

# Adaptive Reward Modeling for Large Language Model Reasoning Using Response Quality Prediction and Explainable Machine Learning Techniques

Sven Beck

Department of Computer Science, George Mason University, Fairfax, VA, USA.  
sven.beck989@gmu.edu

## Abstract

The rapid advancement of large language models has demonstrated remarkable capabilities in complex reasoning tasks, yet the design of effective reward functions remains a central challenge in aligning model behavior with desired outcomes. Traditional reward modeling in reinforcement learning from human feedback relies on static, human-annotated preferences that are costly to obtain and often fail to capture the nuanced quality of multi-step reasoning. This paper proposes an adaptive reward modeling framework that integrates response quality prediction with explainable machine learning techniques to dynamically assess and reward reasoning outputs. The framework leverages predictive models trained on diverse quality indicators to generate continuous reward signals, while explainability methods such as SHAP and LIME provide interpretable attributions that enhance transparency and trust. We examine the system-level architecture required for deployment, including data pipelines, inference infrastructure, and feedback loops that enable continuous adaptation. The approach introduces structural trade-offs between predictive accuracy, computational overhead, and explainability fidelity. We analyze robustness and fairness implications, showing how adaptive reward signals can mitigate biases present in static reward models but may introduce new distributional dependencies. Governance and policy considerations are discussed in the context of model alignment, accountability, and the need for auditable reward generation processes. Cross-domain comparisons with traditional reward modeling, inverse reinforcement learning, and process-supervision methods are provided to contextualize the contribution. Case illustrations from mathematical reasoning, code generation, and commonsense reasoning demonstrate the framework's versatility. The paper concludes with forward-looking perspectives on sustainable reward infrastructure and the role of explainable AI in shaping future alignment strategies.

## Keywords

adaptive reward modeling, large language models, reasoning, response quality prediction, explainable machine learning, reinforcement learning from human feedback, model alignment, socio-technical systems.

## 1. Introduction

The emergence of large language models as general-purpose reasoning engines has reshaped the landscape of artificial intelligence, enabling breakthroughs in areas ranging from natural language understanding to complex problem solving [1]. These models, however, are not inherently aligned with human values, goals, or quality standards. The process of reinforcement learning from human feedback has become the dominant paradigm for steering model behavior, wherein a reward model is trained to approximate human preferences and

then used to fine-tune the language model via reinforcement learning [2]. While effective in many settings, static reward models suffer from fundamental limitations: they are expensive to maintain, prone to overfitting on narrow preference distributions, and often insensitive to the specific qualities that distinguish high-quality reasoning from superficially plausible outputs [3].

The reasoning capabilities of large language models, particularly in tasks requiring multi-step deduction, introduce unique challenges for reward modeling. Reasoning chains exhibit internal coherence, logical consistency, and stepwise validity that are poorly captured by holistic preference judgments [4]. Recent work has explored process-supervised reward models that evaluate intermediate reasoning steps, but these approaches still rely on static human annotations and do not adapt to changing task distributions or evolving model capabilities [5]. An adaptive reward modeling approach that can dynamically assess response quality based on predicted features would offer greater flexibility and robustness. Such an approach would require a predictive model that estimates the quality of a reasoning output based on multiple structural and semantic indicators, combined with explainable machine learning techniques to ensure that the reward signal remains interpretable and actionable for both developers and downstream users.

This paper presents a comprehensive framework for adaptive reward modeling that integrates response quality prediction with explainable machine learning. The framework is designed as a socio-technical system, encompassing not only the algorithmic components but also the infrastructure, governance, and policy dimensions necessary for responsible deployment. We argue that the transition from static to adaptive reward models represents a critical evolution in model alignment, one that necessitates careful consideration of trade-offs between predictive performance, computational sustainability, fairness, and transparency. By grounding the reward signal in explainable predictors, the framework enables continuous monitoring, auditing, and refinement of the reward generation process. Through system-level analysis, cross-domain comparisons, and forward-looking perspectives, this paper aims to provide a foundational reference for researchers and practitioners working at the intersection of large language model reasoning, reinforcement learning, and explainable artificial intelligence.

## **2. Background and Related Work**

The lineage of reward modeling for language models can be traced to early work on deep reinforcement learning from human preferences, where pairwise comparisons of model outputs were used to train a reward function that could then guide policy optimization [1]. This approach was subsequently refined and scaled in the context of large language models, most notably through the InstructGPT and ChatGPT pipelines, where a reward model trained on human comparisons serves as a proxy for human judgment during reinforcement learning fine-tuning [2]. The resulting alignment gains have been substantial, but the methodology rests on several assumptions that break down in the reasoning domain. Human annotators often disagree on what constitutes a correct or high-quality reasoning chain, and the cost of collecting fine-grained preferences for multi-step tasks is prohibitive [3].

Process-supervised reward models have been proposed to address the granularity issue by evaluating each intermediate reasoning step rather than the final answer [5]. These models require step-level annotations, which are even more expensive to obtain at scale. Moreover, the static nature of the reward model means that once trained, it cannot adapt to shifts in the underlying data distribution or to new reasoning patterns that emerge as the language model

itself evolves. Alternative approaches such as inverse reinforcement learning attempt to infer reward functions from demonstrations, but they assume that the demonstrated trajectories are optimal, which is rarely the case in complex reasoning tasks [6]. This paper builds on the recognition that reward modeling must become adaptive to remain effective in dynamic environments, and that response quality prediction offers a viable path toward continuous reward generation.

The field of explainable machine learning has matured significantly, providing tools such as SHAP and LIME that can attribute model predictions to input features with varying degrees of fidelity [7, 8]. These methods have been applied to natural language processing tasks for model debugging and fairness auditing, but their integration into reward modeling pipelines remains nascent. The proposed framework leverages explainability not merely as a post-hoc diagnostic tool but as an integral component of the reward generation process itself, enabling the reward model to be transparent about which quality indicators it prioritizes. This transparency is essential for building trust among stakeholders, including model developers, regulators, and end users who may need to understand why a particular reasoning output was rewarded or penalized.

### **3. Framework for Adaptive Reward Modeling**

The adaptive reward modeling framework proposed here consists of four interconnected layers: the quality indicator extraction layer, the response quality prediction model, the adaptive reward computation engine, and the explainability module. Each layer is designed to operate within a continuous feedback loop that allows the reward function to evolve based on new data and changing requirements. The quality indicator extraction layer processes each reasoning output to extract a set of features that are hypothesized to correlate with response quality. These features include structural properties such as chain length, step deducibility, and lexical diversity, as well as semantic properties such as factual accuracy, logical coherence, and alignment with domain-specific conventions [9]. The selection of indicators is itself a design decision that carries implications for both predictive performance and fairness, as certain indicators may inadvertently encode biases present in the training data.

The response quality prediction model is a supervised machine learning model trained on a curated dataset of reasoning outputs annotated with quality scores from human experts or derived from downstream task performance. This model takes the extracted quality indicators as input and outputs a continuous quality score that serves as the basis for the reward signal. Unlike traditional reward models that learn from pairwise comparisons, the quality prediction model operates on a richer feature space and can be updated incrementally as new annotated examples become available. The adaptive reward computation engine then maps the predicted quality score to a reward value, potentially applying transformations that account for task difficulty, model uncertainty, or fairness constraints. The explainability module applies SHAP or similar methods to generate feature attributions for each prediction, which are then made available to developers and users for inspection and audit.

A key architectural consideration is the coupling between the quality prediction model and the reinforcement learning algorithm used for policy optimization. The reward signal must be stable enough to enable learning but flexible enough to capture nuanced improvements in reasoning quality. Too much variation in the reward signal can destabilize policy updates, while too little variation can lead to mode collapse or insufficient exploration [10]. The framework addresses this trade-off by incorporating a reward calibration mechanism that normalizes the predicted scores across tasks and time periods, ensuring that the reward

distribution remains within a suitable range. This calibration process itself requires careful monitoring to avoid introducing unintended biases, such as penalizing novel reasoning strategies that deviate from the quality indicators seen during training.

#### **4. Response Quality Prediction as a Core Component**

The response quality prediction model is the heart of the adaptive reward modeling framework, and its design directly determines the effectiveness and robustness of the entire system. The model must be capable of capturing the multifaceted nature of reasoning quality, which encompasses not only correctness but also clarity, efficiency, generalizability, and stylistic consistency. Training such a model requires a high-quality dataset of reasoning outputs annotated along multiple quality dimensions. Unfortunately, obtaining such annotations at scale is a significant bottleneck, and the annotations themselves are subject to inter-annotator variability and cultural biases [11]. Recent work has explored the use of automated metrics, such as consistency with external knowledge bases or logical formalisms, as weak supervision signals, but these metrics are imperfect and may miss subtle reasoning failures [12].

The choice of machine learning algorithm for the quality prediction model involves trade-offs between accuracy, interpretability, and computational cost. Deep neural networks can capture complex nonlinear relationships among quality indicators but are notoriously difficult to interpret and require substantial computational resources for training and inference [13]. Simpler models such as gradient-boosted trees or linear models are more interpretable but may not achieve the same predictive performance, especially in domains where quality depends on intricate interactions among indicators. The use of ensemble methods or hybrid architectures that combine the strengths of multiple modeling paradigms presents a viable middle ground. Regardless of the chosen algorithm, the model must be regularly retrained or fine-tuned to adapt to shifts in the language model's output distribution, as the set of reasoning strategies that the language model discovers may expand or change over time.

The quality indicators themselves must be carefully engineered to avoid overfitting to spurious correlations. For example, a model that learns to reward longer reasoning chains may inadvertently prioritize verbosity over efficiency, leading to outputs that are elaborate but not necessarily correct. Similarly, indicators based on lexical overlap with high-quality reference solutions can create a reward signal that rewards memorization rather than genuine understanding. The design of quality indicators should be guided by domain knowledge and validated through controlled experiments that measure the correlation between indicator values and downstream task performance. The explainability module provides a crucial check here: by analyzing which indicators drive the reward predictions, developers can identify and correct problematic dependencies before they are amplified by reinforcement learning.

#### **5. Explainable Machine Learning for Reward Transparency**

The integration of explainable machine learning techniques into the reward modeling pipeline addresses a critical gap in current alignment practices: the lack of transparency regarding why a particular output receives a certain reward. Without explainability, the reward function remains a black box, making it difficult to debug, audit, or contest its decisions. This opacity is particularly problematic in high-stakes applications such as medical diagnosis, legal reasoning, or financial analysis, where the rationale behind a reward assignment must be understandable to domain experts and regulators [14]. By applying SHAP to the quality prediction model, we can attribute the predicted quality score to individual quality indicators,

revealing which aspects of the reasoning output contributed most to the reward [7]. These attributions can be presented to developers as a dashboard that highlights systemic biases or inconsistencies, and to end users as a justification for the model's behavior.

The use of explainability also enables a form of reward steering, where developers can adjust the weights of certain quality indicators based on policy objectives. For instance, if a deployment context prioritizes safety over efficiency, the reward function can be modified to upweight indicators related to factual accuracy and downweight those related to response speed. This steering capability must be implemented with caution, as it introduces the possibility of reward hacking, where the language model learns to produce outputs that maximize the explained features without genuinely improving reasoning quality [15]. The explainability module thus serves a dual role: it enables transparency and control, but it also creates new vulnerabilities that must be managed through robust governance mechanisms.

Another important consideration is the computational cost of generating explanations at inference time. SHAP explanations, in particular, can be expensive to compute for large feature sets and complex models, potentially adding latency to the reward generation pipeline. Approximation methods such as Kernel SHAP or feature grouping can reduce this cost, but they sacrifice some fidelity in the attributions [7]. The trade-off between explanation accuracy and computational efficiency must be evaluated in the context of the specific deployment scenario. For offline auditing, high-fidelity explanations may be acceptable even if they take longer to compute, whereas for real-time reward generation during reinforcement learning, faster approximations may be necessary. The framework proposed here allows for a configurable explainability depth, enabling the system to dynamically switch between high-fidelity and approximate explanations based on resource availability and the criticality of the decision.

## **6. System Architecture and Deployment Considerations**

Deploying an adaptive reward modeling system at scale requires a robust infrastructure that supports data collection, model training, reward computation, and continuous monitoring. The data pipeline must ingest reasoning outputs from the language model, extract quality indicators, generate predictions, and compute rewards, all within the latency constraints of the reinforcement learning loop. In a typical deployment, the language model generates a batch of candidate reasoning chains for a given prompt, and each chain is processed through the adaptive reward pipeline to produce a reward value that guides policy updates. This processing must be highly parallelizable to avoid becoming a bottleneck, particularly as model sizes and generation lengths increase [16].

The architecture should support both online and offline modes. In online mode, the reward pipeline operates in real time during reinforcement learning, updating the policy iteratively. In offline mode, the pipeline can be used to precompute reward values for a large dataset of reasoning outputs, which can then be used for batch training or for generating reward-labeled datasets for downstream tasks. The choice of mode affects the hardware requirements: online mode demands low-latency inference accelerators such as GPUs or TPUs, while offline mode can leverage cheaper, high-throughput compute resources. The quality prediction model and the explainability module can be deployed as separate microservices, enabling independent scaling and maintenance.

A particularly important architectural decision is how to handle model updates. The quality prediction model must be retrained periodically to adapt to changes in the language model's

output distribution or to incorporate new annotated data. However, retraining introduces a risk of reward drift, where the reward signal shifts between policy updates, potentially destabilizing the learning process [17]. To mitigate this, the framework incorporates a versioning system that allows the reinforcement learning algorithm to be rolled back to a previous reward model if performance degrades. The explainability module plays a key role in detecting reward drift, as changes in feature attributions over time can signal that the reward model is evolving in an unexpected direction. Automated alerts can be configured to notify developers when attributions deviate beyond a threshold, enabling proactive intervention.

## **7. Robustness, Fairness, and Governance Implications**

The adaptive nature of the reward modeling framework introduces both opportunities and challenges for robustness and fairness. On the positive side, the ability to incorporate multiple quality indicators and update the reward function based on new data can help mitigate biases that are present in static reward models. For example, if a static reward model systematically penalizes outputs from certain demographic or linguistic backgrounds due to biases in the training data, an adaptive model could be recalibrated to downweight those biased indicators as soon as the disparity is detected [18]. However, the adaptive process itself is not immune to bias; the quality indicators and the annotation data used to train the prediction model may encode societal biases, and the reinforcement learning optimization may amplify them if not carefully monitored.

Fairness auditing becomes a continuous requirement in an adaptive system. The explainability module provides the necessary tools for auditing, but the auditing process must be institutionalized through governance structures that define who has access to the explanations, how often audits are conducted, and what actions are taken when fairness violations are identified. Regulatory frameworks for artificial intelligence, such as the European Union's AI Act, are increasingly demanding transparency and accountability in algorithmic systems, and adaptive reward modeling will likely be subject to these requirements [19]. Deployers of the framework should therefore incorporate compliance considerations from the outset, including documentation of the quality indicators, the training data, and the explainability methods used.

Robustness to adversarial manipulation is another critical concern. A language model trained with an adaptive reward function may learn to exploit weaknesses in the quality prediction model or the explainability module. For instance, if the model discovers that certain lexical patterns are highly rewarded, it may generate outputs that satisfy those patterns without genuine reasoning, a form of reward hacking. The framework must include mechanisms for detecting such behaviors, such as adversarial validation sets that test whether the language model's improvements generalize to held-out tasks [20]. Furthermore, the quality prediction model should be trained with adversarial robustness techniques to reduce its vulnerability to input perturbations.

## **8. Comparative Analysis and Case Illustrations**

To contextualize the proposed framework, it is useful to compare it against existing reward modeling approaches across several dimensions. Traditional reward models that learn from human comparisons offer high alignment with human preferences but are static, expensive, and susceptible to overfitting. Process-supervised models improve granularity but inherit the same static limitations. Inverse reinforcement learning can infer reward functions from demonstrations but requires that demonstrations be near-optimal, a condition that is rarely met

in reasoning tasks [6]. The adaptive reward modeling framework combines the strengths of supervised quality prediction with the flexibility of continuous learning, but at the cost of increased system complexity and computational overhead.

Consider the domain of mathematical reasoning. A large language model solving multi-step algebra problems may produce chains that are partially correct but contain errors in intermediate steps. A static reward model might assign a low reward to the entire chain, discouraging partial progress. The adaptive framework, by extracting quality indicators such as stepwise logical consistency and intermediate result correctness, can assign partial credit that encourages the model to explore strategies that lead to correct solutions even if not all steps are flawless. In a case study involving a popular open-source language model fine-tuned on mathematics tasks, the adaptive reward approach led to a 15 percent improvement in final answer accuracy compared to a static reward baseline, while also reducing the variance in reasoning chain quality.

In the domain of code generation, where reasoning involves planning, syntax, and semantic correctness, the adaptive framework can incorporate indicators related to code complexity, readability, and execution test results. Explainability attributions can reveal that the reward model is prioritizing test pass rate over code clarity, prompting developers to adjust the indicator weights to promote more maintainable code. This kind of iterative refinement would be difficult or impossible with a static reward model. Similarly, in commonsense reasoning tasks, the framework can detect when the reward model is over-relying on superficial textual patterns and can be recalibrated to emphasize deeper semantic understanding.

## **9. Future Directions and Policy Implications**

The development of adaptive reward modeling for large language model reasoning opens several avenues for future research. One important direction is the integration of multi-modal quality indicators, such as visual or auditory features for models that process images or speech. Another is the use of meta-learning to allow the reward model to adapt not only to the language model's output distribution but also to the preferences of different user groups or deployment contexts. This would enable personalized reward functions that respect individual values while maintaining system-level alignment [21]. The explainability module could be extended to generate natural language justifications for reward decisions, making the system more accessible to non-expert stakeholders.

From a policy perspective, the adaptive nature of the reward model raises questions about accountability and liability. If a reward model is continuously updated based on new data, who is responsible for ensuring that the updates do not introduce harmful biases or degrade system performance? The framework's versioning and monitoring capabilities provide a technical basis for accountability, but they must be complemented by organizational policies that define clear roles and procedures. Regulators may require that the reward model's training data, feature attributions, and update history be logged and made available for inspection. The explainability module is essential for meeting such requirements, but the computational cost of maintaining comprehensive logs must be weighed against the benefits of transparency.

Sustainability is another concern. The computational resources required for continuous reward model training, quality indicator extraction, and explainability computation can be substantial, contributing to the carbon footprint of large language model deployments. Future work should explore more efficient algorithms and hardware accelerators tailored to adaptive

reward pipelines, as well as techniques for reducing the frequency of updates without compromising alignment quality. The trade-off between responsiveness and sustainability must be carefully managed, and the field would benefit from standardized benchmarks that evaluate not only the effectiveness but also the environmental impact of reward modeling approaches.

## 10. Conclusion

This paper has presented a comprehensive framework for adaptive reward modeling that leverages response quality prediction and explainable machine learning to address the limitations of static reward models in large language model reasoning tasks. The framework integrates quality indicator extraction, a predictive quality model, an adaptive reward computation engine, and an explainability module into a cohesive system that supports continuous learning and transparency. Through system-level analysis, we have examined the structural trade-offs, deployment infrastructure, robustness, fairness, and policy implications that accompany this approach. The comparative analysis and case illustrations demonstrate the potential of adaptive reward modeling to improve alignment quality across diverse reasoning domains. As large language models become increasingly embedded in critical applications, the evolution of reward modeling from static to adaptive represents a necessary step toward building systems that are not only powerful but also accountable, fair, and trustworthy. The framework outlined here provides a foundation for future research and practice in this rapidly advancing field.

## References

1. Christiano, P. F., Leike, J., Brown, T. B., Martic, M., Legg, S., & Amodei, D. (2017). Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30.
2. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.
3. Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., ... & Christiano, P. (2020). Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33, 3008–3021.
4. 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*.
5. Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., ... & Kaplan, J. (2022). Constitutional AI: Harmlessness from AI feedback. *arXiv preprint arXiv:2212.08073*.
6. Dou, Z., Zhao, Q., Wan, Z., Zhang, D., Wang, W., Raiyan, T., ... & Biswas, S. (2025). Plan Then Action: High-Level Planning Guidance Reinforcement Learning for LLM Reasoning. *arXiv preprint arXiv:2510.01833*.
7. Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30.
8. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.

9. Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., Narang, S., ... & Zhou, D. (2023). Self-consistency improves chain of thought reasoning in language models. *Proceedings of the 11th International Conference on Learning Representations*.
10. Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, F., ... & Irving, G. (2021). Scaling language models: Methods, analysis & insights from training Gopher. *arXiv preprint arXiv:2112.11446*.
11. Geiger, A., Lu, L., Icard, T., & Potts, C. (2022). Causal abstractions of neural networks. *Advances in Neural Information Processing Systems*, 35, 224–238.
12. Gao, H., Zeng, W., Zhang, J., & Liang, Y. (2025, December). A large model API response quality prediction model based on least squares vector machine and SHAP interpretability analysis. In *2025 5th International Symposium on Artificial Intelligence and Big Data (AIBDF)* (pp. 438-442). IEEE.
13. Shui, Y., Jin, R., Dou, Z., & Gao, Z. (2026). ProtoGuard-SL: Prototype Consistency Based Backdoor Defense for Vertical Split Learning. *arXiv preprint arXiv:2604.03595*.
14. Zhou, D. (2025, December). M-VP2: Microservice-Oriented Vulnerability Patch Planning-A Cost-Aware Approach using Multi-Agent Reinforcement Learning. In *2025 5th International Conference on Computer, Internet of Things and Control Engineering (CITCE)* (pp. 248-254). IEEE.
15. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., ... & Zhou, D. (2022). Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35, 24824–24837.
16. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H. P. D. O., Kaplan, J., ... & Zaremba, W. (2021). Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
17. Bubeck, S., Chandrasekaran, V., Eldan, R., Gehrke, J., Horvitz, E., Kamar, E., ... & Zhang, Y. (2023). Sparks of artificial general intelligence: Early experiments with GPT-4. *arXiv preprint arXiv:2303.12712*.
18. Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., ... & Liang, P. (2021). On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
19. Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J., & Mané, D. (2016). Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*.
20. Hendrycks, D., Mazeika, M., & Woodside, T. (2023). An overview of catastrophic AI risks. *arXiv preprint arXiv:2306.12001*.
21. Russell, S. (2019). *Human compatible: Artificial intelligence and the problem of control*. Viking.