# Paper Critique: RLHF: Training Language Models to Follow Instructions with Human Feedback

Dan Peng

Department of Computer Science

Sep 28, 2024

danpeng@unc.edu

## 1. Research Problem

### 1.1. What research problem does the paper address?

This paper addresses the fundamental problem of aligning large language models (LLMs) with human intent. The authors observe that while scaling models improves capabilities, it doesn't inherently make them better at following user instructions or make them more helpful, honest, and harmless. The paper specifically tackles the challenge of developing methods that can effectively fine-tune pretrained language models to follow a wide range of natural language instructions.

### 1.2. What is the motivation of the research work?

The motivation stems from the observation that large language models trained solely on next-token prediction objectives often produce outputs that are untruthful, toxic, or unhelpful to users. Given that these models are increasingly deployed in real-world applications, addressing these alignment issues is crucial. The authors are motivated to find a scalable approach that doesn't sacrifice model capabilities while ensuring the models behave in accordance with human intentions and preferences.

## 2. Technical Novelty

### 2.1. What are the key technical challenges identified by the authors?

The authors identify several key technical challenges:

- Scaling reinforcement learning from human feedback (RLHF) to diverse instruction-following tasks rather than just specific domains like summarization

- Developing efficient methods for collecting high-quality human feedback across a broad distribution of tasks

- Mitigating the "alignment tax" - performance regressions on standard NLP benchmarks that can occur during alignment fine-tuning

- Creating evaluation methods that accurately measure model alignment on dimensions like helpfulness, harmlessness, and truthfulness

- Maintaining generalization to tasks outside the training distribution

### 2.2. How significant is the technical contribution of the paper? If you think that the paper is incremental, please provide references to the most similar work

The paper's technical contribution is highly significant, not merely incremental. While reinforcement learning from human feedback (RLHF) had been applied to language models before in specific domains like summarization [**?**, **?**], this work represents the first comprehensive application of RLHF to align general-purpose language models across diverse tasks. The paper demonstrates that this technique can be effectively scaled to align state-of-the-art language models (175B parameters) with human preferences while avoiding significant capability losses.

The work is particularly significant because it establishes a practical methodology for aligning foundation models that has since become standard practice in the industry. It combines three key components—supervised fine-tuning, reward modeling, and reinforcement learning—into an effective pipeline for alignment.

### 2.3. Identify 1-5 main strengths of the proposed approach.

- **Impressive empirical results with minimal parameter overhead**: The 1.3B InstructGPT model outperforms the 175B GPT-3 model in human evaluations despite having over 100x fewer parameters, demonstrating the efficiency of the alignment approach.

- **Reduced alignment tax**: The authors' PPO-ptx approach (mixing pretraining gradients during RLHF fine-tuning) effectively mitigates performance regressions on standard NLP benchmarks that typically occur during alignment.

- **Generalization beyond the training distribution**: The models demonstrate the ability to follow instructions in domains that are rare in the training data, such as code-related tasks and non-English instructions, showing promising out-of-distribution generalization.

- **Real-world validation**: The approach was tested on actual user prompts from the OpenAI API, ensuring practical relevance and addressing genuine use cases rather than contrived academic examples.

- **Comprehensive evaluation framework**: The authors developed a nuanced evaluation methodology addressing multiple dimensions of alignment (helpfulness, truthfulness, harmlessness) across diverse tasks.

## 2.4. Identify 1-5 main weaknesses of the proposed approach.

- **Lingering alignment issues**: Despite improvements, the models still make simple mistakes, fabricate facts, and sometimes generate harmful content, indicating that full alignment remains an unsolved problem.

- **Limited diversity in the labeler pool**: The human feedback comes from a relatively small group of 40 contractors who may not represent the full spectrum of user preferences and cultural backgrounds, potentially creating biased or culturally narrow alignment.

- **Prioritization of helpfulness over safety**: The training procedure prioritizes helpfulness to users, which can lead the model to generate potentially harmful content if explicitly requested, highlighting tensions between competing alignment objectives.

- **Challenging scalability**: While the approach works well for the model sizes tested, it remains unclear how efficiently this methodology will scale to much larger models where reward model training and RLHF fine-tuning may become prohibitively expensive.

## 3. Empirical Results

## 3.1. Identify 1-5 key experimental results, and explain what they signify.

- **Human preference for smaller aligned models over larger unaligned ones**: The 1.3B InstructGPT model outputs were preferred to those from the 175B GPT-3 model, demonstrating that alignment can overcome raw parameter count advantages. This signifies that alignment techniques can be more important than scale for producing outputs that humans find valuable.

- **Generalization to held-out labelers**: The models performed equally well when evaluated by labelers who didn't contribute to the training data, suggesting that the alignment generalizes reasonably well across different human evaluators within similar demographic groups.

- **Improvements in truthfulness**: InstructGPT models demonstrated approximately twice the truthfulness rate on the TruthfulQA benchmark compared to GPT-3 and showed a 50% reduction in hallucination rates on closed-domain tasks. This signifies that RLHF can effectively reduce factual fabrication without explicit factual supervision.

- **Modest improvements in toxicity**: InstructGPT models generate about 25% fewer toxic outputs than GPT-3 when prompted to be respectful, showing that alignment techniques can reduce harmful outputs, though not eliminate them entirely.

- **PPO-ptx preserves capabilities**: The PPO-ptx variant that mixes pretraining updates during RLHF fine-tuning significantly reduced performance regressions on standard NLP benchmarks compared to standard PPO, demonstrating that the alignment tax can be substantially mitigated.

## 3.2. Are there any weaknesses in the experimental section (i.e., unfair comparisons, missing ablations, etc)?

Several weaknesses exist in the experimental section:

- **Limited demographic diversity in evaluation**: The paper acknowledges that both training and evaluation were conducted primarily with English-speaking labelers from the US and Southeast Asia, leaving open questions about how well the alignment generalizes across diverse cultural contexts.

- **Missing comparison with alternative alignment methods**: The paper compares primarily against unaligned GPT-3 and supervised fine-tuning, but lacks comparisons with other alignment techniques like constitutional AI or chain-of-thought reasoning.

- **Lack of adversarial evaluation**: While the paper evaluates on some datasets like TruthfulQA, it doesn't include comprehensive adversarial testing to probe the boundaries and failure modes of the aligned models.

- **Insufficient exploration of the PPO-ptx hyper-parameters**: While the PPO-ptx approach shows promise, the paper could benefit from more thorough ablation studies on the mixing ratio and alternative formulations to better understand its properties.

- **Bias evaluation limitations**: The evaluation of model bias is limited to existing datasets like Winogender and CrowS-Pairs, which may not capture the full range of potential biases, particularly in the instruction-following context.

## 4. Summary

The paper presents a significant contribution to the field of AI alignment by demonstrating a practical approach to making language models more helpful, honest, and harmless through reinforcement learning from human feedback. The authors show that their method can align language models with human preferences across a diverse range of tasks, resulting in models that significantly outperform much larger unaligned models in human evaluations.

I am 85% impressed by this work and 15% concerned about its limitations. I'm particularly impressed by the empirical demonstration that alignment can outweigh raw parameter count in terms of producing outputs that humans prefer, suggesting a path forward where alignment research is as important as scaling. The effectiveness of the PPO-ptx approach in mitigating the alignment tax is also a crucial contribution that addresses a key practical challenge.

My concerns center on the limited diversity of human feedback sources, which may lead to models that are aligned primarily to specific cultural perspectives, and the ongoing challenges with truthfulness and harmfulness that remain despite the improvements. The paper honestly acknowledges these limitations, noting that "our models are neither fully aligned nor fully safe."