

Rule-Bottleneck RL: Learning to Decide and Explain for Sequential Resource Allocation via LLM Agents

Guojun Xiong*

John A. Paulson School of Engineering and Applied Sciences, Harvard University
Cambridge, MA, United States
xiongj1@gmail.com

Haichuan Wang

John A. Paulson School of Engineering and Applied Sciences, Harvard University
Cambridge, MA, United States
haichuan_wang@g.harvard.edu

Mauricio Tec

John A. Paulson School of Engineering and Applied Sciences, Harvard University
Cambridge, MA, United States
mauriciogtec@gmail.com

Milind Tambe

John A. Paulson School of Engineering and Applied Sciences, Harvard University
Cambridge, MA, United States
milind_tambe@harvard.edu

ABSTRACT

Deep Reinforcement Learning (RL) has demonstrated remarkable success in solving sequential resource allocation problems, but often suffers from limited explainability and adaptability—barriers to integration with human decision-makers. In contrast, LLM agents, powered by large language models (LLMs), provide human-understandable reasoning but may struggle with effective sequential decision making. To bridge this gap, we introduce Rule-Bottleneck RL (RBRL), **the first LLM agent framework for resource allocation problems that jointly optimizes language-based decision policy and explainability.** At each step within RBRL, an LLM first generates candidate rules—language statements capturing decision priorities tailored to the current state. RL then optimizes rule selection to maximize environmental rewards and explainability, with the LLM acting as a judge. Finally, the LLM chooses the action (optimal allocation) based on the rule. We provide conditions for RBRL performance guarantees as well as the finite-horizon evaluation gap of the learned RBRL policy. Furthermore, we provide evaluations in real-world scenarios, particularly in public health, showing that RBRL not only improves the performance of baseline LLM agents, but also approximates the performance of Deep RL while producing more desirable human-readable explanations. We conduct a human survey validating the improvement in the quality of the explanations.

ACM Reference Format:

Guojun Xiong, Mauricio Tec, Haichuan Wang, and Milind Tambe. 2026. Rule-Bottleneck RL: Learning to Decide and Explain for Sequential Resource Allocation via LLM Agents. In *Proc. of the 25th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2026)*, Paphos, Cyprus, May 25 – 29, 2026, IFAAMAS, 9 pages.

1 INTRODUCTION

Sequential resource allocation is a fundamental problem in many domains, including healthcare, finance, public policy, and operations research [2, 5, 10, 60]. This task involves allocating limited

resources over time while accounting for dynamic changes and competing demands. Deep reinforcement learning (RL) is an effective method to optimize decision-making in resource allocation offering scalable high-reward policies [47, 57, 59], albeit generally providing action recommendations without human-readable reasoning and explanations. Such lack of interpretability poses a major challenge in critical high-stake domains where decisions must be transparent, justifiable, and in line with human decision-makers to ensure trust and compliance with ethical and regulatory standards.

For example, in healthcare settings, doctors may need to decide whether to prioritize intervention for Patient A or Patient B based on their current vital signs [5]. An RL algorithm might suggest: “*Intervene with Patient A*” with the implicit goal of maximizing the value function. However, the underlying reasoning may not be clear to the doctors, leaving them uncertain about the factors influencing the decision [26]. For doctors, a more effective suggestion could be risk-based with specific information, e.g., “*Patient A’s vital signs are likely to deteriorate leading to higher potential risk compared to Patient B, so intervention with Patient A is prioritized*” [4, 18].

LLM agents [45], on the other hand, leverage large language models (LLMs) for multi-step decision-making using reasoning techniques like chain of thought (CoT) [52]. They enable natural language goal specification [14] and enhance human understanding [22, 44]. However, agents based solely on LLM reasoning often struggle with complex sequential decision-making out of the box [17], making RL a crucial tool for grounding to specific tasks [6, 48, 53, 61].

Consequently, aiming to combine the strengths of both deep RL and LLM agents, we pose the following question:

Can we design an LLM agent framework that can simultaneously perform sequential resource allocation and provide human-readable explanations?

Similar to the celebrated index policy for prioritizing arms in resource allocation problems [54], we explore the potential of using rules-based prioritization in resource allocation tasks. In the context of LLM agents, rules are defined as “structured statements” that capture prioritization among choices in a given state, aligning with the agent’s goals [44]. Building on this, we propose a novel LLM agent framework called Rule-Bottleneck Reinforcement Learning

*Correspondence to: Guojun Xiong (xiongj1@gmail.com).

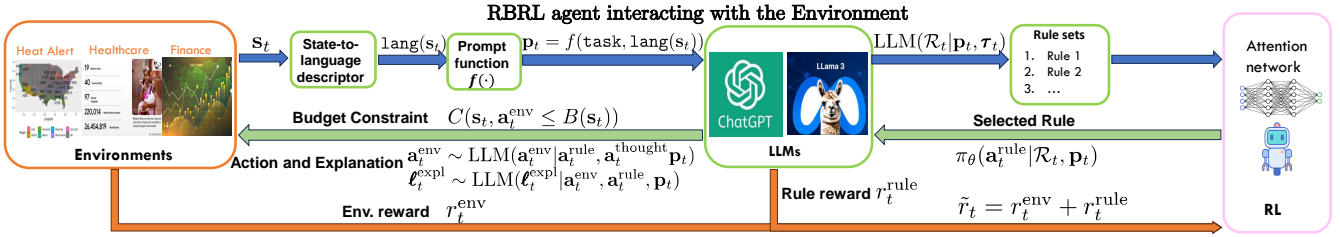


Figure 1: Overview of the framework of RBRL for joint sequential decision-making and explanation generation at time instance t . Starting with current state s_t , a state-to-language descriptor generates $\text{lang}(s_t)$, which is used to create the input prompt p_t . The LLM processes p_t to produce a thought τ_t and a set of candidate rules \mathcal{R}_t . An attention-based policy network selects a rule a_t^{rule} obeying the budget constraint $B(s_t)$, which is used by LLM to derive an executable action a_t^{env} for the environment and a human-readable explanation ℓ_t^{expl} , while also providing a rule reward r_t^{rule} . The environment transitions to the next state s_{t+1} , returning an environment reward r_t^{env} . This process is repeated iteratively at subsequent time steps.

(RBRL), which integrates the strengths of LLM and RL to bridge the gap between decision-making and interpretability. RBRL provides an agent framework (as shown in Figure 1) that *simultaneously* makes sequential resource allocation decisions and provides human-readable explanations, in contrast to prior work that generates post-hoc explanations for a learned policy [26, 30]. RBRL leverages LLMs to generate candidate rules and employs RL to optimize policy, allowing the creation of effective decision policies while simultaneously providing human-understandable explanations. RBRL aims to increase efficiency and avoid the computational cost of directly fine-tuning LLM agents, which can be highly challenging in interactive environments due to the heavy computational costs and the complexity of token-level optimization [33].

Our contributions are summarized as follows. *First*, LLMs are leveraged to generate a diverse set of rules according to the environment state, where each rule serves as a prioritization strategy for individuals in resource allocation, enhancing interpretability in decision-making. *Second*, we extend the conventional environmental state-action space by integrating the rules into states generated by LLMs, creating a novel framework that enables RL to operate on a richer, more interpretable decision structure. *Third*, we introduce an attention-based training framework that maps states and rules to queries and keys of a cross-attention network. The rule selection process is optimized by a policy network trained using the Soft Actor-Critic (SAC) algorithm [21], ensuring robust and efficient decision-making. In particular, the LLM also acts as a feedback mechanism, providing guidance during RL exploration to improve policy optimization and promote more effective learning. To the best of our knowledge, this is the first work to jointly optimize decision-making and explanation generation in constrained RL tasks.

We evaluate our method in environments from three real-world domains: *HeatAlerts*, where resources are allocated to mitigate extreme heat events; *WearableDeviceAssignment*, for distributing monitoring devices to patients; and *BinPacking*, which models allocating limited space in containers under constraints to optimize space utilization and minimize overflows. Using cost-effective LLMs such as gpt-4o-mini [27] and Llama 3.1 8B [25], we first assess decision performance by comparing RBRL with pure RL methods and language agent baselines. We then evaluate explanation quality through a human survey conducted under IRB approval. The

results demonstrate RBRL’s effectiveness in both decision quality and interpretability.

2 RELATED WORK

Our work intersects with three distinct areas within the RL literature. We discuss related work in each of these domains.

RL for Sequential Resource Allocation RL has been widely studied for constrained resource allocation across domains. In maternal health, [5] apply RL to a restless multi-armed bandit (RMAB) problem [54] to compute stochastic intervention probabilities. Also in an RMAB setting, Xiong et al. [56] propose a model-based RL approach that prioritizes users via an index and allocates resources under budget constraints. In public health, [10] propose RL to optimize extreme heat warnings under a budget on the number of possible alerts. Other works include multi-agent RL for robotic warehouse allocation [39] and exogenous MDPs for cloud resource management [42]. While these methods optimize rewards effectively, they often lack interpretability—critical for deployment in sensitive domains requiring trust, transparency, and accountability.

RL and LLM Agents One stream of research in LLM agents [45] has developed somewhat independently of RL, with works like ReAct prompting [58] extending chain-of-thought (CoT) [52] to action settings. These works have focused on tasks such as open-ended web navigation [32], social simulations [29], and virtual assistants [50]. Meanwhile, LLM agents have also been proposed for dealing with complex Markov decision processes such as GLAM [6], TWO-SOME [48], BAD [53], and AgentGym [55], which use LLM fine-tuning techniques in RL environments with a reward function. While our work is related to hierarchical methods that leverage LLMs for high-level planning [46, 51], our framework is novel in its objective. Unlike prior work that uses language solely to guide a policy toward high task rewards, RBRL is the first to treat the language-based “rule” as a primary output, jointly optimizing for both decision-making performance and the rule’s quality as a human-readable explanation via a dedicated reward signal.

Explainable RL (XRL) Early XRL relied on methods like decision trees and concept-based explanations [12], but these struggled with scalability in dynamic environments [31]. Recent advances introduced LLMS for post-hoc explanations, such as explaining decision paths from policy trees [62] or adding language descriptions to RL policies [9]. However, these approaches focus on interpreting

pre-existing policies rather than enabling LLMs to generate inherently explainable decisions, with challenges in aligning explanations to human reasoning [43]. By contrast, inherently (also known as intrinsically) interpretable policies are those that have internal representation that allow explanations [26, 30]. Our work sits within this literature by using LLM reasoning traces as the basis for environment action selection. Other methods like EDGE [20] and RICE [8] are primarily attribution-based; they identify which inputs (e.g., pixels or state features) were most critical to a decision. Similarly, SelfIE [7] provides a post-hoc, mechanistic explanation of an LLM’s internal mechanics (hidden states). In contrast, RBRL generates high-level, user-facing policy rules that are functional components within the RL loop, providing an explanation of the agent’s intent.

3 IMPACT STATEMENT AND LIMITATIONS

This work advances the development of transparent AI systems for high-stakes decision-making in domains like healthcare, public policy, industry, and many other applications. By enabling LLM agents to generate human-readable rules and explanations while attaining reward maximization via RL, RBRL improves trust and accountability, critical for ethical deployment in settings where lives and resources are at stake. While the framework prioritizes alignment with human reasoning, potential risks include over-reliance on imperfect LLM-generated rules or explanations that may inadvertently obscure biases in training data. Mitigation requires rigorous validation of rules by domain experts and ongoing monitoring of LLM outputs. Additionally, RBRL’s reliance on LLMs raises computational and accessibility challenges in resource-constrained environments. By addressing these considerations, this research contributes to safer, more equitable AI systems that empower—rather than replace—human decision-makers. To further validate the interpretability of our method, we obtained IRB approval and conducted a human subject study to evaluate the quality of the generated explanations.

Notice that the Uganda dataset used in this study is derived from a simulator that models vital sign trajectories of patients, as provided by [5]. Importantly, this simulator only replicates vital sign transitions and does not include any feature information or identifying details of real patients. Thus, the data generated by the simulator cannot be traced to or represent actual individuals, ensuring privacy and ethical compliance. We emphasize that this is purely a simulated patient study; and recognize that for any next steps towards real world use, there is a need to conduct rigorous simulation studies on a large scale with real patient data, with detailed assessments of potential biases, verification of policy convergence and its robustness to distribution shifts in patient populations, and making necessary adjustments. Beyond that, there will be a need to obtain ethics and regulatory approval to test the policy in a real-world setting for future comprehensive field testing, addressing issues of participant consent and privacy; and ultimately there would need to be sufficient human oversight for any future deployment.

Interpretability vs. Performance Tradeoff Various works acknowledge the trade-off between interpretability and performance [35]. In practice, prioritizing interpretability is crucial in practice in many high-stake applications: an approach that we subscribe to in this work. For example, in the clinical AI domain, high-performing

Step 1: Generate Thoughts
<p>Two example thoughts:</p> <ul style="list-style-type: none"> - There are only four warnings remaining in the budget. - The current heat index is high, and issuing alert could raise public awareness.
Step 2: Generate Rules Based on Thoughts and the Current State
<p>An example rule:</p> <ul style="list-style-type: none"> - Background: Maintaining a balance in warning issuance is crucial for future effectiveness - Rule: If there are 3 or more warnings remaining, issue a warning when the heat index is above 105 F. - State Relevance: There are 4 warnings remaining, allowing for proactive issuance given the current heat index of 107 F.

(a) Examples of generated rules for the Heat Alert Issuance task.

Example Language Wrapper for Heat Alert Issuance
<p>Task: Assist policymakers in deciding when to issue public warnings to protect against heatwaves. Your goal is to minimize the long-term impact on health and mortality. Your decision should be based on the remaining budget, weather conditions, day of the week, past warning history, and remaining warnings for the season. The goal is to issue warnings when they are most effective, minimizing warning fatigue and optimizing for limited resources.</p> <p>Action: A single integer value representing the decision: 1 = issue a warning, 0 = do not issue a warning. Warning can only be issued if the 'Remaining number of warnings/budget' is positive. Response in JSON format. For example: {'action': 1}.</p> <p>State: Remaining warning budget: 4, - Current date and day of summer: 2008-07-10, - Current heat index: 107 F.</p>

(b) Examples of language wrapper, containing task description, available actions and current state.

Figure 2: Examples of task prompts and generated rules for HeatAlerts domain.

black-box systems often face rejection in clinical workflows due to distrust [15, 40] as physicians require transparency to validate recommendations and uphold ethical accountability, as mandated by regulatory frameworks (e.g., [11, 16]). Interpretable models enable clinicians to audit biases, adapt decision logic to local contexts, and iteratively refine recommendations—fostering collaborative decision-making over reliance on inflexible oracles—whereas opaque policies are prone to failure under real-world distribution shifts [13, 15, 35, 40]. Empirical surveys show clinicians favor models that enable shared decision-making, error accountability, and ethical oversight despite modest performance penalties [40]—a critical stance in high-stakes healthcare environments where trust and adaptability outweigh narrow efficiency gains.

4 PRELIMINARY, KEY CONCEPTS, AND PROBLEM FORMULATION

4.1 Preliminary: Resource-Constrained Allocation

Resource-constrained allocation tasks are usually formulated as a special type of constrained Markov Decision Process, which is defined by the tuple $\langle \mathcal{S}, \mathcal{A}, P, R, C, h, \gamma \rangle$, where \mathcal{S} denotes a state space and \mathcal{A} denotes a finite action space. The transition probability function, specifying the probability of transitioning to state $s' \in \mathbb{R}^{d_1}$ after taking action $a \in \mathbb{R}^{d_2}$ in state s , is $P(s'|s, a) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \Delta(\mathcal{S})$, $R(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ represents the reward function, defining the immediate reward received after taking action a in state s , and we let $C(s, a) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^{d_3}$ be the immediate cost incurred after taking action a in state s . Often, each dimension $i \in [d_2]$ in a is either 0 or 1 in resource-constrained allocation tasks. In addition, h is the time horizon and $\gamma \in [0, 1]$ denotes the discount factor, which determines the present value of future rewards.

The goal is to find a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that maximizes the expected cumulative discounted reward while satisfying the cost constraints with a budget function $B : \mathcal{S} \rightarrow \mathbb{R}^{d_3}$:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} J(\pi) := \left[\sum_{t=1}^h \gamma^{t-1} R(s_t, a_t) \right], \quad (1)$$

$$s.t. \forall t \in [h] : C(s_t, a_t) \leq B(s_t).$$

4.2 Key Concepts for Rule-based LLM Agents

Our challenge is to design a rule-based LLM agent that jointly optimizes a language policy to both solve the optimization problem and improve explanation quality—a direction rarely explored. We next introduce the key concepts and terminologies underlying our main contribution.

LLM Agent For our LLM agent, the action space includes internal language actions $\tilde{\mathcal{A}} = \mathcal{A} \cup \mathcal{L}$ [58]. The LLM agent has two types of internal language actions: First, *thoughts* $\mathbf{a}^{\text{thought}} \in \mathcal{L}$, are reasoning traces from the current problem state used to inform environment action selection $\mathbf{a}^{\text{env}} \in \mathcal{A}$. Second, *explanations* ℓ^{expl} , are generated from actions and thoughts to enhance human trust and interpretability [62], a focus of this work.

Rules Thoughts are useful to highlight relevant aspects of a problem. However, they often lack detailed information to identify the next optimal action. In this work, we will consider “rules” $\mathbf{a}^{\text{rule}} \in \mathcal{L}$, which are structured language statements derived from thoughts that generally take the form “[if/when][do/prioritize]”. More formally, each rule \mathbf{a}^{rule} consists of a triple

$$(\text{background}, \text{rule_statement}, \text{state_relevance}).$$

Figure 2a shows examples of generated rules from one of the domains used in the experiments.

Task and Constraints Description Language agents require: (1) a language description of the environment and the agent’s goal, denoted task , containing the available actions for the task; (2) a function describing the state of the environment in natural language, denoted $\text{lang} : \mathcal{S} \rightarrow \mathcal{L}$. At each state s_t , these descriptors are used to construct a natural language prompt $\mathbf{p}_t = f(\text{task}, \text{lang}(s_t))$. Figure 2b exemplifies language wrapper generated for one of the environments in our experiments.

Rule-based Language Policy The objective is to jointly optimize the reward and explainability of the environment. Hence, we have an LLM agent-driven policy π_{LLM} for online interaction with the environment:

$$\begin{aligned} \mathbf{a}_t^{\text{thought}} &\sim \pi_{\text{LLM}}(\mathbf{a}_t^{\text{thought}} \mid \mathbf{p}_t), \mathbf{a}_t^{\text{rule}} \sim \pi_{\text{LLM}}(\mathbf{a}_t^{\text{rule}} \mid \mathbf{a}_t^{\text{thought}}, \mathbf{p}_t), \\ \mathbf{a}_t^{\text{env}} &\sim \pi_{\text{LLM}}(\mathbf{a}_t^{\text{env}} \mid \mathbf{a}_t^{\text{rule}}, \mathbf{a}_t^{\text{thought}}, \mathbf{p}_t), \\ \ell_t^{\text{expl}} &\sim \pi_{\text{LLM}}(\ell_t^{\text{expl}} \mid \mathbf{a}_t^{\text{env}}, \mathbf{a}_t^{\text{rule}}, \mathbf{p}_t). \end{aligned} \quad (2)$$

The rule acts as a “bottleneck” to the action and explanation. In the next section, we will introduce RBRL, which allows an RL-based learnable selection policy π_θ choosing among a set of dynamically generated candidate rules.

4.3 Problem Statement

We aim to increase the quality of ℓ^{expl} while also optimizing decision-making by selecting rules that encourage both good quality explanations and high reward. To achieve this goal, we construct a surrogate explainability “rule reward” $R_{\text{LLM}}^{\text{rule}}(\mathbf{a}^{\text{rule}})$ using an LLM as judge [3, 19, 38], which will be detailed in Section 5. Then, we propose the following augmented optimization objective under the joint environment/rule reward as $\tilde{R}(s_t, \mathbf{a}_t^{\text{env}}) = R(s_t, \mathbf{a}_t^{\text{env}}) + R_{\text{LLM}}^{\text{rule}}(\mathbf{a}_t^{\text{rule}})$:

$$\max_{\pi} \mathbb{E}_{\pi} \tilde{J}(\pi) := \left[\sum_{t=1}^h \gamma^{t-1} \tilde{r}_t \right], \text{ s.t. constraint in (1),} \quad (3)$$

where $\tilde{r}_t = \tilde{R}(s_t, \mathbf{a}_t^{\text{env}})$. We emphasize that LLMs cannot fully replace the ultimate human assessment, but they they provide a scalable alternative during the optimization process.

5 RULE-BOTTLENECK REINFORCEMENT LEARNING (RBRL)

In this section, we propose RBRL, a novel LLM agent based on the key concepts in Section 4.2, which leverages the strengths of LLMs and RL to achieve both interpretability and robust sequential decision-making for (3), thereby achieving our goal of jointly optimizing policies and explanations for resource-constrained allocation in (1).

Algorithm 1 RBRL

Require: Rule-selection policy π_θ ; and Replay buffer \mathcal{B} .

- 1: **Initialization:** Initial state \mathbf{s}_0 and task-state prompt \mathbf{p}_0 .
 - 2: **for** $t = 0, \dots, \text{max_iters} - 1$ **do**
 - 3: Generate rule candidates \mathcal{R}_t using CoT from \mathbf{p}_t and $\mathbf{a}_t^{\text{thought}}$.
 // Section 5.1
 - 4: Select rule $\mathbf{a}_t^{\text{rule}}$ using the RL policy π_θ from \mathcal{R}_t and \mathbf{s}_t . // Section 5.2
 - 5: Generate the environment action $\mathbf{a}_t^{\text{env}}$ with the LLM from $\mathbf{a}_t^{\text{rule}}$, \mathbf{p}_t , and previous thoughts.
 - 6: Apply the action in the environment and obtain new state \mathbf{s}_{t+1} and environment reward r_t^{env} .
 - 7: Generate explanation with the LLM from $\mathbf{a}_t^{\text{env}}$, $\mathbf{r}_t^{\text{rule}}$, \mathbf{p}_t , and previous thoughts.
 - 8: Generate rule reward r_t^{rule} with the LLM as judge. // Section 5.3
 - 9: Update the prompt \mathbf{p}_{t+1} from \mathbf{s}_{t+1} , and the constraints C and budget B .
 - 10: Append transition to the replay buffer $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\tilde{\mathbf{s}}_t, \mathbf{a}_t^{\text{rule}}, \tilde{r}_t, \tilde{\mathbf{s}}_{t+1})\}$.
 - 11: Sample from the replay buffer and update the policy network $\pi_\theta(\mathbf{a}_t^{\text{rule}} \mid \tilde{\mathbf{s}}_t)$. // Section 5.4
 - 12: **end for**
-

Algorithm Overview The framework of RBRL shown in Algorithm 1 involves four steps: (1) RULE SET GENERATION (line 3), where the LLM processes the state-task \mathbf{p}_t to create candidate rules \mathcal{R}_t for action selection; (2) RULE SELECTION (line 4), where an attention-based RL policy π_θ selects the best rule $\mathbf{a}_t^{\text{rule}} \in \mathcal{R}$; (3) DECISION, RULE REWARD AND EXPLANATION (lines 5-8), where the LLM generates an environment action $\mathbf{a}_t^{\text{env}}$ and based on the chosen rule $\mathbf{a}_t^{\text{rule}}$ gives a human-readable explanation ℓ_t^{expl} ; (4) REINFORCEMENT LEARNING (line 11), where it updates the policy π_θ based on collected data with standard RL algorithm [21] and the combined environment and rule reward \tilde{r}_t . Algorithm 1 details the entire process. Further sections elaborate on these steps.

5.1 Rule Set Generation

The rule generation process seeks to create interpretable and actionable guidelines for decision-making. Under this framework, a set of candidate rules \mathcal{R}_t is generated according to $\mathcal{R}_t \sim \pi_{\text{LLM}}(\mathcal{R}_t \mid \mathbf{p}_t, \mathbf{a}_t^{\text{thought}})$. To enhance interpretability, each rule is accompanied by a rationale

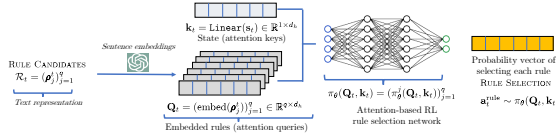


Figure 3: Overview of the RULE SELECTION step. The current state is encoded as a key vector, while candidate rules are encoded as Queries using a text embedding API (e.b., BERT sentence embedding). An attention-based policy network π_θ computes a probability distribution over the candidate rules, enabling the selection of the most suitable rule for decision-making and explanation.

explaining the reasoning behind the decision. The LLM is instructed to generate rules as a JSON format, which is common for integration of LLMs with downstream applications [38]. An example generated rule is given in Figure 2a.

5.2 Rule Selection

In this step, rules are converted from text to vector form, and a trainable attention-based policy network π_θ provides the probability distribution for sampling a rule. Figure 3 illustrates the process. Below are the major components of the architecture of π_θ . We propose to base the architecture on cross-attention layers [1, 49], with the state acting as the keys and values, and the rules as the queries. This allows to learn from the embedding representations of rules, and efficiently handle dynamically changing number of rules if needed.

State Representation The numeric state is projected by a linear layer: $\mathbf{k}_t = \text{Linear}(s_t) \in \mathbb{R}^{1 \times d_h}$, with d_h being to denote the architecture hidden dimension.

Rule Candidate Embedding Each rule in the list of rule candidates $\mathcal{R}_t = \{\rho^1, \rho^2, \dots, \rho^q\}$ is embedded into a numeric representation using a Sentence Embedding language model (e.g., SentenceBERT [34]) and further processed by a projection layer similar to the state representation. This results in a query matrix $\mathbf{Q}_t \in \mathbb{R}^{q \times d_h}$.

Attention-based Policy Network π_θ The vector \mathbf{k}_t , serving as keys, engages with the rule embeddings \mathbf{Q}_t , acting as queries, via a cross-attention mechanism to derive a hidden state representation $\mathbf{h}_t^{(1)} = \text{Attention}(\mathbf{Q}_t, \mathbf{k}_t^T, \mathbf{k}_t^T) \in \mathbb{R}^{q \times d_h}$, computed as $\text{Attention}(\mathbf{Q}_t, \mathbf{k}_t^T, \mathbf{k}_t^T) = \text{softmax}\left(\frac{\mathbf{Q}_t \mathbf{k}_t^T}{\sqrt{d_h}}\right) \mathbf{k}_t^T$, which facilitates the rule candidate vector embeddings in attending to the environment state. Subsequently, we sequentially apply self-attention layers to the hidden representation $\mathbf{h}^{(k+1)} = \text{Attention}(\mathbf{h}_t^{(k)}, \mathbf{h}_t^{(k)}, \mathbf{h}_t^{(k)})$, enabling the rule embeddings to attend to one another to rank an optimal candidate. Ultimately, following $K - 1$ self-attention layers, a final linear layer converts the rule representations into logit vectors $\alpha_\theta^i = \text{Linear}(\mathbf{h}_t^{(k)}) \in \mathbb{R}^q$ used for the computation of the probability of selecting each rule.

Rule Selection The policy distribution over the rules is calculated as: $\pi_{\theta,i}(\mathbf{Q}_t, \mathbf{k}_t) = \frac{\exp(\alpha_{\theta,i}^i(\mathbf{Q}_t, \mathbf{k}_t))}{\sum_{j=1}^q \exp(\alpha_{\theta,j}^i(\mathbf{Q}_t, \mathbf{k}_t))}$, $i = 1, \dots, q$. Therefore, a rule is selected at random from the distribution: $\mathbf{a}_t^{\text{rule}} \sim \text{Categorical}(\mathcal{R}; (\pi_{\theta,i}(\mathbf{Q}_t, \mathbf{k}_t))_{i=1}^q)$.

5.3 Decision, Rule Reward, and Explanation

Upon selection of rule $\mathbf{a}_t^{\text{rule}}$, the LLM determines the action to be applied within the environment $\mathbf{a}_t^{\text{env}} \sim \pi_{\text{LLM}}(\mathbf{a}_t^{\text{env}} | \mathbf{a}_t^{\text{rule}}, \mathbf{a}_t^{\text{thought}}, \mathbf{p}_t)$, ensuring concordance with the chosen strategy. Subsequently, the LLM formulates an explanation $\ell_t^{\text{expl}} \sim \pi_{\text{LLM}}(\ell_t^{\text{expl}} | \mathbf{a}_t^{\text{env}}, \mathbf{a}_t^{\text{rule}}, \mathbf{a}_t^{\text{thought}}, \mathbf{p}_t)$ contingent upon the rule.

This procedure concurrently produces the rule reward $R_{\text{LLM}}^{\text{rule}}(r_t^{\text{rule}})$, used for RL in the next step. This rewards is derived from using the LLM as a judge to answer the following three questions: ER₁. Without providing $\mathbf{a}_t^{\text{env}}$, is $\mathbf{a}_t^{\text{rule}}$ sufficient to predict the optimal action? ER₂. Does $\mathbf{a}_t^{\text{rule}}$ contain enough details about the applicability of the rule to current state? ER₃. Given $\mathbf{a}_t^{\text{env}}$, is $\mathbf{a}_t^{\text{rule}}$ compatible with this selection? Each question scores as 0 if negative or 1 if positive. The rule reward is calculated as $r_t^{\text{rule}} = R_{\text{LLM}}^{\text{rule}}(\mathbf{a}_t^{\text{rule}}) \propto (1/3) \sum_i \text{ER}_i$.

5.4 Policy Update through RL

Augmented state space Traditional RL frameworks fail to directly return a policy based on current environment state due to intermediate steps: generating the rule set \mathcal{R}_t , mapping rules $\mathbf{a}_t^{\text{rule}}$ to actions $\mathbf{a}_t^{\text{env}}$ in an LLM-driven environment. RBRL addresses this issue by creating an augmented state $\tilde{s}_t := (s_t, \mathcal{R}_t)$ with transition dynamics $P(\tilde{s}_{t+1} | \tilde{s}_t, \mathbf{a}_t^{\text{rule}})$, integrating rules into the state space for reasoning over both the environment’s dynamics and decision rules $\mathbf{a}_t^{\text{rule}}$. The following theorem explains the transition computation.

Theorem 5.1. *The state transition of the RBRL MDP can be calculated as*

$$P(\tilde{s}_{t+1} | \tilde{s}_t, \mathbf{a}_t^{\text{rule}}) = P(\mathcal{R}_{t+1} | s_{t+1}) \times \int_{\mathbf{a}} P(s_{t+1} | \mathbf{a}^{\text{env}}, s_t) \cdot P(\mathbf{a}^{\text{env}} | \mathbf{a}_t^{\text{rule}}, s_t) d\mathbf{a}^{\text{env}}, \quad (4)$$

where $P(\mathcal{R}_{t+1} | s_{t+1}) = \pi_{\text{LLM}}(\mathcal{R}_{t+1} | \mathbf{p}_t, \boldsymbol{\tau}_t)$ is the probability of the LLM generating rule set \mathcal{R}_{t+1} provided the state s_{t+1} , $P(s_{t+1} | \mathbf{a}^{\text{env}}, s_t)$ is the original environment dynamics, and $P(\mathbf{a}^{\text{env}} | \mathbf{a}_t^{\text{rule}}, s_t)$ is the probability of the LLM selecting the environment action \mathbf{a}^{env} , equivalent to $\pi_{\text{LLM}}(\mathbf{a}^{\text{env}} | \mathbf{p}_t, \mathbf{a}_t^{\text{rule}})$.

Policy update step The attention-based policy network in Section 5.2 is optimized using the standard SAC algorithm, which balances reward maximization with exploration. The policy network in SAC is updated by minimizing the KL divergence between the policy and the Boltzmann distribution induced by Q networks Q_{ϕ_i} , $\forall i = 1, 2$, which is expressed as

$$L_\pi(\theta) = \mathbb{E}_{\mathcal{D}} \left[\beta \log \pi_\theta(\mathbf{a}_t^{\text{rule}} | \tilde{s}_t) - \min_{i=1,2} Q_{\phi_i}(\tilde{s}_t, \mathbf{a}_t^{\text{rule}}) \right], \quad (5)$$

where β is a temperature parameter. The detailed implementation for SAC update procedure is detailed in Algorithm ?? in Appendix ??.

6 PERFORMANCE GUARANTEE

In this section, we derive and prove conditions under which RBRL can learn the optimal task policy, as well as characterize the potential trade-off between explainability and task performance when rewarding rules for higher explainability.

Proposition 6.1 (Rule Set Coverage). *Let \mathcal{A} be a finite action space and $Q^*(s, \mathbf{a}^{\text{env}})$ the optimal state-action value function, with $\mathbf{a}^{\text{env},*}(s) := \arg \max_{\mathbf{a}^{\text{env}} \in \mathcal{A}} Q^*(s, \mathbf{a}^{\text{env}})$ denoting the optimal action at*

state \mathbf{s} . Given state \mathbf{s} , an LLM samples N rules independently from a conditional distribution $\pi_{LLM}(\cdot | \mathbf{s})$, and each rule ρ_i maps \mathbf{s} to an action $\mathbf{a}_i^{env} \sim \pi_{LLM}(\mathbf{a}_i^{env} | \rho_i, \mathbf{s})$. Assume there exists $\delta > 0$ and $\eta_s \in (0, 1]$ such that: $\mathbb{P}_{\rho_i \sim \pi_{LLM}(\cdot | \mathbf{s})} [Q^*(\mathbf{s}, \mathbf{a}_i^{env}) \geq Q^*(\mathbf{s}, \mathbf{a}^{env,*}(\mathbf{s})) - \delta] \geq \eta_s$. Define the δ -optimal rule set as:

$$\mathcal{R}^\delta(\mathbf{s}) := \{\rho_i : Q^*(\mathbf{s}, \mathbf{a}_i^{env}) \geq Q^*(\mathbf{s}, \mathbf{a}^{env,*}(\mathbf{s})) - \delta\}.$$

Then with high probability over the sampled rules, there at least has a rule ρ_i and the induced action $\mathbf{a}_i^{env} \sim \pi_{LLM}(\mathbf{a}_i^{env} | \rho_i, \mathbf{s})$ that satisfies:

$$\mathbb{E} [Q^*(\mathbf{s}, \mathbf{a}^{env,*}(\mathbf{s})) - Q^*(\mathbf{s}, \mathbf{a}_i^{env})] \leq \delta + \epsilon_{worst} \cdot (1 - \eta_s)^N, \quad (6)$$

where $\epsilon_{worst} := \max_{\rho \notin \mathcal{R}^\delta(\mathbf{s})} (Q^*(\mathbf{s}, \mathbf{a}^{env,*}(\mathbf{s})) - Q^*(\mathbf{s}, \mathbf{a}_i^{env}))$ is the worst-case loss outside the δ -optimal set.

Remark 6.2. Proposition 6.1 states the *rule diversity* property in the rule candidate set such that the best possible action (when $\delta \rightarrow 0$) is included is guaranteed with high probability when number of rules N goes large. This is crucial in guaranteeing that RBRL can learn a near-optimal policy with high probability (with optimality when $\delta = 0$ and $\eta_s = 1$).

Define the T-step value function $V_{\mathcal{M}'}^{\pi, T}(\mathbf{s}_0) = [\sum_{t=0}^{T-1} \gamma^t R_t^{M'}(\mathbf{s}_t, \pi(\mathbf{s}_t))]_{\mathbf{s}_0}$, where $R^{M'}$ is the reward function in \mathcal{M}' . We will denote the original MDP as \mathcal{M} and use $\tilde{\mathcal{M}}$ to denote the MDP for the RBRL agent with transition function as in Theorem 5.1 and reward \tilde{R} . We have the following theorem.

Theorem 6.3. *The evaluation gap $Gap(T, \mathbf{s}_0) := V_{\mathcal{M}}^{\pi^*, T}(\mathbf{s}_0) - V_{\mathcal{M}}^{\pi_{RBRL}, T}(\mathbf{s}_0)$ of RBRL is bounded as*

$$Gap(T, \mathbf{s}_0) = V_{\mathcal{M}}^{\pi^*, T}(\mathbf{s}_0) - V_{\tilde{\mathcal{M}}}^{\pi_{RBRL}, T}(\mathbf{s}_0) + V_{\tilde{\mathcal{M}}}^{\pi_{RBRL}, T}(\mathbf{s}_0) - V_{\mathcal{M}}^{\pi_{RBRL}, T}(\mathbf{s}_0) \leq \lambda \cdot \frac{1 - \gamma^T}{1 - \gamma}, \quad (7)$$

where λ is a constant depending on the magnitude of the rule reward, and, with a slight notational abuse, $V_{\tilde{\mathcal{M}}}^{\pi_{RBRL}, T}$ is the value of the RBRL policy when seen as a policy in the original MDP mapping states to actions (i.e., by integrating out the rule generation and action selection via LLMs.)

Remark 6.4. This analysis focuses on the evaluation gap between the optimal policy π^* under the original MDP \mathcal{M} and the policy π_{RBRL} , captures the suboptimality of using π_{RBRL} instead of the true optimal policy π^* , assuming RBRL is optimized under the extended MDP $\tilde{\mathcal{M}}$ (with same transitions as \mathcal{M} but additional rule-based reward). It can be decomposed into two interpretable terms. The first part captures the optimism of using π_{RBRL} under the extended MDP $\tilde{\mathcal{M}}$ rather than the original MDP, which is non-positive. The second part quantifies the accumulated reward difference induced by the additional explanation rewards when using the same RBRL policy in both MDPs.

7 EXPERIMENTS & HUMAN SURVEY

In this section, we evaluate RBRL and empirically show that it can achieve a joint improvement in both reward and explainability over comparable baselines. We briefly summarize these environments here.

Domains. We evaluate RBRL in three main distinct resource-constrained allocation domains:

- **WearableDeviceAssignment:** We use two environments, Uganda and MimicIII, from the vital sign monitoring domain introduced by [5], modeling the allocation of limited wireless devices among postpartum mothers as an MDP setting.
- **HeatAlerts:** We use the `weather2alert` environment from [10], which formulates issuing heat alerts as a constrained MDP to reduce hospitalization risk from the alert.
- **BinPacking:** We adopt the online stochastic BinPacking: environment introduced by [2], which Sequentially places arriving items into bins with fixed capacity to minimize total waste, following the online stochastic formulation.

7.1 Environment Reward Optimization

We discuss the main results here. Unless otherwise specified, we use `gpt-4o-mini` as LLM due to its reasonable cost and high performance.

Q1. Did RBRL optimize the reward function? RBRL is compared to CoT [52] for language reasoning and PPO [36] for numeric states. Figure 4 indicates RBRL outperforms CoT, showing RL-optimized rule selection improves decision-making. RBRL also exceeds PPO in all environments with equal environment steps, suggesting a better online learning performance. Notice that RBRL is compatible with a baseline LLM trained for advanced reasoning techniques (e.g., GRPO [37]). However, GRPO or similar cannot be used directly in MDPs. Nevertheless, our experiments with the comparable GPT o3 prove that RBRL can also help improve reasoning models in our tasks.

Q2. Did structured rules help optimization? We conduct two ablation studies on structured rules. First, we benchmark the use of structured rules without RL, called baseline `Rule-bottleneck (no RL)`, which is shown in Equation (2)-(5). Next, we compare RBRL with a variant optimizing unstructured thoughts, termed `thoughts-based RL (TBRL)`. The implementation mimics RBRL, utilizing a candidate pool \mathcal{P} with the CoT prompt. Results in Figure 5a show that comparing RBRL with `RulesLLMOnly` highlights RL training gains, suggesting rule generation alone does not explain RBRL’s performance. Additionally, significant improvements over TBRL suggest optimizing structured rules is more effective than optimizing free reasoning traces.

Q3. How does RBRL compare to token-level LLM finetuning with RL? We implement LLM finetuning on a Llama 3.1 8B model, termed `FinetunePPO`. A value head is trained on final hidden states, with KL divergence from a reference model as regularization [63]. CoT is generated, followed by an action query, optimizing both. Training runs for 18 hours on 3 seeds using an A100 40G GPU (1200 steps/seed). For fair comparison, RBRL is also run on Llama 3.1 8B. Figure 5b shows RBRL outperforms the flatter trend of finetuning, indicating better online learning. Moreover, RBRL runs on a regular laptop, whereas `FinetunePPO` requires specialized hardware and takes 4× longer per step. Due to compute limits, results are shown only for the less noisy `MimicII` domain.

Additional comparison with XRL benchmarks. We further compare RBRL against a representative XRL method that also targets joint optimization and intrinsic interpretability: Decision Diffusion

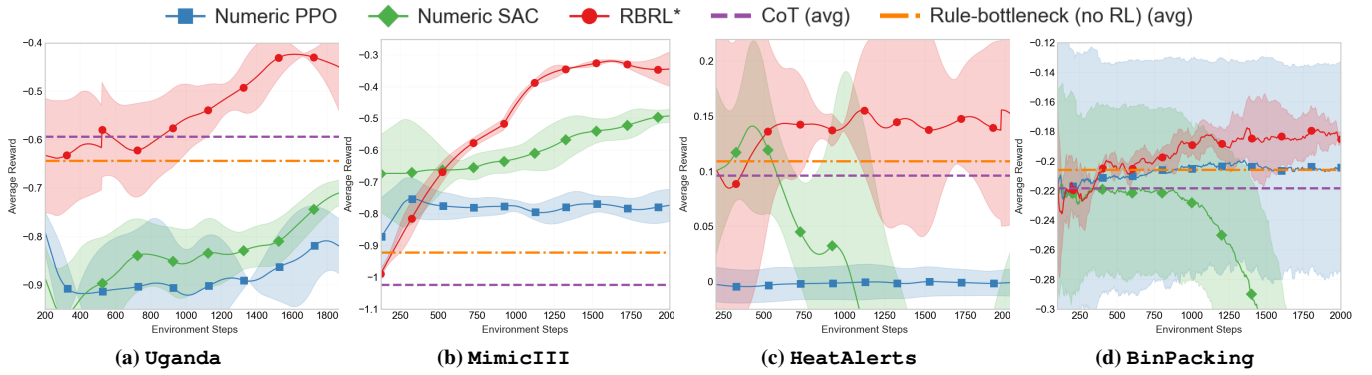


Figure 4: Results from Q1 using ChatGPT 4o-mini. The plots show the mean and standard error across three seeds, using exponentially weighted moving averages ($\lambda_{ema} = 200$).

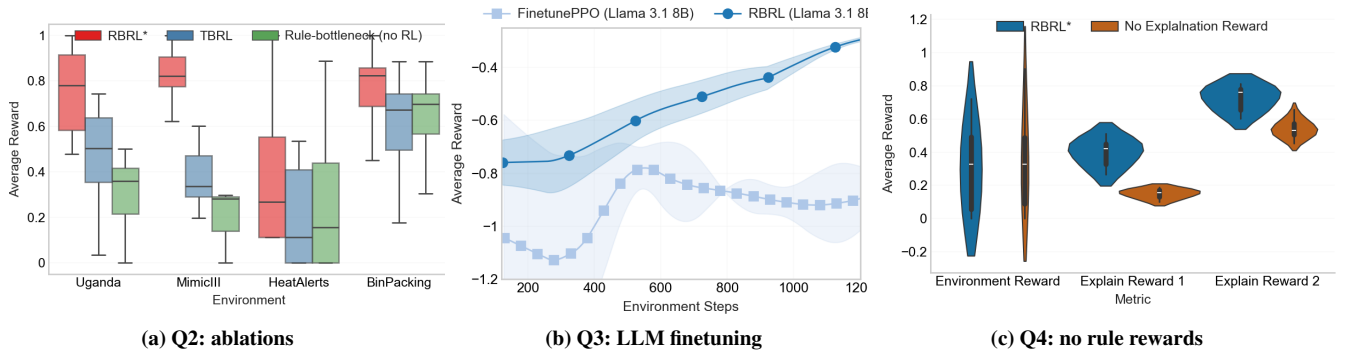


Figure 5: Additional experiments and ablations. (a) Comparison of RBRL with thoughts-based RL (TBRL) and the baseline rule-based LLM without RL training; (b) comparison against LLM finetuning with PPO at the token level from the environment reward with CoT generation for the *Mimic*; (c) shows the effect of removing the rule reward in the *HeatAlerts* environments. For (a) and (c), we show distribution of rewards in the last 20% training steps.

Table 1: XRL Baselines Results Table

Dataset (@steps)	RBRL	SAC	PPO	DDT	DDT w/rules
Uganda (@500)	-0.56 ± 0.18	-0.83 ± 0.14	-0.91 ± 0.14	-1.01 ± 0.20	-1.20 ± 0.31
Uganda (@2500)	-0.60 ± 0.20	-0.75 ± 0.14	-0.74 ± 0.05	-1.28 ± 0.35	-1.20 ± 0.30
MimicIII (@500)	-0.36 ± 0.05	-0.61 ± 0.11	-0.78 ± 0.05	-0.92 ± 0.10	-1.02 ± 0.10
MimicIII (@2500)	-0.39 ± 0.07	-0.43 ± 0.10	-0.64 ± 0.10	-0.97 ± 0.11	-0.99 ± 0.13
HeatAlerts (@500)	0.14 ± 0.11	-0.04 ± 0.33	0.00 ± 0.01	0.22 ± 0.25	0.15 ± 0.29
HeatAlerts (@2500)	0.13 ± 0.14	0.05 ± 0.04	0.00 ± 0.01	0.38 ± 0.57	0.38 ± 0.56
BinPacking (@500)	-0.03 ± 0.00	-0.03 ± 0.00	-0.03 ± 0.00	-0.19 ± 0.03	-0.19 ± 0.04
BinPacking (@2500)	-0.03 ± 0.00	-0.06 ± 0.00	-0.03 ± 0.00	-0.21 ± 0.02	-0.21 ± 0.02

Trees (DDTs) [41]. As shown in Table 1, RBRL is consistently competitive and often outperforms the tree-based baseline across most domains, particularly in the early stages of training, underscoring its sample efficiency. Although DDT achieves a higher average reward than RBRL in *HeatAlerts*, it exhibits substantially higher variance, highlighting the greater stability of RBRL.

7.2 Human Survey and Explainability

Q4. Did RBRL increase the explainability of explanations? A survey with 40 participants was conducted to assess explanation quality. Each prompt included the task, state, and action space as originally given to the LLM, followed by actions and explanations from

the CoT agent and the RBRL agent, without disclosing agent types. Participants were asked to choose preference for explanation A, B, or none. Prompts were split between *WearableDeviceAssignment* and *HeatAlerts* domains. Figure 6 shows results, favoring RBRL’s explanations in both domains. An additional experiment with an LLM judge using a large *gpt-4o* model showed strong agreement with humans, preferring RBRL’s explanations in all cases.

Discussion on Explainability. The trustworthiness of explanations is a core challenge in XAI. Following recent work [23, 24, 28], we highlight three concepts: *Plausibility*: whether an explanation is convincing to humans (validated via our survey, Figure 6). *Consistency*: whether the stated reason logically entails the action. *Faithfulness*: whether the explanation reflects the true decision mechanism.

Our work is motivated by the gap between plausibility and faithfulness in post-hoc methods. By design, RBRL ensures consistency: explanations factually follow the State → Rule → Action pipeline, where the rule is the verifiable cause of the action. While our experiments validate consistency, establishing faithfulness—verifying the LLM’s internal reasoning for rule generation—remains an open challenge.

Q5. What was the effect of the rule reward? During training of RBRL, rules received rewards from two prompts. We examine an

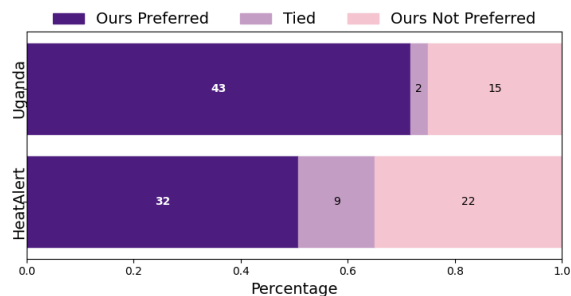


Figure 6: Results from the human survey.

ablation without this reward. Figure 5c illustrates results for the HeatAlerts environment, noted for high variance and a challenging reward function. We extended training to 5k steps to understand these dynamics. Without rule reward, environment reward remains steady (slightly increasing), but explainability scores drop significantly. Refer to Section 5.3 for the definition of the rule reward metrics. A decline in metric 1 indicates that rules are less predictive of the optimal actions. A decline in metric 2 suggests rules lack detailed applicability to the current problem state, indicating more generic rather than specialized rule selection. Metric 3 (not shown) was always 1 in all steps, indicating the limitations of directly evaluating post hoc explanations. Although judged by the LLM, these results are encouraging, as our previous experiment showed alignment between the LLM and human assessments.

ETHICS STATEMENT

The authors of this work adhere to the AAMAS Code of Ethics. Our research involves domains with significant ethical considerations, particularly in healthcare and public policy, which we have carefully addressed. Our work includes a human survey to evaluate the quality of explanations, which was conducted under Institutional Review Board (IRB) approval.

REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *ICLR*.
- [2] Bharathan Balaji, Jordan Bell-Masterson, Enes Bilgin, Andreas Damianou, Pablo Moreno Garcia, Arpit Jain, Runfei Luo, Alvaro Maggiar, Balakrishnan Narayanaswamy, and Chun Ye. 2019. Orl: Reinforcement learning benchmarks for online stochastic optimization problems. *arXiv preprint arXiv:1911.10641* (2019).
- [3] Amrita Bhattacharjee, Raha Moraffah, Joshua Garland, and Huan Liu. 2024. Towards llm-guided causal explainability for black-box text classifiers. In *AAAI 2024 Workshop on Responsible Language Models*.
- [4] Adeline A Boatin, Joseph Ngonzi, Blair J Wylie, Henry M Lugobe, Lisa M Bebell, Godfrey Mugenyi, Sudi Mohamed, Kenia Martinez, Nicholas Musinguzi, Christina Psaros, et al. 2021. Wireless versus routine physiologic monitoring after cesarean delivery to reduce maternal morbidity and mortality in a resource-limited setting: protocol of type 2 hybrid effectiveness-implementation study. *BMC Pregnancy and Childbirth* 21 (2021), 1–12.
- [5] Niclas Boehmer, Yunfan Zhao, Guojun Xiong, Paula Rodriguez-Diaz, Paola Del Cueto Cibrian, Joseph Ngonzi, Adeline Boatin, and Milind Tambe. 2024. Optimizing vital sign monitoring in resource-constrained maternal care: An rl-based restless bandit approach. *arXiv preprint arXiv:2410.08377* (2024).
- [6] Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer. 2023. Grounding large language models in interactive environments with online reinforcement learning. In *ICML*.
- [7] Haozhe Chen, Carl Vondrick, and Chengzhi Mao. 2024. Selfie: Self-interpretation of large language model embeddings. *arXiv preprint arXiv:2403.10949* (2024).
- [8] Zelei Cheng, Xian Wu, Jiahao Yu, Sabrina Yang, Gang Wang, and Xinyu Xing. 2024. Rice: Breaking through the training bottlenecks of reinforcement learning

with explanation. *arXiv preprint arXiv:2405.03064* (2024).

- [9] Cédric Colas, Tristan Karch, Clément Moulin-Frier, and Pierre-Yves Oudeyer. 2022. Language and culture internalization for human-like autotelic AI. *Nature Machine Intelligence* 4, 12 (2022), 1068–1076.
- [10] Ellen M Considine, Rachel C Nethery, Gregory A Wellenius, Francesca Dominici, and Mauricio Tec. 2025. Optimizing Heat Alert Issuance with Reinforcement Learning. In *AAAI*.
- [11] CSL. 2024. Senate Bill No. 896: Generative Artificial Intelligence Accountability Act. <https://legiscan.com/CA/text/SB896/id/3023382>. Accessed: 2025-02-06.
- [12] Devleena Das, Sonia Chernova, and Been Kim. 2023. State2Explanation: Concept-Based Explanations to Benefit Agent Learning and User Understanding. *NeurIPS* (2023).
- [13] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).
- [14] Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. 2023. Guiding pretraining in reinforcement learning with large language models. In *ICML*. 8657–8677.
- [15] Kelly N DuBois. 2019. Deep medicine: How artificial intelligence can make healthcare human again. *Perspectives on Science and Christian Faith* 71, 3 (2019), 199–201.
- [16] EPC. 2024. Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence. <http://data.europa.eu/eli/reg/2024/1689/oj>. Accessed: 2025-02-06.
- [17] Hiroki Furuta, Yutaka Matsuo, Aleksandra Faust, and Izzeddin Gur. 2024. Exposing Limitations of Language Model Agents in Sequential-Task Compositions on the Web. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.
- [18] Georges Gebrael, Kamal Kant Sahu, Beverly Chigarira, Nishita Tripathi, Vinay Mathew Thomas, Nicolas Sayegh, Benjamin L Maughan, Neeraj Agarwal, Umang Swami, and Haoran Li. 2023. Enhancing triage efficiency and accuracy in emergency rooms for patients with metastatic prostate cancer: a retrospective analysis of artificial intelligence-assisted triage using ChatGPT 4.0. *Cancers* 15, 14 (2023), 3717.
- [19] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. 2024. A Survey on LLM-as-a-Judge. *arXiv preprint arXiv:2411.15594* (2024).
- [20] Wenbo Guo, Xian Wu, Usman Khan, and Xinyu Xing. 2021. Edge: Explaining deep reinforcement learning policies. *Advances in Neural Information Processing Systems* 34 (2021), 12222–12236.
- [21] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 1856–1865.
- [22] Hengyuan Hu and Dorsa Sadigh. 2023. Language instructed reinforcement learning for human-AI coordination. In *ICML*.
- [23] Alon Jacovi and Yoav Goldberg. 2020. Towards faithfully interpretable NLP systems: How should we define and evaluate faithfulness? *arXiv preprint arXiv:2004.03685* (2020).
- [24] Jenny Kunz and Marco Kuhlmann. 2024. Properties and challenges of llm-generated explanations. *arXiv preprint arXiv:2402.10532* (2024).
- [25] Meta AI. 2024. Introducing Llama 3.1: Our most capable models to date. *AI at Meta* (2024). <https://ai.meta.com/blog/meta-llama-3-1/>
- [26] Stephanie Milani, Nicholay Topin, Manuela Veloso, and Fei Fang. 2024. Explainable reinforcement learning: A survey and comparative review. *Comput. Surveys* 56, 7 (2024), 1–36.
- [27] OpenAI. 2024. GPT-4o Mini: Advancing Cost-Efficient Intelligence. *OpenAI Blog* (2024).
- [28] Letitia Parcalabescu and Anette Frank. 2023. On measuring faithfulness or self-consistency of natural language explanations. *arXiv preprint arXiv:2311.07466* (2023).
- [29] Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*. 1–22.
- [30] Xiangyu Peng, Mark Riedl, and Prithviraj Ammanabrolu. 2022. Inherently Explainable Reinforcement Learning in Natural Language. In *NeurIPS*, Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (Eds.).
- [31] Eleonora Poeta, Gabriele Ciravegna, Eliana Pastor, Tania Cerquitelli, and Elena Baralis. 2023. Concept-based explainable artificial intelligence: A survey. *arXiv preprint arXiv:2312.12936* (2023).
- [32] Pranav Putta, Edmund Mills, Naman Garg, Sumeet Motwani, Chelsea Finn, Divyansh Garg, and Rafael Rafailov. 2024. Agent q: Advanced reasoning and learning for autonomous ai agents. *arXiv preprint arXiv:2408.07199* (2024).
- [33] Ahmad Rashid, Ruotian Wu, Julia Grosse, Agustinus Kristiadi, and Pascal Poupart. 2024. A Critical Look At Tokenwise Reward-Guided Text Generation. *arXiv preprint arXiv:2406.07780* (2024).
- [34] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on EMNLP-IJCNLP*. Association for Computational Linguistics, 3980–3990.

- [35] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence* 1, 5 (2019), 206–215.
- [36] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).
- [37] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300* (2024).
- [38] Yiran Shen, Aditya Emmanuel Arokiaraj John, and Brandon Fain. 2024. Explainable Rewards in RLHF Using LLM-as-a-Judge. <https://openreview.net/forum?id=FaOeBrlPst>
- [39] Yi Shen, Benjamin McClosky, and Michael Zavlanos. 2023. Multi-agent reinforcement learning for resource allocation in large-scale robotic warehouse sortation centers. (2023).
- [40] Daria Shevtsova, Anam Ahmed, Iris WA Boot, Carmen Sanges, Michael Hudecek, John JL Jacobs, Simon Hort, Hubertus JM Vrijhoef, et al. 2024. Trust in and Acceptance of Artificial Intelligence Applications in Medicine: Mixed Methods Study. *JMIR Human Factors* 11, 1 (2024), e47031.
- [41] Andrew Silva, Matthew Gombolay, Taylor Killian, Ivan Jimenez, and Sung-Hyun Son. 2020. Optimization methods for interpretable differentiable decision trees applied to reinforcement learning. In *International conference on artificial intelligence and statistics*. PMLR, 1855–1865.
- [42] Sean R. Sinclair, Felipe Vieira Frujeri, Ching-An Cheng, Luke Marshall, Hugo Barbalho, Jingling Li, Jennifer Neville, Ishai Menache, and Adith Swaminathan. 2023. Hindsight Learning for MDPs with Exogenous Inputs. In *ICML*.
- [43] Chandan Singh, Jeevana Priya Inala, Michel Galley, Rich Caruana, and Jianfeng Gao. 2024. Rethinking interpretability in the era of large language models. *arXiv preprint arXiv:2402.01761* (2024).
- [44] Megha Srivastava, Cédric Colas, Dorsa Sadigh, and Jacob Andreas. 2024. Policy learning with a language bottleneck. In *RLC Workshop on Training Agents with Foundation Models*.
- [45] Theodore Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas Griffiths. 2024. Cognitive Architectures for Language Agents. *TMLR* (2024). <https://openreview.net/forum?id=Ii6ZCvflQJ>
- [46] Andrew Szot, Max Schwarzer, Harsh Agrawal, Bogdan Mazouze, Rin Metcalf, Walter Talbott, Natalie Mackrath, R Devon Hjelm, and Alexander T Toshev. 2023. Large language models as generalizable policies for embodied tasks. In *ICLR*.
- [47] Fatma M Talaat. 2022. Effective deep Q-networks (EDQN) strategy for resource allocation based on optimized reinforcement learning algorithm. *Multimedia Tools and Applications* 81, 28 (2022), 39945–39961.
- [48] Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. 2024. True Knowledge Comes from Practice: Aligning Large Language Models with Embodied Environments via Reinforcement Learning. In *ICLR*.
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Kukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, Vol. 30.
- [50] Alexander Sasha Vezhnevets, John P Agapiou, Avia Aharon, Ron Ziv, Jayd Matyas, Edgar A Duéñez-Guzmán, William A Cunningham, Simon Osindero, Danny Karmon, and Joel Z Leibo. 2023. Generative agent-based modeling with actions grounded in physical, social, or digital space using Concordia. *arXiv preprint arXiv:2312.03664* (2023).
- [51] Zihao Wang, Shaofei Cai, Guanzhou Chen, Anji Liu, Xiaojian Ma, Yitao Liang, and Team CraftJarvis. 2023. Describe, explain, plan and select: interactive planning with large language models enables open-world multi-task agents. In *NeurIPS*.
- [52] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [53] Muning Wen, Ziyu Wan, Jun Wang, Weinan Zhang, and Ying Wen. 2024. Reinforcing LLM Agents via Policy Optimization with Action Decomposition. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- [54] Peter Whittle. 1988. Restless bandits: Activity allocation in a changing world. *Journal of applied probability* 25, A (1988), 287–298.
- [55] Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, et al. 2024. AgentGym: Evolving Large Language Model-based Agents across Diverse Environments. *arXiv preprint arXiv:2406.04151* (2024).
- [56] Guojun Xiong, Xudong Qin, Bin Li, Rahul Singh, and Jian Li. 2022. Index-aware reinforcement learning for adaptive video streaming at the wireless edge. In *Proceedings of the Twenty-Third International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing*. 81–90.
- [57] Guojun Xiong, Shufan Wang, Gang Yan, and Jian Li. 2023. Reinforcement learning for dynamic dimensioning of cloud caches: A restless bandit approach. *IEEE/ACM Transactions on Networking* 31, 5 (2023), 2147–2161.
- [58] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. 2023. ReAct: Synergizing Reasoning and Acting in Language Models. In *ICLR*.
- [59] Chao Yu, Jiming Liu, Shamim Nemati, and Guosheng Yin. 2021. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)* 55, 1 (2021), 1–36.
- [60] Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yupeng Cao, Zhi Chen, Jordan W Suchow, Rong Liu, Zhenyu Cui, Zhaozhuo Xu, et al. 2024. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *arXiv preprint arXiv:2407.06567* (2024).
- [61] Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr, Saining Xie, Yann LeCun, Yi Ma, et al. 2024. Fine-Tuning Large Vision-Language Models as Decision-Making Agents via Reinforcement Learning. *arXiv preprint arXiv:2405.10292* (2024).
- [62] Xijia Zhang, Yue Guo, Simon Stepputtis, Katia Sycara, and Joseph Campbell. 2023. Understanding Your Agent: Leveraging Large Language Models for Behavior Explanation. *arXiv preprint arXiv:2311.18062* (2023).
- [63] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593* (2019).