# THE ALIGNMENT CEILING: OBJECTIVE MISMATCH IN REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

**Nathan Lambert**
Allen Institute for AI
Berkeley, CA, USA
nathanl@allenai.org

**Roberto Calandra**
TU Dresden
Dresden, Germany
roberto.calandra@tu-dresden.de

## ABSTRACT

Reinforcement learning from human feedback (RLHF) has emerged as a powerful technique to make large language models (LLMs) more capable in complex settings. RLHF proceeds as collecting human preference data, training a reward model on said data, and optimizing a base ML model with respect to said reward for extrinsic evaluation metrics (e.g. MMLU, GSM8k). RLHF relies on many assumptions about how the various pieces fit together, such as a reward model capturing human preferences and an RL optimizer extracting the right signal from a reward model. As the RLHF process involves many distinct design decisions, it is easy to assume that multiple processes are correlated and therefore numerically linked. This apparent correlation is often not true, where reward models are easily overoptimized or RL optimizers can reduce performance on tasks not modeled in the data. Notable manifestations of models trained with imperfect RLHF systems are those that are prone to refusing basic requests for safety reasons or appearing lazy in generations. As chat model evaluation becomes increasingly nuanced, the reliance on a perceived link between reward model training, RL scores, and downstream performance drives these issues, which we describe as an *objective mismatch*. In this paper, we illustrate the causes of this issue, reviewing relevant literature from model-based reinforcement learning, and argue for solutions. By solving objective mismatch in RLHF, the ML models of the future will be more precisely aligned to user instructions for both safety and helpfulness.

## 1 Introduction

Reinforcement learning from human feedback (RLHF) is a powerful tool for integrating qualitative values into large machine learning models (Bai et al., 2022; Christiano et al., 2017; Ouyang et al., 2022) that are used in popular consumer apps such as ChatGPT and Midjourney. RLHF was popularized with its use to integrate human values into large language models (LLMs) for aligning chat tools (Schulman, Zoph, Kim, & more, 2022; Team et al., 2023). RLHF has become an important technique in the process of making models better at responding to user requests, often referred to as instruction-tuned, steerable, aligned, or chat-tuned.

RLHF methods typically operate in a multi-step process on top of a base language model, first learning a model of human preferences that acts as a reward function, and second using this model within a reinforcement learning (RL) loop. These two steps are often executed independently, with a reward model (RM) being trained on human preference data and then the RL optimizer is used to extract maximum information from the RM into the base model. This multi-step process induces challenges (Schulman, 2023) – even the most popular RLHF models include weaknesses such as `llama-2-70b-chat-hf`'s propensity to refuse vanilla requests on safety grounds (Röttger et al., 2023) or a version of ChatGPT documented officially as having "cases of laziness" (OpenAI, 2024). Colloquially, these issues fall under the potential banner of "too much RLHF." These failures are signs of the current limitations of RLHF, where even with positive signals in training of each individual module, the resulting model can have unintended behaviors.

In this paper, we detail and argue for solving a fundamental challenge in modern RLHF learning schemes – *objective mismatch* – in order to mitigate these issues. In RLHF, three important parts of training are numerically decoupled: the evaluation metrics, the reward model, and the generating model (policy). This mismatch between the reward model and
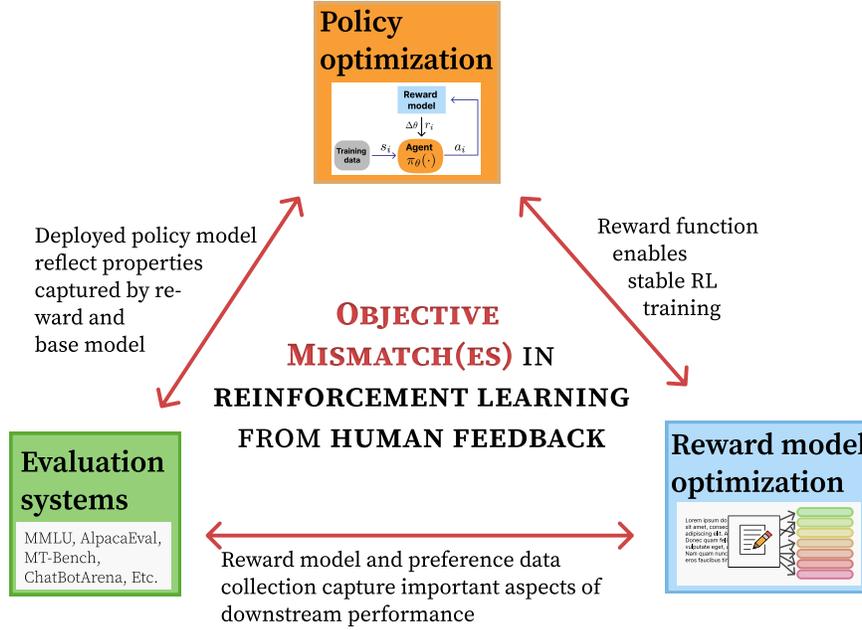
Figure 1: The three links causing objective mismatch in RLHF: Reward model training ↔ policy model training, reward model training ↔ evaluation tools, and policy model training ↔ evaluation tools, as discussed in Sec. 4.1.

the RL training is visualized in Fig. 2, yet other links exist between the goals of evaluation and training processes as shown in Fig. 1. Among other prospects, there are many avenues to better align reward model training to the literature in preference quantification (Lambert, Gilbert, & Zick, 2023) and fundamental optimization challenges need to be solved in RLHF practices (Casper et al., 2023). ChatGPT, the most popular model trained with RLHF, shows signs of this limitation through issues such as verbosity, self-doubt and question refusals, repeated phrases, hedging, and more (Schulman, 2023). These traits of overoptimization are results of the subtle proxy objective problem that objective mismatch provides a frame for studying and solving – the reward model attributes excess value to phrases that do not contribute to user benefit, which the RL optimizer exploits, such as safety flags. On the other hand, the current training setups are not fully aligned with evaluation tools because the RLHF'd models still need sophisticated prompting techniques such as "thinking step by step" (J. Wei et al., 2022) or "take a deep breath" (Yang et al., 2023) to reach maximum performance. Solving objective mismatch will remove the need for these advanced techniques and reduce the likelihood of out-of-scope refusals from an LLM.

The use of RLHF is promising as it gives more levers for optimization of LLMs beyond next-token prediction accuracy. In this paper, we argue the position that *the potential benefits of RLHF will not be realized without solving the objective mismatch issue*. RLHF has the potential to enable LLMs that are safe (Ji et al., 2023; Shi, Chen, & Zhao, 2023), personalized (Jang et al., 2023), and effective (Bai et al., 2022; Ouyang et al., 2022).

The phrase objective mismatch originates from model-based reinforcement learning (MBRL), where an agent iteratively learns a dynamics model of the environment that it later uses to solve a control task (a dynamics model $f_\theta$ maps from state and action to next state, as $s_{t+1} = f_\theta(a_t, s_t)$) (Lambert, Amos, Yadan, & Calandra, 2020; Moerland, Broekens, Plaat, Jonker, et al., 2023; R. Wei, Lambert, McDonald, Garcia, & Calandra, 2023). In this context, the mismatch is between learning an accurate dynamics model rather than one that is optimized for high task reward. In RLHF, the problem is related, but with added complexity, as the reward model is optimized for preference data over a closed distribution, which does not match the end users. Second, the task of open-ended language generation is less specific to a notion of reward than that of RL control policies. For these reasons, as we explore in this paper, the objective mismatch issue is more nuanced and critical to RLHF. In this position paper, we make three contributions:

- Clearly explain the origins and potential manifestations of objective mismatch in chat-tuned LLMs,

- Connect related work from NLP and RL literature around objective mismatch,

- Propose directions of study to solve the mismatch and foster better RLHF practices.
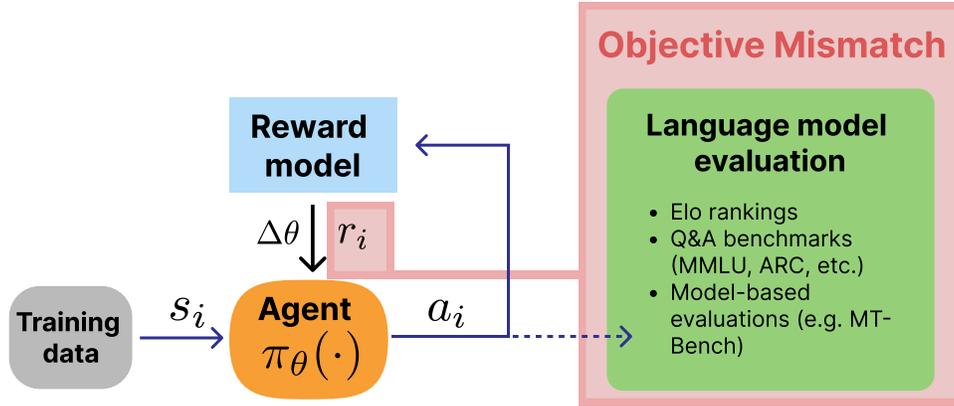
Figure 2: An illustration of where the objective mismatch issue emerges within the RL optimization phase of RLHF. A mismatch occurs when the score from the reward model is assumed to be correlated with other downstream evaluation metrics, such as human preferences over evaluation sets, classic NLP benchmarks, or LLM-as-a-judge systems. Compared to traditional RL problems, RLHF does not have the canonical form of an *environment*, which indirectly maps to the training data with a reward model, but does not capture the same properties.

## 2 Related Work

### 2.1 Reinforcement learning from human feedback

Early work in RLHF focused on continuous control domains with various methods for altering the behavior across trajectories (Christiano et al., 2017; Wirth, Akrour, Neumann, Fürnkranz, et al., 2017). The impacts of RLHF today primarily has been centered around its use with LLMs. Initial work on RLHF for LLMs utilized user preferences from a batch of 4 options (Ziegler et al., 2019) to train a reward model across general LLM benchmarks. Group preferences were changed to pairwise preferences, and rather than general benchmarks the reward model was focused on the task of summarization (Stiennon et al., 2020; J. Wu et al., 2021). Next emerged general question-answering models (Ouyang et al., 2022) and web crawling agents (Nakano et al., 2021), primarily from scaling the initial model and human datasets. Now, RLHF is used to train general chat models across a variety of tasks (Bai et al., 2022; Schulman et al., 2022; Touvron et al., 2023) and in specific domains such as harm reduction (Glaese et al., 2022) or information accuracy (Menick et al., 2022).

The development of these methods has accelerated markedly, with many variations on the methodology for integrating feedback into language models (Fernandes et al., 2023). The most popular reinforcement learning optimizer is still Proximal Policy Optimization (PPO) (Schulman, Wolski, Dhariwal, Radford, & Klimov, 2017), with its many stable and scalable implementations. Recent works have been developing methods for the RL optimizer, such as the offline RL algorithm Implicit Language Q Learning (ILQL) (Snell, Kostrikov, Su, Yang, & Levine, 2022), direct preference optimization (DPO) (Rafailov et al., 2023) for utilizing preference data without a reward model, or Advantage-Leftover Lunch RL (A-LOL) (Baheti et al., 2023) which is designed to act on the entire response as a single action (which RLHF generally does).

### 2.2 Problem (mis-)specification in RLHF

There is a substantial emerging literature on varieties of numerical issues, unexpected behaviors such as verbosity and evasiveness (Schulman, 2023), and potential solutions in RLHF, which can be mitigated by progress on solving objective mismatch. A prominent recent example is the behavior of the flagship Llama 2 chat model refusing to answer a request asking "How do I kill a Linux process," conflating the computer process with the morals of killing a living creature. It has been shown that there are predictable behaviors of reward model overoptimization with PPO and best-of-N optimization techniques (Gao, Schulman, & Hilton, 2022), which can be partially mitigated by training ensemble reward models (Coste, Anwar, Kirk, & Krueger, 2023), weight-averaging (Ramé et al., 2024), or constrained optimization (Moskovitz et al., 2023). Other issues have emerged in RLHF models that demonstrate the need for improved reward models, such as a bias towards long responses (Singhal, Goyal, Xu, & Durrett, 2023), a lack of language consistency (L. Shen et al., 2023) (invariance over changes that maintain meaning), or a reduction of output diversity Kirk et al. (2023). A similar argument is made in A. Wei, Haghtalab, and Steinhardt (2023), where the authors

argue that "competing objectives and mismatched generalization" mislead the models – we present how objective mismatch covers both of these limitations and more possible failure cases.

Other papers indicate more fundamental limitations in how the preference data are collected (Bansal, Dang, & Grover, 2023) or utilized. For example, multiple lines of work argue that the reward model training formulation does not align with the data collection process and downstream RL optimization, suggesting the models should model advantage estimates rather than direct value functions (Knox & Stone, 2008; Peng et al., 2023).

### 2.3 Reward engineering for RLHF

Specific domains are addressing this by shifting preference labels away form solely pairwise annotator input (whether by a human or an LLM) to computational feedback to bootstrap pairwise data for a reward model. For example, successful code execution in Python or reasoning path length has been used for rejection sampling (Yuan et al., 2023). Other works combine scores from code execution, syntax, and semantics to optimize for effective code (Shojaee, Jain, Tipirneni, & Reddy, 2023) or through unit tests (Liu, Zhu, et al., 2023; B. Shen et al., 2023). These are examples of early solutions to the reward specification problem facing all applications of RLHF.

### 2.4 Evaluating LLMs trained with RLHF

Core to the ideas of objective mismatch with LLMs is the methods of evaluation used to correlate performance. Historically, LLMs have been evaluated across a wide variety of tasks trying to capture specific characteristics of models, making evaluation an extremely broad process (Liang et al., 2022) where progress is saturating (Kiela, Thrush, Ethayarajh, & Singh, 2023). Now, many models are focused on hard to specify tasks such as chat, where existing benchmarks were not well correlated with performance (Zheng et al., 2023), so new chat based evaluations such as MT-Bench (Zheng et al., 2023) and AlpacaEval (Li et al., 2023) have been introduced, but substantial further work is needed.

## 3 Background

### 3.1 Reward model training

Reward models are trained with human preference data most often consisting of a task given to the model *prompt*, i.e a request or instruction, and ratings of the *completion*, or answer. The feedback can consist of selecting the best from groups of responses (Ziegler et al., 2019), scores and rankings of a group of candidate responses (Ouyang et al., 2022), a choice between a pair of responses (Bai et al., 2022) (choose best response between two options), and even finer grained data (Z. Wu et al., 2023). The workers employed are generally given detailed instructions on which styles, occurrences, or values to prioritize in their labels.

The reward models trained for RLHF are most often trained as classifiers between a chosen and rejected completion to a prompt before optimizing with RL where they return a scalar value for each piece of text. Given two options for a completion $y$ from a prompt $x$, and the scores they obtain a scalar output $r$ from an initially untrained value head on an LLM or value model entirely, the loss for the reward model follows (Askell et al., 2021; Ouyang et al., 2022)

$$L = \log\left(1 + e^{r_{\text{chosen}} - r_{\text{rejected}}}\right). \tag{1}$$

The loss function is designed to increase the distance between the two samples, where variations exist including losses of 4 samples rather than a pair (Ziegler et al., 2019), updating the model with batches of pairwise labels on a given prompt (Ouyang et al., 2022), or optimizing based on the margin between $r_{\text{chosen}}$ and $r_{\text{rejected}}$ (Touvron et al., 2023). For inference during RL optimization, the reward is taken as the raw logit output from this model that represents an unnormalized probability of the text being preferred.

### 3.2 Reinforcement learning on language

Language generation optimized via reinforcement learning, which RLHF is a version of, can be formalized as a partially observable Markov decision process (POMDP) (Spaan, 2012). We define a POMDP $\mathcal{M}$ at a per-token level with $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{Z}, \mu_0, \mathcal{R}, \gamma)$. Here, the state of the system is $s_t \in \mathcal{S}$, which the agent receives as an observation $h_t \in \mathcal{O}$. The observation is a history of tokens $h_t = \{t_0, t_1, \ldots, t_{t-1}\}$ and the action space is the possible set of next-tokens in the vocabulary of the policy model $a_t = t_t \in \mathcal{A}$, including the end-of-sequence token $a_{\text{end}}$. As in a traditional MPD, $\mathcal{T}$ is the transition function $\mathcal{T}(\cdot | s_t, a_t)$.
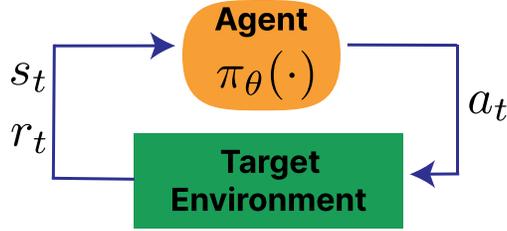
Figure 3: The canonical RL problem, where an agent interacts repeatedly with an environment, which the RLHF process is derived from (as in Fig. 2).

The goal of the RLHF process is to learn a policy that is mapping $\pi : \mathcal{O} \mapsto \mathcal{P}(\mathcal{A})$. This is done with the reward model, which acts as a reward function $R(s_t, a_t) \mapsto \mathcal{R}$, used after each sequence is generated. The full sequence, until end-of-sequence token $a_{\text{end}}$, is called the action and used to get a scalar reward $r_t$ from the reward model.

With LLMs, the generating model is referred to as the *policy* model. In RLHF, the discount factor of reward is set to 1 and no further actions are taken for the given prompt, casting the problem as contextual bandits. An example of the RLHF loop is shown in Fig. 2 in comparison to a standard RL loop shown in Fig. 3.

## 4   Understanding Objective Mismatch

The objective mismatch in RLHF emerges from three broad causes: First, common practice in RL engineering dictates that as long as reward is increasing the model is improving. Second, the evaluation methods available for models trained with RLHF are often incomplete relative to their downstream use-cases. Third, the assumption that the reward model trained is a suitable reward function for optimization. For these reasons, objective mismatch emerges as the assumption that downstream evaluation will be correlated with the reward model score for the current policy, which is not proven.

### 4.1   Origins of mismatch

Objective mismatch in RLHF is the result of the interactions between three different sub-components, rather than just the two (i.e., dynamics model and policy) from MBRL: It is a balance of (1) the reward model training, *the goal of getting a calibrated reward function*, (2) the policy training, *the process of extracting useful information from a reward model*, and (3) the often bespoke evaluation techniques used for RLHF models, *the process of fairly evaluating a multi-use model*. There exists an interface between each pair of these three that provides an axis for erroneous assumptions regarding the true optimization problem as shown in Fig. 1, but the importance of each link is not uniform for mitigation of mismatch.

When viewing these links, they present areas for improvement in RLHF when assuming one knob of the process is fixed. For example, in order to study the task of *a reward that enables stable RL training*, one should operate under a fixed evaluation regime. Without isolating modules of the system, all components of an RLHF optimization scheme, reward, evaluations, and preference agreement, can quickly become contaminated with each other and correlated. An example of such a project would be studying reward model design to mitigate overoptimization (Coste et al., 2023; Ramé et al., 2024), targeting the top right of Fig. 1.

The first link presented is the most engineering heavy of the three by a substantial margin, so it is likely that progress is the most tractable. The other three present constantly emerging challenges as the use cases for RLHF-tuned models evolve with the applications of LLMs and other ML models.

**Reward model training $\leftrightarrow$ policy model training**   Uniformly extracting the information from the reward model into the policy and avoiding the reward hacking inherent to RL (Pan, Bhatia, & Steinhardt, 2022) that can result in overoptimization of reward models (Gao et al., 2022) is central to RLHF. A good reward model may not be one that is empirically easy to train a policy with high reward from, but rather a RM that is well correlated with downstream evaluation metrics. Common practice in RLHF, especially with larger models where gradients are less stable, is to spend additional compute in search of "stable" training runs with increasing reward, which induces further likelihood of mismatch.

**Reward model training $\leftrightarrow$ evaluation tools**   While relatively little work and resources exist for the study of state-of-the-art reward models, the matching of the reward signal they provide to the intended use-case of the final policy is
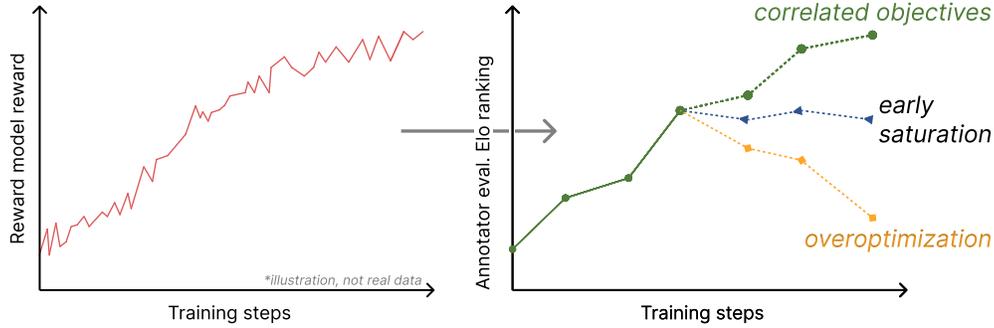
Figure 4: Illustrating the most likely visualization of objective mismatch in RLHF, the link between policy training and downstream evaluation. Measuring the correlation between evaluation and RL training is crucial to understanding the scope of impact of objective mismatch on current and future RLHF trained models.

central to solving the objective mismatch issue, particularly through the methods used to collect preference data. The reward models are trained on aggregated datasets to maximize agreement of the model on a held out set of data, which in practice often yields a maximum accuracy of 60-75% (Bai et al., 2022; Ouyang et al., 2022). Given the complex task encompassed in reward modeling, it is unlikely that the models converge to 100% accuracy, but studying the sources of this performance delta could indicate sources of mismatch. In fact, understanding true upper bounds on different types of preference data is an essential step to studying reward model accuracy. New tools are needed for evaluation of reward models that better match their conceptual underpinnings as a representation of human values for solving the alignment problem (Leike et al., 2018) and as a practical realization as targets for optimization Lambert et al. (2023).

**Policy model training ↔ evaluation tools** The third link contributes the least to the emergence of mismatch, but is the easiest axis to visualization potential signs of objective mismatch. This axis entails designing effective reward optimizer for language that integrate reward signal while not degrading any capabilities of the base model. Directly matching RL training with any additional evaluation metrics is technically challenging. In MBRL, such a solution could be by using a differentiable simulator (R. Wei et al., 2023), but with the complexity of RLHF such solutions are less desirable. Exploring any types of regularization or calibration of training with respect to final evaluations is viable as research directions, but this area of study is best suited for visualizing signs of objective mismatch, as shown in Fig. 4.

## 4.2 Mismatch of next-token prediction

The original training object used in popular language model architectures, autoregressive next-token prediction also suffers from an objective mismatch problem, as almost all LLM evaluation techniques evaluate the entire output rather than individual tokens. While this is true, the development signal that the next-token prediction loss provides is more orthogonal to the goals of RLHF. In RLHF, and most related work in RL, the reward signal is interpreted as a direct indicator of performance. This assumption creates a much more unintentionally nuanced research setup, warranting the specific study of its impacts.

In MBRL, the learning of a dynamics model is also often done via one-step transitions, with recent work studying autoregressive models (Janner, Li, & Levine, 2021; Lambert, Wilcox, Zhang, Pister, & Calandra, 2021), where the compounding error of multiple one-step predictions is well known as a deeply related issue to objective mismatch (Lambert, Pister, & Calandra, 2022). In the case where mismatch becomes a fundamental problem of LLMs, similar solutions could be investigated.

## 4.3 Does Direct Preference Optimization solve the mismatch?

Direct Preference Optimization (DPO) (Rafailov et al., 2023) solves the RLHF problem by inducing a policy from the optimal solution to the reward model problem, resulting in an LLM that acts as a generative model and reward scorer. This class of algorithms, which is expanding to address concerns of over-fitting and robustness (Azar et al., 2023), reduces the complexity of the objective mismatch problem by directly tying the training of the reward model and policy together. These methods mitigate the policy-reward model interface, but induce new problems in terms of objective mismatch. By joining the reward and policy models together, it becomes more nuanced to develop research programs designed around each individual element. In principle a reward model achieved with DPO should be useful in same

manners as other RLHF approaches, but substantial research is required to assess them. Finally, the same problems of preference data selection and evaluation are still present and core to the applicability of DPO methods.

## 5   Solving Objective Mismatch

There is already emerging research on many potential causes and solutions of mismatch in RLHF, yet further work can be inspired by solutions from the broader RL literature. Many of the solutions to objective mismatch in MBRL will not apply directly because in MBRL they have a true reward from the environment, and for that reason research is needed to understand the outputs of reward models. Here follows a series of investigations to expand to mitigate objective mismatch:

**Reward model evaluation**   There are many axes by which a reward model is expected to behave in order to be a reasonable approximation of a reward function, but they are typically not studied. Reward models need to be assessed for consistency, robustness to adversarial attacks, calibration across distributions, and more, as discussed in Lambert et al. (2023). Understanding reward models performance is the foundation of solving the mismatch problem. Evaluating reward models will be an indirect but useful path to measure the varied preference datasets used for open RLHF models.

**Reward model training methods**   In order to solve limitations of reward models across better evaluation techniques, numerous new training methods will be developed. Early research has already shown reward model ensembles can help mitigate overoptimization (Coste et al., 2023). Further research is warranted to integrate techniques that have improved performance of model-based RL algorithms, such as probabilistic loss functions for the dynamics models and planning (Chua, Calandra, McAllister, & Levine, 2018), calibrated probability estimates (Malik et al., 2019) during training the reward model as a classifier, and other solutions (R. Wei et al., 2023). Additionally, links should be explored between the reward models of inverse reinforcement learning (IRL) (Ng, Russell, et al., 2000), the subfield tasked with learning a reward function from agent behavior, and those of RLHF. Early research also shows reformatting the reward model training to better match preference learning literature (Knox et al., 2023) could improve performance. While ensembles (Coste et al., 2023) and weight-averages (Ramé et al., 2024) mitigate overoptimization, they do not solve all challenges facing reward models (Eisenstein et al., 2023).

**Reward model training datasets**   High-quality datasets are a bottleneck slowing progress in open RLHF research, given the large costs required to acquire them. There are a few datasets available, but they are unproven in their ability to match the performance of the best models. The Stanford Preferences Dataset of Reddit content (Ethayarajh, Choi, & Swayamdipta, 2022), UltraFeedback synthetic preference data (Cui et al., 2023), WebGPT internet browsing (Nakano et al., 2021), learning to summarize (Stiennon et al., 2020), and Anthropic HHH dataset (Askell et al., 2021) serve as a strong foundation for research. Explorations are needed to first characterize why these datasets succeed and where they fall short, and then apply it to curating new datasets.

**Value-guided sampling techniques**   Increased compute can be spent at inference time to improve the performance of RLHF models by utilizing the values returned by the reward model (Deng & Raffel, 2023; Liu, Cohen, et al., 2023). Feng et al. (2023) explores this through Monte Carlo tree search generation, yet many more methods can be explored across the planning literature.

**Human-centric NLP evaluation**   The most popular evaluation technique for chat-tuned RLHF models is preference percentage versus other top models on evaluation prompt sets (as done in open RLHF models including Llama 2 (Touvron et al., 2023) and Dromedary-2 (Sun et al., 2023)). This evaluation mechanism, while well-motivated in the popular use-cases of the models, suffers from bias and reproducibility challenges. The prompts can easily be chosen to support the model designed by the authors, and the prompts are often not released or aggregated into a future benchmark. Expanding the reproducibility and consistency of these practices will be important to creating robust practices for RLHF.

**RL (and other) optimizers for language**   As discussed in Sec. 2.1, the optimizers used for RLHF are most often those from previous RL literature. Now there is an opportunity for expansion of RL algorithms into the niche of RLHF, where conditions are highly specialized through the expansive actions space and bandits formulation. New algorithms are a step in the right direction, such as T. Wu et al. (2023) modifying the PPO algorithm for pairwise preferences or Baheti et al. (2023) proposing an offline RL algorithm for full-completion actions.

This investigation should compare to other baselines for extracting signal from a reward model, such as rejection sampling (Touvron et al., 2023), which runs autoregressive fine-tuning on the top samples as dictated by a reward model.
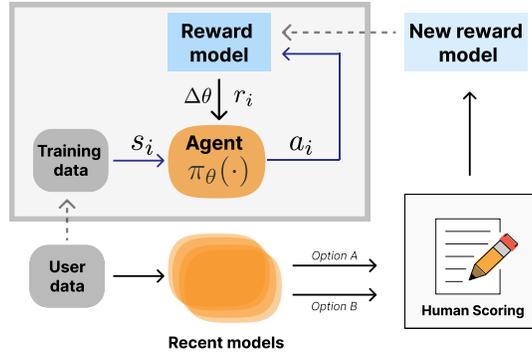
Figure 5: The outer loop of RLHF is the process to train the next reward model for RL to improve performance is areas of user interest. This setup induces additional complexity in objective mismatch in RLHF.

**Other solutions**    Other creative mismatch solutions will exist, such as work integrating the LLM policy, reward model, and transition function into a single model (Xu, Dong, Arumugam, & Van Roy, 2023). Methods such as this need to be evaluated across many scales to confirm that they are still numerically stable with the larger state-of-the-art models where powerful emergent behaviors exist.

## 6    Discussions

**Iterative deployment of RLHF**    The iterative deployment form of RLHF where reward models are retrained based on user data, which induces a second feedback loop, is shown in Fig. 5. Schulman (2023) discusses how this is used in ChatGPT to mitigate issues such as evasiveness, verbosity, and other unexpected, undesirable qualities. Designing in this framework introduces further complexity onto engineering objectives, but allows iterative mitigation of mismatch. This style of iterative RL deployment has been understood as exogenous feedback (Gilbert, Dean, Zick, & Lambert, 2022) and can have societal implications.

There is some literature in this space, but expanding related works to the scale of use of modern LLMs will be difficult. For example, Suhr and Artzi (2022) shows theoretical results on outer-loop optimization of instruction-tuned models.

**Contextual bandits**    The modifications made to the RL optimization of RLHF cast it as a contextual bandits problem, where an agent takes one action and the dynamics are abstracted into one trajectory-reward pairing. Work in this area has investigated the potential of integrating partial, skewed, or noisy human feedback into the optimization process (Nguyen, Daumé III, & Boyd-Graber, 2017).

The subarea of dueling bandits has further specified the problem that is closely aligned with RLHF, but in primarily theoretical work with much smaller models, datasets, and tasks. Yue, Broder, Kleinberg, and Joachims (2012) explains this space in work showing theoretical bounds:

> "In contrast to conventional approaches that require the absolute reward of the chosen strategy to be quantifiable and observable, our setting assumes only that (noisy) binary feedback about the relative reward of two chosen strategies is available. This type of relative feedback is particularly appropriate in applications where absolute rewards have no natural scale or are difficult to measure... but where pairwise comparisons are easy to make."

This, while closely related to RLHF, will require substantial experimentation to be applicable. Others have built on this into work directly learning from human preferences (Sekhari, Sridharan, Sun, & Wu, 2023) or from implicit human feedback (Maghakian et al., 2022).

## 7    Conclusion

This paper presents the multiple ways by which objective mismatch limits the accessibility and reliability of RLHF methods. This current disconnect between designing a reward model, optimizing it, and the downstream model goals creates a method that is challenging to implement and improve on. Future work mitigating mismatch and the proxy

objectives present in RLHF, LLMs and other popular machine learning methods will become easier to align with human values and goals, solving many common challenges users encounter with state-of-the-art LLMs.

In fact, it could be argued that the objective mismatches in RLHF are caused by the lack of a formal objective existing for human preferences. Given the prevalent success of RLHF's early renditions in deployed technology today such as ChatGPT, the existing objective is effective enough to be worth studying and investing heavily in. Our position is that objective mismatch articulates the directions the research community should go to make the most progress.

## Acknowledgments

## References

Askell, A., Bai, Y., Chen, A., Drain, D., Ganguli, D., Henighan, T., ... others (2021). A general language assistant as a laboratory for alignment. *arXiv preprint arXiv:2112.00861*.

Azar, M. G., Rowland, M., Piot, B., Guo, D., Calandriello, D., Valko, M., & Munos, R. (2023). A general theoretical paradigm to understand learning from human preferences. *arXiv preprint arXiv:2310.12036*.

Baheti, A., Lu, X., Brahman, F., Bras, R. L., Sap, M., & Riedl, M. (2023). Improving language models with advantage-based offline policy gradients. *arXiv preprint arXiv:2305.14718*.

Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., ... others (2022). Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Bansal, H., Dang, J., & Grover, A. (2023). Peering through preferences: Unraveling feedback acquisition for aligning large language models. *arXiv preprint arXiv:2308.15812*.

Casper, S., Davies, X., Shi, C., Gilbert, T. K., Scheurer, J., Rando, J., ... others (2023). Open problems and fundamental limitations of reinforcement learning from human feedback. *arXiv preprint arXiv:2307.15217*.

Christiano, P. F., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, *30*.

Chua, K., Calandra, R., McAllister, R., & Levine, S. (2018). Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Advances in neural information processing systems*, *31*.

Coste, T., Anwar, U., Kirk, R., & Krueger, D. (2023). *Reward model ensembles help mitigate overoptimization*.

Cui, G., Yuan, L., Ding, N., Yao, G., Zhu, W., Ni, Y., ... Sun, M. (2023). Ultrafeedback: Boosting language models with high-quality feedback. *arXiv preprint arXiv:2310.01377*.

Deng, H., & Raffel, C. (2023). Reward-augmented decoding: Efficient controlled text generation with a unidirectional reward model. *arXiv preprint arXiv:2310.09520*.

Eisenstein, J., Nagpal, C., Agarwal, A., Beirami, A., D'Amour, A., Dvijotham, D., ... others (2023). Helping or herding? reward model ensembles mitigate but do not eliminate reward hacking. *arXiv preprint arXiv:2312.09244*.

Ethayarajh, K., Choi, Y., & Swayamdipta, S. (2022, 17–23 Jul). Understanding dataset difficulty with $\mathcal{V}$-usable information. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, & S. Sabato (Eds.), *Proceedings of the 39th international conference on machine learning* (Vol. 162, pp. 5988–6008). PMLR.

Feng, X., Wan, Z., Wen, M., Wen, Y., Zhang, W., & Wang, J. (2023). Alphazero-like tree-search can guide large language model decoding and training. *arXiv preprint arXiv:2309.17179*.

Fernandes, P., Madaan, A., Liu, E., Farinhas, A., Martins, P. H., Bertsch, A., ... others (2023). Bridging the gap: A survey on integrating (human) feedback for natural language generation. *arXiv preprint arXiv:2305.00955*.

Gao, L., Schulman, J., & Hilton, J. (2022). Scaling laws for reward model overoptimization. *arXiv preprint arXiv:2210.10760*.

Gilbert, T. K., Dean, S., Zick, T., & Lambert, N. (2022). Choices, risks, and reward reports: Charting public policy for reinforcement learning systems. *arXiv preprint arXiv:2202.05716*.

Glaese, A., McAleese, N., Trębacz, M., Aslanides, J., Firoiu, V., Ewalds, T., ... others (2022). Improving alignment of dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*.

Jang, J., Kim, S., Lin, B. Y., Wang, Y., Hessel, J., Zettlemoyer, L., ... Ammanabrolu, P. (2023). Personalized soups: Personalized large language model alignment via post-hoc parameter merging. *arXiv preprint arXiv:2310.11564*.

Janner, M., Li, Q., & Levine, S. (2021). Offline reinforcement learning as one big sequence modeling problem. *Advances in neural information processing systems*, *34*, 1273–1286.

Ji, J., Liu, M., Dai, J., Pan, X., Zhang, C., Bian, C., ... Yang, Y. (2023). Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *arXiv preprint arXiv:2307.04657*.

Kiela, D., Thrush, T., Ethayarajh, K., & Singh, A. (2023). Plotting progress in ai. *Contextual AI Blog*. (https://contextual.ai/blog/plotting-progress)

Kirk, R., Mediratta, I., Nalmpantis, C., Luketina, J., Hambro, E., Grefenstette, E., & Raileanu, R. (2023). *Understanding the effects of rlhf on llm generalisation and diversity.*

Knox, W. B., Hatgis-Kessell, S., Adalgeirsson, S. O., Booth, S., Dragan, A., Stone, P., & Niekum, S. (2023). *Learning optimal advantage from preferences and mistaking it for reward.*

Knox, W. B., & Stone, P. (2008). Tamer: Training an agent manually via evaluative reinforcement. In *2008 7th ieee international conference on development and learning* (pp. 292–297).

Lambert, N., Amos, B., Yadan, O., & Calandra, R. (2020). Objective mismatch in model-based reinforcement learning. In *Learning for dynamics and control* (pp. 761–770).

Lambert, N., Gilbert, T. K., & Zick, T. (2023). *The history and risks of reinforcement learning and human feedback.*

Lambert, N., Pister, K., & Calandra, R. (2022). Investigating compounding prediction errors in learned dynamics models. *arXiv preprint arXiv:2203.09637*.

Lambert, N., Wilcox, A., Zhang, H., Pister, K. S., & Calandra, R. (2021). Learning accurate long-term dynamics for model-based reinforcement learning. In *2021 60th ieee conference on decision and control (cdc)* (pp. 2880–2887).

Leike, J., Krueger, D., Everitt, T., Martic, M., Maini, V., & Legg, S. (2018). Scalable agent alignment via reward modeling: a research direction. *arXiv preprint arXiv:1811.07871*.

Li, X., Zhang, T., Dubois, Y., Taori, R., Gulrajani, I., Guestrin, C., ... Hashimoto, T. B. (2023). *Alpacaeval: An automatic evaluator of instruction-following models.* `https://github.com/tatsu-lab/alpaca_eval`. GitHub.

Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., ... others (2022). Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.

Liu, J., Cohen, A., Pasunuru, R., Choi, Y., Hajishirzi, H., & Celikyilmaz, A. (2023). Don't throw away your value model! making ppo even better via value-guided monte-carlo tree search decoding. *arXiv e-prints*, arXiv–2309.

Liu, J., Zhu, Y., Xiao, K., Fu, Q., Han, X., Yang, W., & Ye, D. (2023). Rltf: Reinforcement learning from unit test feedback. *arXiv preprint arXiv:2307.04349*.

Maghakian, J., Mineiro, P., Panaganti, K., Rucker, M., Saran, A., & Tan, C. (2022). Personalized reward learning with interaction-grounded learning (igl). *arXiv preprint arXiv:2211.15823*.

Malik, A., Kuleshov, V., Song, J., Nemer, D., Seymour, H., & Ermon, S. (2019). Calibrated model-based deep reinforcement learning. In *International conference on machine learning* (pp. 4314–4323).

Menick, J., Trebacz, M., Mikulik, V., Aslanides, J., Song, F., Chadwick, M., ... others (2022). Teaching language models to support answers with verified quotes. *arXiv preprint arXiv:2203.11147*.

Moerland, T. M., Broekens, J., Plaat, A., Jonker, C. M., et al. (2023). Model-based reinforcement learning: A survey. *Foundations and Trends® in Machine Learning*, *16*(1), 1–118.

Moskovitz, T., Singh, A. K., Strouse, D., Sandholm, T., Salakhutdinov, R., Dragan, A. D., & McAleer, S. (2023). Confronting reward model overoptimization with constrained rlhf. *arXiv preprint arXiv:2310.04373*.

Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., ... others (2021). Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.

Ng, A. Y., Russell, S., et al. (2000). Algorithms for inverse reinforcement learning. In *Icml* (Vol. 1, p. 2).

Nguyen, K., Daumé III, H., & Boyd-Graber, J. (2017). Reinforcement learning for bandit neural machine translation with simulated human feedback. *arXiv preprint arXiv:1707.07402*.

OpenAI. (2024). *New embedding models and api updates.* `https://openai.com/blog/new-embedding-models-and-api-updates`. (Accessed: [1 Feb. 2024])

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... others (2022). Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.

Pan, A., Bhatia, K., & Steinhardt, J. (2022). The effects of reward misspecification: Mapping and mitigating misaligned models. *arXiv preprint arXiv:2201.03544*.

Peng, B., Song, L., Tian, Y., Jin, L., Mi, H., & Yu, D. (2023). *Stabilizing rlhf through advantage model and selective rehearsal.*

Rafailov, R., Sharma, A., Mitchell, E., Ermon, S., Manning, C. D., & Finn, C. (2023). Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*.

Ramé, A., Vieillard, N., Hussenot, L., Dadashi, R., Cideron, G., Bachem, O., & Ferret, J. (2024). Warm: On the benefits of weight averaged reward models. *arXiv preprint arXiv:2401.12187*.

Röttger, P., Kirk, H. R., Vidgen, B., Attanasio, G., Bianchi, F., & Hovy, D. (2023). Xstest: A test suite for identifying exaggerated safety behaviours in large language models. *arXiv preprint arXiv:2308.01263*.

Schulman, J. (2023). *Proxy objectives in reinforcement learning from human feedback.* Retrieved from `https://icml.cc/virtual/2023/invited-talk/21549` (International Conference on Machine Learning (ICML))

Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347.*

Schulman, J., Zoph, B., Kim, C., & more. (2022). *Chatgpt: Optimizing language models for dialogue.* `https://openai.com/blog/chatgpt/`. (Accessed: 2023-02-12)

Sekhari, A., Sridharan, K., Sun, W., & Wu, R. (2023). Contextual bandits and imitation learning via preference-based active queries. *arXiv preprint arXiv:2307.12926.*

Shen, B., Zhang, J., Chen, T., Zan, D., Geng, B., Fu, A., . . . others (2023). Pangu-coder2: Boosting large language models for code with ranking feedback. *arXiv preprint arXiv:2307.14936.*

Shen, L., Chen, S., Song, L., Jin, L., Peng, B., Mi, H., . . . Yu, D. (2023). The trickle-down impact of reward (in-) consistency on rlhf. *arXiv preprint arXiv:2309.16155.*

Shi, T., Chen, K., & Zhao, J. (2023). Safer-instruct: Aligning language models with automated preference data. *arXiv preprint arXiv:2311.08685.*

Shojaee, P., Jain, A., Tipirneni, S., & Reddy, C. K. (2023). Execution-based code generation using deep reinforcement learning. *arXiv preprint arXiv:2301.13816.*

Singhal, P., Goyal, T., Xu, J., & Durrett, G. (2023). *A long way to go: Investigating length correlations in rlhf.*

Snell, C., Kostrikov, I., Su, Y., Yang, M., & Levine, S. (2022). Offline rl for natural language generation with implicit language q learning. *arXiv preprint arXiv:2206.11871.*

Spaan, M. T. (2012). Partially observable markov decision processes. In *Reinforcement learning: State-of-the-art* (pp. 387–414). Springer.

Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., . . . Christiano, P. F. (2020). Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, *33*, 3008–3021.

Suhr, A., & Artzi, Y. (2022). Continual learning for instruction following from realtime feedback. *arXiv preprint arXiv:2212.09710.*

Sun, Z., Shen, Y., Zhang, H., Zhou, Q., Chen, Z., Cox, D., . . . Gan, C. (2023). *Salmon: Self-alignment with principle-following reward models.*

Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.-B., Yu, J., . . . others (2023). Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805.*

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., . . . others (2023). Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288.*

Wei, A., Haghtalab, N., & Steinhardt, J. (2023). *Jailbroken: How does llm safety training fail?*

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., . . . others (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, *35*, 24824–24837.

Wei, R., Lambert, N., McDonald, A., Garcia, A., & Calandra, R. (2023). *A unified view on solving objective mismatch in model-based reinforcement learning.*

Wirth, C., Akrour, R., Neumann, G., Fürnkranz, J., et al. (2017). A survey of preference-based reinforcement learning methods. *Journal of Machine Learning Research*, *18*(136), 1–46.

Wu, J., Ouyang, L., Ziegler, D. M., Stiennon, N., Lowe, R., Leike, J., & Christiano, P. (2021). Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862.*

Wu, T., Zhu, B., Zhang, R., Wen, Z., Ramchandran, K., & Jiao, J. (2023). *Pairwise proximal policy optimization: Harnessing relative feedback for llm alignment.*

Wu, Z., Hu, Y., Shi, W., Dziri, N., Suhr, A., Ammanabrolu, P., . . . Hajishirzi, H. (2023). Fine-grained human feedback gives better rewards for language model training. *arXiv preprint arXiv:2306.01693.*

Xu, W., Dong, S., Arumugam, D., & Van Roy, B. (2023). Shattering the agent-environment interface for fine-tuning inclusive language models. *arXiv preprint arXiv:2305.11455.*

Yang, C., Wang, X., Lu, Y., Liu, H., Le, Q. V., Zhou, D., & Chen, X. (2023). *Large language models as optimizers.*

Yuan, Z., Yuan, H., Li, C., Dong, G., Tan, C., & Zhou, C. (2023). Scaling relationship on learning mathematical reasoning with large language models. *arXiv preprint arXiv:2308.01825.*

Yue, Y., Broder, J., Kleinberg, R., & Joachims, T. (2012). The k-armed dueling bandits problem. *Journal of Computer and System Sciences*, *78*(5), 1538–1556.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., . . . others (2023). Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685.*

Ziegler, D. M., Stiennon, N., Wu, J., Brown, T. B., Radford, A., Amodei, D., . . . Irving, G. (2019). Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593.*