Tong Mu, Alex Helyar, Johannes Heidecke, Joshua Achiam, Andrea Vallone, Ian D. Kivlichan, Molly Lin, Alex Beutel, John Schulman, and Lilian Weng. Rule based rewards for language model safety. *ArXiv*, abs/2411.01111, 2024. URL https://api.semanticscholar.org/CorpusID:273812284.
- [37] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv* preprint arXiv:2209.11895, 2022.
- [38] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.
- [39] OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander MÄĚdry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guarraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaresi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson,

649

650

651

652

653

654

655

656

657

658

659

661

662

663

665

667

668

669

670

671

672

673

674

675

676

677

679

680

684

685

686

687

688

689

690

691

692

696

697

700

Chak Li, Chan Jun Shern, Channing Conger, Charlotte Barette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, Dane Sherburn, Daniel Kappler, Daniel Levin, Daniel Levy, David Carr, David Farhi, David Mely, David Robinson, David Sasaki, Denny Jin, Dev Valladares, Dimitris Tsipras, Doug Li, Duc Phong Nguyen, Duncan Findlay, Edede Oiwoh, Edmund Wong, Ehsan Asdar, Elizabeth Proehl, Elizabeth Yang, Eric Antonow, Eric Kramer, Eric Peterson, Eric Sigler, Eric Wallace, Eugene Brevdo, Evan Mays, Farzad Khorasani, Felipe Petroski Such, Filippo Raso, Francis Zhang, Fred von Lohmann, Freddie Sulit, Gabriel Goh, Gene Oden, Geoff Salmon, Giulio Starace, Greg Brockman, Hadi Salman, Haiming Bao, Haitang Hu, Hannah Wong, Haoyu Wang, Heather Schmidt, Heather Whitney, Heewoo Jun, Hendrik Kirchner, Henrique Ponde de Oliveira Pinto, Hongyu Ren, Huiwen Chang, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian O'Connell, Ian Osband, Ian Silber, Ian Sohl, Ibrahim Okuyucu, Ikai Lan, Ilya Kostrikov, Ilya Sutskever, Ingmar Kanitscheider, Ishaan Gulrajani, Jacob Coxon, Jacob Menick, Jakub Pachocki, James Aung, James Betker, James Crooks, James Lennon, Jamie Kiros, Jan Leike, Jane Park, Jason Kwon, Jason Phang, Jason Teplitz, Jason Wei, Jason Wolfe, Jay Chen, Jeff Harris, Jenia Varavva, Jessica Gan Lee, Jessica Shieh, Ji Lin, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joanne Jang, Joaquin Quinonero Candela, Joe Beutler, Joe Landers, Joel Parish, Johannes Heidecke, John Schulman, Jonathan Lachman, Jonathan McKay, Jonathan Uesato, Jonathan Ward, Jong Wook Kim, Joost Huizinga, Jordan Sitkin, Jos Kraaijeveld, Josh Gross, Josh Kaplan, Josh Snyder, Joshua Achiam, Joy Jiao, Joyce Lee, Juntang Zhuang, Justyn Harriman, Kai Fricke, Kai Hayashi, Karan Singhal, Katy Shi, Kavin Karthik, Kayla Wood, Kendra Rimbach, Kenny Hsu, Kenny Nguyen, Keren Gu-Lemberg, Kevin Button, Kevin Liu, Kiel Howe, Krithika Muthukumar, Kyle Luther, Lama Ahmad, Larry Kai, Lauren Itow, Lauren Workman, Leher Pathak, Leo Chen, Li Jing, Lia Guy, Liam Fedus, Liang Zhou, Lien Mamitsuka, Lilian Weng, Lindsay McCallum, Lindsey Held, Long Ouyang, Louis Feuvrier, Lu Zhang, Lukas Kondraciuk, Lukasz Kaiser, Luke Hewitt, Luke Metz, Lyric Doshi, Mada Aflak, Maddie Simens, Madelaine Boyd, Madeleine Thompson, Marat Dukhan, Mark Chen, Mark Gray, Mark Hudnall, Marvin Zhang, Marwan Aljubeh, Mateusz Litwin, Matthew Zeng, Max Johnson, Maya Shetty, Mayank Gupta, Meghan Shah, Mehmet Yatbaz, Meng Jia Yang, Mengchao Zhong, Mia Glaese, Mianna Chen, Michael Janner, Michael Lampe, Michael Petrov, Michael Wu, Michele Wang, Michelle Fradin, Michelle Pokrass, Miguel Castro, Miguel Oom Temudo de Castro, Mikhail Pavlov, Miles Brundage, Miles Wang, Minal Khan, Mira Murati, Mo Bavarian, Molly Lin, Murat Yesildal, Nacho Soto, Natalia Gimelshein, Natalie Cone, Natalie Staudacher, Natalie Summers, Natan LaFontaine, Neil Chowdhury, Nick Ryder, Nick Stathas, Nick Turley, Nik Tezak, Niko Felix, Nithanth Kudige, Nitish Keskar, Noah Deutsch, Noel Bundick, Nora Puckett, Ofir Nachum, Ola Okelola, Oleg Boiko, Oleg Murk, Oliver Jaffe, Olivia Watkins, Olivier Godement, Owen Campbell-Moore, Patrick Chao, Paul McMillan, Pavel Belov, Peng Su, Peter Bak, Peter Bakkum, Peter Deng, Peter Dolan, Peter Hoeschele, Peter Welinder, Phil Tillet, Philip Pronin, Philippe Tillet, Prafulla Dhariwal, Qiming Yuan, Rachel Dias, Rachel Lim, Rahul Arora, Rajan Troll, Randall Lin, Rapha Gontijo Lopes, Raul Puri, Reah Miyara, Reimar Leike, Renaud Gaubert, Reza Zamani, Ricky Wang, Rob Donnelly, Rob Honsby, Rocky Smith, Rohan Sahai, Rohit Ramchandani, Romain Huet, Rory Carmichael, Rowan Zellers, Roy Chen, Ruby Chen, Ruslan Nigmatullin, Ryan Cheu, Saachi Jain, Sam Altman, Sam Schoenholz, Sam Toizer, Samuel Miserendino, Sandhini Agarwal, Sara Culver, Scott Ethersmith, Scott Gray, Sean Grove, Sean Metzger, Shamez Hermani, Shantanu Jain, Shengjia Zhao, Sherwin Wu, Shino Jomoto, Shirong Wu, Shuaiqi, Xia, Sonia Phene, Spencer Papay, Srinivas Narayanan, Steve Coffey, Steve Lee, Stewart Hall, Suchir Balaji, Tal Broda, Tal Stramer, Tao Xu, Tarun Gogineni, Taya Christianson, Ted Sanders, Tejal Patwardhan, Thomas Cunninghman, Thomas Degry, Thomas Dimson, Thomas Raoux, Thomas Shadwell, Tianhao Zheng, Todd Underwood, Todor Markov, Toki Sherbakov, Tom Rubin, Tom Stasi, Tomer Kaftan, Tristan Heywood, Troy Peterson, Tyce Walters, Tyna Eloundou, Valerie Qi, Veit Moeller, Vinnie Monaco, Vishal Kuo, Vlad Fomenko, Wayne Chang, Weiyi Zheng, Wenda Zhou, Wesam Manassra, Will Sheu, Wojciech Zaremba, Yash Patil, Yilei Qian, Yongjik Kim, Youlong Cheng, Yu Zhang, Yuchen He, Yuchen Zhang, Yujia Jin, Yunxing Dai, and Yury Malkov. Gpt-40 system card, 2024. URL https://arxiv.org/abs/2410.21276.

[40] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, and Others. Training language models to follow instructions with

- human feedback. *Advances in Neural Information Processing Systems (NeurIPS)*, 35, 2022. doi: 10.48550/arxiv.2203.02155. URL https://arxiv.org/abs/2203.02155.
 - [41] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*, 2022.
 - [42] Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*, 2023. URL https://arxiv.org/abs/2312.06681.
 - [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
 - [44] Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*, 2022.
 - [45] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018. URL https://api.semanticscholar.org/CorpusID: 49313245.
 - [46] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. arXiv preprint arXiv:2305.18290, 2023. doi: 10.48550/arxiv.2305.18290. URL https://arxiv.org/abs/2305.18290.
 - [47] Govind Ramesh, Yao Dou, and Wei Xu. GPT-4 jailbreaks itself with near-perfect success using self-explanation. *arXiv* preprint arXiv:2405.13077, 2024.
 - [48] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. doi: 10.48550/arxiv. 1707.06347. URL https://arxiv.org/abs/1707.06347.
 - [49] Sarah Schwettmann. Finding alignments between interpretable causal variables and distributed neural representations. In *arXiv preprint arXiv:2303.02536*, 2023.
 - [50] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
 - [51] Anikait Singh, Sheryl Hsu, Kyle Hsu, Eric Mitchell, Stefano Ermon, Tatsunori Hashimoto, Archit Sharma, and Chelsea Finn. FSPO: Few-shot preference optimization of synthetic preference data in LLMs elicits effective personalization to real users. arXiv preprint arXiv:2502.19312, 2025. doi: 10.48550/arxiv.2502.19312. URL https://arxiv.org/abs/2502.19312.
 - [52] Stanford Center for Research on Foundation Models (CRFM). Alpaca: A strong, replicable instruction-following model. https://crfm.stanford.edu/2023/03/13/alpaca.html, 2023. Accessed: 2025-05-20.
 - [53] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, Paul Christiano, Jan Leike, and Others. Learning to summarize from human feedback. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. doi: 10.48550/arxiv.2009.01325. URL https://arxiv.org/abs/2009.01325.
 - [54] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020.

- [55] Nishant Subramani, Nivedita Suresh, and Matthew Peters. Extracting latent steering vectors from pretrained language models. In *Findings of the Association for Computational Linguistics:* ACL 2022, pp. 566–581, Dublin, Ireland, May 2022. Association for Computational Linguistics.
 - [56] Granite Vision Team, Leonid Karlinsky, Assaf Arbelle, Abraham Daniels, Ahmed Nassar, Amit Alfassi, Bo Wu, Eli Schwartz, Dhiraj Joshi, Jovana Kondic, et al. Granite vision: a lightweight, open-source multimodal model for enterprise intelligence. *arXiv preprint arXiv:2502.09927*, 2025.
 - [57] Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function vectors in large language models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
 - [58] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023.
 - [59] Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. Steering language models with activation engineering. *arXiv* preprint *arXiv*:2308.10248, 2024. URL https://arxiv.org/abs/2308.10248.
 - [60] Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. Large language models are not fair evaluators. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 5605–5620, 2024.
 - [61] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
 - [62] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 2004.
 - [63] Zhaofeng Wu, Michihiro Yasunaga, Andrew Cohen, Yoon Kim, Asli Celikyilmaz, and Marjan Ghazvininejad. rewordbench: Benchmarking and improving the robustness of reward models with transformed inputs. *arXiv preprint arXiv:2503.11751*, 2025.
 - [64] Michihiro Yasunaga, Luke Zettlemoyer, and Marjan Ghazvininejad. Multimodal reward-bench: Holistic evaluation of reward models for vision-language models. *arXiv preprint* arXiv:2502.14191, 2025. doi: 10.48550/arxiv.2502.14191. URL https://arxiv.org/abs/2502.14191.
 - [65] Michihiro Yasunaga, Luke S. Zettlemoyer, and Marjan Ghazvininejad. Multimodal rewardbench: Holistic evaluation of reward models for vision language models. *ArXiv*, 2025.
 - [66] Kayo Yin and Jacob Steinhardt. Which attention heads matter for in-context learning? In *Proceedings of the 42nd International Conference on Machine Learning (ICML)*, 2025. URL https://arxiv.org/abs/2502.14010.
 - [67] Chao Yu, Qixin Tan, Hong Lu, Jiaxuan Gao, Xinting Yang, Yu Wang, Yi Wu, and Eugene Vinitsky. ICPL: Few-shot in-context preference learning via LLMs. *arXiv preprint arXiv:2410.17233*, 2024. doi: 10.48550/arxiv.2410.17233. URL https://arxiv.org/abs/2410.17233.
 - [68] Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. Rlhf-v: Towards trustworthy mllms via behavior alignment from fine-grained correctional human feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13807–13816, 2024.
 - [69] Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, Songfang Huang, and Fei Huang. Rrhf: Rank responses to align language models with human feedback without tears. *arXiv preprint arXiv:2304.05302*, 2023.

- [70] Hanning Zhang, Juntong Song, Juno Zhu, Yuanhao Wu, Tong Zhang, and Cheng Niu. Ragreward: Optimizing rag with reward modeling and rlhf. *arXiv preprint arXiv:2501.13264*, 2025.
- [71] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. *ArXiv*, 2024.
- [72] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2025.
- [73] Yi-Fan Zhang, Tao Yu, Haochen Tian, Chaoyou Fu, Peiyan Li, Jianshu Zeng, Wulin Xie, Yang Shi, Huanyu Zhang, Junkang Wu, et al. Mm-rlhf: The next step forward in multimodal llm alignment. *arXiv preprint arXiv:2502.10391*, 2025.
- [74] Bolei Zhou, David Bau, Aude Oliva, and Antonio Torralba. Interpreting deep visual representations via network dissection. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, volume 41, pp. 2131–2145, 2018.
- [75] Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019. doi: 10.48550/arxiv.1909.08593. URL https://arxiv.org/abs/1909.08593.