BEYOND SIMPLE SUM OF DELAYED REWARDS: NON-MARKOVIAN REWARD MODELING FOR REIN FORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement Learning (RL) empowers agents to acquire various skills by learning from reward signals. Unfortunately, designing high-quality instance-level rewards often demands significant effort. An emerging alternative, RL with delayed reward, focuses on learning from rewards presented periodically, which can be obtained from human evaluators assessing the agent's performance over sequences of behaviors. However, traditional methods in this domain assume the existence of underlying Markovian rewards and that the observed delayed reward is simply the sum of instance-level rewards, both of which often do not align well with realworld scenarios. In this paper, we introduce the problem of RL from Composite Delayed Reward (RLCoDe), which generalizes traditional RL from delayed rewards by eliminating the strong assumption. We suggest that the delayed reward may arise from a more complex structure reflecting the overall contribution of the sequence. To address this problem, we present a framework for modeling composite delayed rewards, using a weighted sum of non-Markovian components to capture the different contributions of individual steps. Building on this framework, we propose Composite Delayed Reward Transformer (CoDeTr), which incorporates a specialized in-sequence attention mechanism to effectively model these contributions. We conduct experiments on challenging locomotion tasks where the agent receives delayed rewards computed from composite functions of observable step rewards. The experimental results indicate that CoDeTr consistently outperforms baseline methods across evaluated metrics. Additionally, we demonstrate that it effectively identifies the most significant time steps within the sequence and accurately predicts rewards that closely reflect the environment feedback. Code is available at an anonymous link: https://anonymous.4open.science/r/CoDe-67E8/.

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

033

1 INTRODUCTION

Reinforcement Learning (RL) has become a powerful paradigm for solving sequential decisionmaking problems by enabling agents to learn optimal policies through interactions with their environments (Kaelbling et al., 1996; Arulkumaran et al., 2017). A critical component of RL is the reward function, which guides the learning process by indicating the desirability of different states and actions. In complex real-world applications such as autonomous driving (Sallab et al., 2017; Kiran et al., 2021), financial trading (Yang et al., 2020; Hambly et al., 2023), and healthcare (Coronato et al., 2020; Yu et al., 2021), designing detailed reward signals for every possible state-action pair is often impractical due to the high dimensionality and complexity of these domains.

To address the challenge of specifying fine-grained rewards, researchers have explored the use of delayed rewards. Instead of assigning rewards at every individual step or action, feedback is given based on the outcome of a sequence of actions (Liu et al., 2019; Gangwani et al., 2020; Efroni et al., 2021; Raposo et al., 2021; Ren et al., 2021). These works generally aim to simplify the reward design process while still allowing for effective policy learning. For example, Liu et al. (2019) leveraged sequence modeling with expert demonstrations, while Gangwani et al. (2020) and Ren et al. (2021) proposed iterative and randomized approaches to refine reward credit assignment. While these approaches reduce the burden of reward engineering, they typically rely on two key assumptions. First, they assume the existence of underlying *Markovian* rewards, meaning that the



Figure 1: Illustration of our framework. The Composite Delayed Reward Transformer generates the predicted non-Markovian rewards \hat{r}_t along with the corresponding importance weights w_t for each sequence. The final composite reward for each sequence is calculated as a weighted sum of the predicted rewards.

068

069

054

056

060 061

062

063

reward at any step depends solely on the current state and action. Second, they posit that the structure of the delayed reward is a *simple sum* of individual rewards associated with each state-action pair, with equal weighting across the sequence.

However, these assumptions may not hold in many practical scenarios. The Markovian assumption 071 neglects the fact that rewards can depend on sequences of states or actions that are not fully captured by the current state (Bacchus et al., 1996; 1997). Furthermore, the assumption of equal-weighted 073 summation does not align with how humans evaluate experiences. Psychological studies have shown 074 that people tend to assign disproportionate importance to remarkable or significant moments within 075 an experience (Kahneman, 2000; Newell et al., 2022). This suggests that in tasks involving human feedback, certain states or actions within a sequence may contribute more heavily to the overall 076 assessment than others. An example of this phenomenon is observed in high-stake environments 077 studied by Klein (2008): experts such as firefighters, pilots, and emergency medical personnel often focus intensely on critical cues and pivotal moments that can significantly affect outcomes. These 079 professionals rely on recognizing patterns and key indicators to make rapid decisions, effectively assigning greater weight to crucial information rather than treating all information equally. There-081 fore, previous methods that assume equal-weighted summation of Markovian rewards may fail to 082 capture the true nature of the feedback. As a result, these methods may not yield promising results 083 in scenarios where the reward structure is inherently non-Markovian and where critical moments 084 disproportionately influence the overall evaluation. 085

Building on the above observations, we emphasize the need for RL frameworks that can address the RL from Composite Delayed Reward (RLCoDe) problem by accommodating non-Markovian rewards and capturing the disproportionate weighting inherent in delayed feedback. To meet this need, we propose the Composite Delayed Reward Transformer (CoDeTr), which introduces a non-Markovian reward model that predicts rewards more accurately reflecting the underlying reward structure as perceived by humans. Additionally, CoDeTr incorporates an *in-sequence attention mechanism* to model the reward aggregation process, thereby capturing critical instances within a sequence and reflecting the varying importance that human feedback assigns to these moments. By integrating these non-Markovian rewards into the policy learning process, the agent can learn more effectively in environments where traditional assumptions about rewards do not hold.

- ⁵ Our main contributions can be summarized as follows:
- 096

098

099

102

103

- We identify the limitations of existing delayed reward frameworks relying on Markovian assumption and equal-weighted summation, and proposed RLCoDe to address these issues.
- We propose CoDeTr to accommodate non-Markovian reward functions by capturing the disproportionate weight of feedback through an in-sequence attention mechanism.
- We demonstrate that our approach outperforms state-of-the-art delayed reward methods in environments where rewards depend on the sequence of visited states and where critical moments have a greater impact on the overall evaluation.
- We verify that our method effectively learns the rewards corresponding to the agent's actions, while also identifying which specific steps within the sequence are given more emphasis by the composite delayed reward.

108 2 PRELIMINARIES 109

110 In this section, we revisit conventional RL and RL from delayed rewards under the assumption 111 of Markovian rewards. Building on these concepts, we introduce the problem of RL from Com-112 posite Delayed Reward (RLCoDe), which generalizes these settings by removing the Markovian 113 assumption and allowing for non-Markovian reward structures with flexible, non-uniform weighting 114 of contributions across the sequence.

116 2.1 STANDARD REINFORCEMENT LEARNING

In standard RL (Bellman, 1966), the environment is modeled as a Markov Decision Process (MDP):

Definition 1 (MDP) An MDP is defined by the tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P, r, \mu)$, where

- *S* is a finite set of states.
- *A is a finite set of actions.*
- $P: S \times A \times S \rightarrow [0,1]$ is the state transition probability function, where $P(\mathbf{s}'|\mathbf{s}, \mathbf{a})$ gives the probability of transitioning to state s' from state s after taking action a.
- $r: S \times A \rightarrow \mathbb{R}$ is the immediate reward function.
- *μ* is the initial state distribution.

130 An agent interacts with the environment by following a policy $\pi: \mathcal{S} \times \mathcal{A} \to [0, 1]$, where $\pi(\mathbf{a}|\mathbf{s})$ is the probability of taking action a in state s. The objective is to find an optimal policy π^* that maximizes the expected cumulative reward: 132

$$I(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} r(\mathbf{s}_t, \mathbf{a}_t) \right],\tag{1}$$

where the expectation is over trajectories induced by the policy π and the transition dynamics P.

2.2 REINFORCEMENT LEARNING FROM DELAYED REWARDS

J

141 In many real-world applications, immediate rewards are not readily available or are difficult to spec-142 ify (Devidze et al., 2022; Tang et al., 2024). Instead, the agent may receive delayed rewards, often at the end of a sequence or trajectory. This setting is common in domains where the outcome of ac-143 tions is not immediately observable or when feedback is provided by human evaluators who assess 144 performance over extended periods (Shen & Chi, 2016; Krishnan et al., 2019; Gao et al., 2024). 145

146 In RL with delayed rewards, the environment is modeled as an MDP with a cumulative reward 147 $R(\tau)$ provided for a sequence τ or trajectory \mathcal{T} . Let $\mathcal{T} = \{(\mathbf{s}_0, \mathbf{a}_0), (\mathbf{s}_1, \mathbf{a}_1), \dots, (\mathbf{s}_{T-1}, \mathbf{a}_{T-1})\}$ represent an agent trajectory over T time steps, encompassing all experienced states and actions 148 within an episode. A sequence τ refers to a portion of the trajectory \mathcal{T} , starting at time step i and 149 consisting of n_i state-action pairs: 150

153

154

155 156 157

158

159 160

$$\tau = \{ (\mathbf{s}_i, \mathbf{a}_i), (\mathbf{s}_{i+1}, \mathbf{a}_{i+1}), \dots, (\mathbf{s}_{i+n_i-1}, \mathbf{a}_{i+n_i-1}) \}.$$

A common assumption is that there exists an underlying Markovian reward function $r(\mathbf{s}, \mathbf{a})$ such that the sequence-level reward can be decomposed as

- $R(\tau) = \sum_{t=i}^{i+n_i-1} r(\mathbf{s}_t, \mathbf{a}_t).$ (2)
- This assumption simplifies the problem by allowing standard RL algorithms to be applied after 161 redistributing the cumulative reward $R(\tau)$ back to individual time steps. However, this approach

119 120

115

117

118

121

122

123 124

125

126

127 128

129

- 131
- 133 134 135
- 136 137

138 139

may not be suitable in situations where the reward at each time step depends on the context of the entire sequence, or when the delayed reward places greater emphasis on certain critical time steps.

2.3 REINFORCEMENT LEARNING FROM COMPOSITE DELAYED REWARD

To address the limitations of assuming Markovian and additive reward structure, we introduce the Composite Delayed Reward Markov Decision Process (CoDeMDP), which generalizes the conventional MDP to allow for non-Markovian rewards and non-uniform weighting of contributions.

Definition 2 (CoDeMDP) A CoDeMDP is defined by the tuple (S, A, P, R_{co}, μ) , where

- *S* is the set of states.
- *A is the set of actions.*
- $P: S \times A \times S \rightarrow [0, 1]$ is the state transition probability function.
- $R_{co} : \tau \to \mathbb{R}$ is the composite delayed reward function, defined on sequences τ within a trajectory \mathcal{T} .
- μ is the initial state distribution.

This composite delayed reward is assigned periodically, based on the sequence of states and actions between two reward points, similar to traditional delayed reward settings (Gangwani et al., 2020; Ren et al., 2021). In this framework, R_{co} is a reward function of the entire sequence τ , which may depend on complex interactions among states and actions, without assuming Markovian properties for instance-level rewards or requiring the sequence-level reward to be additive by instance-level rewards. This generalization allows us to model scenarios where the reward depends on the sequence history, future outcomes, or non-linear aggregation of individual contributions. The agent's objective is to maximize the expected cumulative composite reward over trajectories:

$$J(\pi) = \mathbb{E}_{\mathcal{T}(\pi)} \left[\sum_{\tau \subseteq \mathcal{T}} R_{\rm co}(\tau) \right],\tag{3}$$

where $\mathcal{T}(\pi)$ denotes the distribution over trajectories induced by policy π .

3 PROPOSED METHOD

In this section, we introduce our proposed method to address the challenges of RLCoDe. We begin by presenting a general framework for modeling delayed rewards using a weighted sum of non-Markovian components. Following this, we propose a transformer-based architecture called CoDeTr, specifically designed to predict instance-level non-Markovian rewards and integrate them into a sequence-level composite delayed reward.

165

166 167

168

169

170 171

172

173 174

175 176

177

178

179

181

196 197

198 199

200

201

202

203

3.1 SEQUENCE-LEVEL REWARD DECOMPOSITION

207 As discussed above, most prior work assumes that the instance-level reward is Markovian and that 208 sequence-level feedback is based on an equal-weighted sum of immediate rewards. This assump-209 tion may not hold in many real-world scenarios if the reward depends on the sequence of states 210 and actions, including historical context and future outcomes, rather than solely on the current state 211 and action (Bacchus et al., 1996; 1997; Early et al., 2022). Meanwhile, some certain states or ac-212 tions within a sequence have a disproportionate impact on the overall reward, reflecting the human 213 tendency to assign greater importance to critical moments (Kahneman, 2000; Newell et al., 2022). These limitations necessitate a framework that can capture complex reward dependencies and vary-214 ing importance of different parts of a trajectory. To address these challenges, we propose modeling 215 the sequence-level reward using a non-Markovian reward function \hat{r} with learnable weights w:

²⁰⁴ 205 206



Figure 2: The architecture of the proposed Composite Delayed Reward Transformer. The model processes the sequence of state-action pairs using a causal transformer, where the embeddings xrepresent the context information from the initial time step to current time step. The in-sequence attention mechanism computes the non-Markovian rewards $\{\hat{r}\}_{\tau}$, the queries $\{\mathbf{q}\}_{\tau}$, and keys $\{\mathbf{k}\}_{\tau}$, in a sequence τ . The query and key vectors are multiplied and passed through a softmax operation to compute attention weights. The attention-weighted sum of instance-level rewards is then aggregated via sum pooling to generate the final sequence-level reward $\hat{R}_{co}(\tau)$ to approximate $R_{co}(\tau)$.

$$\hat{R}_{\rm co}(\tau) = \sum_{t=i}^{i+n_i-1} w \big(\{ (\mathbf{s}_{t'}, \mathbf{a}_{t'}) \}_{t'=i}^{i+n_i-1}; \psi \big)_t \cdot \hat{r} \big(\{ (\mathbf{s}_{t'}, \mathbf{a}_{t'}) \}_{t'=i}^t; \psi \big)_t, \tag{4}$$

where ψ represents the learnable parameters for both the weight function w and the reward function \hat{r} , ensuring they can be jointly optimized for more effective learning of the reward structure. Unlike traditional Markovian rewards that depend only on the current state and action, \hat{r} considers the historical sequence of state-action pairs $\{(\mathbf{s}_{t'}, \mathbf{a}_{t'})\}_{t'=i}^{t}$ up to the current time step t. The weight function w takes the entire sequence $\tau = \{(\mathbf{s}_{t'}, \mathbf{a}_{t'})\}_{t'=i}^{t}$ as input, allowing the model to evaluate the importance of each time step in the context of the whole sequence. This formulation defines how the sequence-level reward is related to instance-level contributions and their associated weights.

247 248 249

265 266 267

228

229

230

231

232

233

234 235

237 238 239

3.2 ARCHITECTURE

To implement the proposed approach, we design a transformer-based architecture called Composite
 Delayed Reward Transformer (CoDeTr), which functions as a reward model for predicting instance level non-Markovian rewards and representing the composite delayed reward.

253 **Instance-Level Reward Prediction.** We adopt the transformer network (Vaswani et al., 2017) 254 as the backbone for instance-level reward prediction. Specifically, we employ the GPT architec-255 ture (Radford et al., 2018), which utilizes a causal self-attention mechanism. This causal trans-256 former ensures the chronological order of state-action pairs is preserved in our proposed reward 257 model. For each time step t in a sequence of M time steps, the causal transformer, represented as 258 a function g, processes the input sequence $\sigma = \{(\mathbf{s}_0, \mathbf{a}_0), \dots, (\mathbf{s}_{M-1}, \mathbf{a}_{M-1})\}$, generating outputs 259 $\{\mathbf{x}_t\}_{t=0}^{M-1} = g(\sigma)$. By aligning the output \mathbf{x}_t with the action token \mathbf{a}_t , we directly model the con-260 sequences of actions, which are pivotal in computing immediate rewards and predicting subsequent 261 states, thereby helping the model better understand environmental dynamics.

After obtaining the output embeddings $\{\mathbf{x}_t\}_{t=0}^{M-1}$ from the causal transformer, a linear transformation is applied to each \mathbf{x}_t to compute the corresponding instance-level reward \hat{r}_t :

$$\hat{r}_t = \text{Linear}(\mathbf{x}_t). \tag{5}$$

Since the underlying instance-level reward depends on the previous state-action pairs, this allows us
 to model the immediate rewards associated with each action in the sequence, capturing the essential dynamics at the instance level.

Composite Delayed Reward Representation. To represent the composite delayed reward based on instance-level rewards, we use an in-sequence attention mechanism inspired by traditional attention models (Vaswani et al., 2017). The in-sequence attention mechanism operates only within the sequence τ corresponding to the composite delayed reward, ensuring that the attention mechanism focuses solely on the instance-level rewards that are relevant to the specific sequence. This attention is bi-directional, allowing the model to consider the entire sequence when determining the importance of each time step.

277 Specifically, we apply two linear transformations to each embedding \mathbf{x}_t , resulting in query embed-278 dings $\mathbf{q}_t \in \mathbb{R}^d$ and key embeddings $\mathbf{k}_t \in \mathbb{R}^d$, where *d* is the embedding dimension. The attention 279 weight for each time step is calculated using the dot product between the query and key embeddings, 280 followed by a softmax operation to normalize the weights:

281 282 283

284 285 286

287

294

295

296

297 298

299

$$\hat{R}_{\rm co}(\tau) = \sum_{i \in \tau} \sum_{t \in \tau} \operatorname{softmax}\left(\frac{\{\langle \mathbf{q}_i, \mathbf{k}_{t'} \rangle\}_{t' \in \tau}}{\sqrt{d}}\right) \hat{r}_t,\tag{6}$$

where $\langle \cdot, \cdot \rangle$ denotes the dot product, and the scaling factor \sqrt{d} is used to prevent extremely small gradients (Vaswani et al., 2017). The importance weight for each time step t is given by:

$$w_t = \sum_{i \in \tau} \operatorname{softmax}\left(\frac{\{\langle \mathbf{q}_i, \mathbf{k}_{t'} \rangle\}_{t' \in \tau}}{\sqrt{d}}\right).$$
(7)

In this formulation, the attention mechanism captures the relationships among all instances within the sequence, aligning with the requirements of sequence-level reward prediction. By assigning different weights to each time step, the model can focus on critical moments that disproportionately influence the overall evaluation.

3.3 LEARNING PROCESS

We train the proposed reward model by minimizing the loss between the observable composite re-300 ward $R_{co}(\tau)$ and the predicted reward $R_{co}(\tau)$ using the mean square error $(R_{co}(\tau) - R_{co}(\tau))^2$. 301 For RL, the learned reward function is used to label all state-action pairs. Since we are 302 training a non-Markovian reward function, we provide the model with the past H transitions, 303 $\{(\mathbf{s}_{t-H+1}, \mathbf{a}_{t-H+1}), \dots, (\mathbf{s}_t, \mathbf{a}_t)\}$, and use the output at time step t from the attention layer as the 304 reward for that time step. This approach allows the agent to learn policies that consider the sequence 305 context and assign appropriate credit. The training process alternates between updating the reward 306 model and training the policy: the agent generates delayed reward sequences from environment in-307 teractions to update the reward model, which then labels instance-level rewards to refine the policy. 308 This iterative process leads to mutual improvement. See Appendix B for algorithm details. 309

4 EXPERIMENT

310 311 312

In this section, we begin by evaluating the empirical performance of our proposed method across a variety of benchmark tasks from MuJoCo (Todorov et al., 2012) and the DeepMind Control Suite (Tassa et al., 2018), each featuring different types and lengths of composite delayed rewards. We then assess the effectiveness of proposed method on traditional sum-form delayed rewards, examining whether removing the sum-form assumption affects training compared to established baselines. Lastly, we analyze the relationship between the predicted rewards, the learned attention weights, and the actual environment conditions, providing insights into the interpretability and accuracy of our reward model.

320 321

- 4.1 COMPARE WITH BASELINE METHODS
- **Experiment Setting.** We evaluated our method on benchmark tasks from the MuJoCo locomotion suite (Ant-v2, HalfCheetah-v2, and Walker2d-v2) and the DeepMind Control Suite (fish-upright



Figure 3: A performance comparison of composite delayed rewards, SumSquare (upper), Square-Sum (meddle), and Max (lower), across MuJoCo and DeepMind Control Suite environments with six different delay lengths (5, 25, 50, 100, 200, and 500). The normalized scores are averaged over 3 trials, with the mean and standard deviation computed across a total of 1e6 time steps.

and cartpole-swingup). Differing from standard environments where rewards are assigned at each step, our approach involved assigning a composite reward at the end of sequence while assigning a reward of zero to all other state-action pairs. Unlike previous work, which typically assumes that the sequence-level reward is simply the sum of individual step rewards, our composite reward encompasses more complex forms. Specifically, we included the following composite delayed reward structures, computed over a segment $\tau = \{(\mathbf{s}_t, \mathbf{a}_t)\}_{t=i}^{i+n_i-1}$, where n_i represents the length of the delayed steps:

- SumSquare: The composite delayed reward is the sum of squared step rewards, placing more emphasis on larger rewards: $R_{co} = \sum_{t=i}^{i+n_i-1} \operatorname{abs}(r_t) \cdot r_t$.
- SquareSum: The composite delayed reward is the square of the sum of step rewards, highlighting overall sequence performance: $R_{co} = abs(\sum_{t=i}^{i+n_i-1} r_t) \cdot (\sum_{t=i}^{i+n_i-1} r_t)$.
- Max: The composite delayed reward is a softmax-weighted sum of step rewards, giving more attention to the highest rewards: $R_{co} = \sum_{t=i}^{i+n_i-1} \frac{n_i \cdot e^{\beta r_t}}{\sum_{t'=i}^{i+n_i-1} e^{\beta r_{t'}}} \cdot r_t$,
 - where β is a scaling parameter that controls the sharpness of the softmax distribution.

367 Here, $r_t = r(\mathbf{s}_t, \mathbf{a}_t)$ represents the unobservable Markovian reward at time step t from the original 368 environment. Additionally, we investigated the impact of varying levels of reward delay by setting 369 the delay steps n_i to 5, 25, 50, 100, 200, and 500. These composite delayed rewards complicate credit assignment to state-action pairs and capturing sequence dependencies. Longer delays weaken 370 the correlation between actions and rewards, making it harder to trace rewards back to specific ac-371 tions. Both the complexity of composite rewards and extended delays present significant challenges 372 for effective learning. Our experiments are based on the Soft Actor-Critic (SAC) (Haarnoja et al., 373 2018) algorithm, with the maximum length for each episode fixed at 1000 steps across all tasks. 374 Further experimental details are provided in the Appendix A. 375

376

344

345

347 348

357

359

360

361 362

364 365

366

Baselines. In the comparative analysis, our framework was rigorously evaluated against several leading algorithms in the domain of RL with delayed reward:

7



Figure 4: Performance comparison of sum-form delayed rewards in MuJoCo and DeepMind Control Suite environments with three different delay lengths. The mean and standard deviation of the normalized scores are reported over 6 trials, spanning a total of 1e6 time steps.

- **SAC** (Haarnoja et al., 2018): It directly utilized the original delayed reward information for policy training using the SAC algorithm.
- LIRPG (Zheng et al., 2018): It learned an intrinsic reward function to complement sparse environmental feedback, training policies to maximize combined extrinsic and intrinsic rewards. We used the same code provided by the paper.
- HC (Han et al., 2022): The HC-decomposition framework was utilized to train the policy using a value function that operates on sequences of data. We used the original implementation as provided by the paper.
 - **IRCR** (Gangwani et al., 2020): It adopted a non-parametric uniform delayed reward redistribution. We used the code supplied by the original paper.
- **RRD** (Ren et al., 2021): It employed a reward model trained with a randomized return decomposition loss. We used the same code provided by the paper.
 - **RBT** (Tang et al., 2024): It was employed under the sum-form reward with uniform weight assumption for reward redistribution, utilizing a transformer-based reward model. We employed the code as the original paper.

Evaluation Metric. For evaluation, we report the normalized average accumulative reward across
3 seeds with random initialization to demonstrate the performance of evaluated methods. Higher
accumulative reward in evaluation indicates better performance. Details of the normalization procedure are provided in the Appendix A.3.

Overall Performance Comparison. In this part, the delayed rewards no longer align with the sum-form assumption used in previous work. In Fig. 3, the upper, medium, and lower rows shows the SumSquare, SquareSum, and Max composite delayed reward result, respectively. For SumSquare, our method consistently outperforms baseline methods across different environments, maintaining strong performance even as the delay in reward steps increases. In the SquareSum setting, the delayed reward is more complex, as it emphasizes the cumulative effect of multiple steps rather than focusing on individual large rewards, making it more difficult to attribute rewards to specific actions. Despite this complexity, our method still holds a clear advantage over most baselines, even as perfor-mance slightly decreases with longer delays. In the Max setting, where rewards are sparse and only


Figure 5: Comparison of mean of observed delayed rewards (blue line), predicted rewards (green line), and learned weights (red line) in the Ant environment under two different delayed reward structures: Sum (left) and Max (right). In the Max setting, the blue points indicate the steps with the highest rewards in the original environment. Every 25 steps form a delayed reward sequence, after which a composite delayed reward is assigned. The images below correspond to the behavior of agent at the time steps highlighted by the black frames in the plots.

453 a few key steps contribute to the overall outcome, our method manages to learn effectively, showing 454 that it can still identify and focus on the most critical instances in the sequence, even with limited 455 reward signals. Under these conditions, methods relying on the sum-form Markovian rewards with 456 uniform weights assumption experience a significant drop in performance. This indicates that these 457 methods are heavily dependent on this assumption. It is also worth noting that the HC method does 458 not rely on this assumption, but its effectiveness is limited to shorter delays. As the delay steps 459 increase, its performance drops sharply.

Overall, our proposed method shows a strong ability to handle various composite delayed reward
 structures and adapt to increasing delays, outperforming baseline methods in most environments.
 Moreover, these results highlight the importance of further exploring how to learn reward models
 and train policies in environments with composite delayed rewards, making it worthy of deeper in vestigation. Additional experimental results for varying delayed lengths are provided in Appendix C.

4.2 PERFORMANCE ON TRADITIONAL DELAYED REWARD

Fig. 4 illustrates the results of the traditional sum-form delayed reward setting, where the observed composite delayed rewards are simply the sum of the Markovian rewards provided by the original environment. It is important to note that this experiment is consistent with those conducted in prior works and aligns with the assumptions of baseline methods (Ren et al., 2021; Arjona-Medina et al., 2019; Tang et al., 2024). The figure presents results across short (25 steps), medium (100 steps), and long (500 steps) delays, showing that our method remains robust across different delay levels, performing comparably to the state-of-the-art baselines. The error bars indicate the variance across multiple runs, highlighting the stability of our approach. This demonstrates that our reward model is capable of learning useful information effectively, even without the restrictive sum-form assumption.

477 4.3 CASE STUDY

In this section, we analyze the relationship between the predicted rewards and weights from our reward model and the real rewards from environment under Sum and Max settings of delayed reward composition. We chose these two settings because they provide contrasting scenarios for evaluating the effectiveness of our reward model.

In Fig. 5, we analyze the attention weights learned under both the sum-form and max-form delayed reward settings. For the sum-form delayed reward, where the true reward represents the sum of instance-level contributions, the learned attention weights consistently hover around 1, suggesting that our reward model effectively assigns uniform importance to each time step within the sequence.

On the other hand, in the max-form delayed reward setting, where the reward is determined by the
 highest instance-level contribution, the learned weights emphasize the critical points in the sequence
 corresponding to the highest rewards. Notably, within each sequence, the peaks of the learned
 weights align closely with the peaks of rewards in the true environment (blue points), demonstrating
 that our reward model effectively identifies the most critical instances.

As shown in the images below, when the agent performs stably in both settings, the predicted rewards are high at the corresponding time steps. Conversely, when the performance of agent deteriorates or fails, the predicted rewards drop significantly, even approaching zero. This indicates that the predicted rewards align well with the agent's behavior, confirming that our reward model has effectively learned a reasonable reward allocation.

The match between predicted rewards and agent behavior shows that our reward model closely aligns with the agent's performance. The learned attention weights capture dependencies between individual and overall rewards, demonstrating robustness in delayed and long-term dependencies. This flexibility allows the model to adapt across tasks and reward structures, making it effective in real-world scenarios where traditional Markovian assumptions fail.

501 502

5 RELATED WORK

503 504

505 RL from Delayed Rewards. Learning from aggregated or delayed rewards has been extensively studied, especially in situations where immediate feedback is impractical. Traditional methods rely 506 on trajectory-level feedback, assuming that the reward is the sum of individual step rewards. For 507 example, IRCR (Gangwani et al., 2020) and RRD (Ren et al., 2021) redistribute rewards under 508 Markovian assumptions with equal weighting, while RUDDER (Arjona-Medina et al., 2019) uses 509 recurrent neural networks for credit assignment. Extensions, including expert demonstrations and 510 language models (Liu et al., 2019; Widrich et al., 2021; Patil et al., 2022), further improve preci-511 sion. The Reward Bag Transformer (RBT) (Tang et al., 2024) models the sequential feedback as a 512 bagged reward and employs a transformer to redistribute the bagged reward into instance-level re-513 wards. Although these methods effectively redistribute rewards, they are limited by the assumption 514 that all states contribute equally, which may not reflect the varying importance of different parts of 515 a trajectory in many real-world tasks. Additionally, Han et al. (2022) proposed to directly modifiy 516 RL algorithms by introducing a novel definition of the Q-function to leverage sequence-level infor-517 mation. However, this approach is only suitable for shorter sequence learning and struggles to adapt to more complex, long-term delayed reward scenarios. 518

519

520 **Transformers for RL.** Transformers (Vaswani et al., 2017) have shown significant effectiveness in 521 RL, especially for sample-efficient and generalizable learning, as demonstrated in StarCraft (Vinyals 522 et al., 2019; Zambaldi et al., 2019) and DMLab30 (Parisotto et al., 2020). In offline RL, transformers 523 have been utilized for sequential modeling of RL problems (Chen et al., 2021; Janner et al., 2021; Kim et al., 2023). Luo et al. (2021) combined deep convolution transfer-learning models and inverse 524 RL for reward function acquisition, while Zhang et al. (2023) transformed non-Markovian reward 525 processes into Markovian ones, enhancing online interaction efficiency. In our work, we employ 526 Transformers as reward models to capture non-Markovian reward dependencies in the composite 527 delayed reward problem. 528

529

⁵³⁰ 6 CONCLUSION

531

In this paper, we addressed the challenge of composite delayed reward structures in RL, a problem that extends beyond the traditional sum-form assumption commonly used in existing methods. We proposed an effective solution by introducing a reward model capable of flexibly handling various composite delayed reward structures, incorporating non-Markovian dependencies through an attention mechanism. Our approach consistently outperformed baseline methods across a range of environments and delay settings. While our results show significant improvements, this work opens several avenues for future research. Further exploration is warranted to better model and learn from complex reward structures in more diverse and real-world scenarios. Investigating how to scale our method to even longer delays or more intricate dependency patterns could provide deeper insights.

540 7 ETHICS STATEMENT

This research focuses on developing RL method with composite delayed rewards, with experiments conducted solely on simulated environments such as MuJoCo and DeepMind Control Suite, and without involving any direct human subjects. We acknowledge that reinforcement learning techniques can be applied to sensitive real-world scenarios like healthcare, autonomous driving, and so-cial decision-making, which may have significant ethical implications. Notably, the use of composite delayed rewards in our work inherently reduces the granularity of reward information, thereby help-ing to protect individual privacy when applied to real-world scenarios. Throughout this research, we emphasize fairness and transparency in our methodologies to mitigate the risks of unintended con-sequences from applying the learned policies in such settings. Our model is trained and evaluated on well-established benchmarks, ensuring the reproducibility and reliability of our findings. We do not foresee any discrimination, bias, or privacy concerns arising from this research. Additionally, we adhere to all guidelines regarding dataset use, licensing, and attribution, and have disclosed any potential conflicts of interest. We believe our work aligns with the code of ethics and does not raise ethical concerns related to harmful outcomes or unethical applications.

8 REPRODUCIBILITY STATEMENT

We have made substantial efforts to ensure the reproducibility of our work. Specifically, the im-plementation details of our proposed method, including hyperparameters, training setups, and en-vironment configurations, are provided in the Section 4 and further documented in Appendix A. The datasets and data processing steps are also thoroughly explained in the Appendix A to support replication. All algorithms, including the policy optimization with CoDeTr, are described in de-tail in Section 3 and further supported by pseudocode in Appendix B. Additionally, we provide an anonymous link to the downloadable source code in the abstract to facilitate the replication of our experiments. We believe that these resources collectively contribute to the reproducibility of our results.

594 REFERENCES 595

596 597 598	Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brand- stetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. In <i>The Thirty-</i> <i>second Annual Conference on Advances in Neural Information Processing Systems</i> , 2019.					
599 600	Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep reinforcement learning: A brief survey. <i>IEEE Signal Processing Magazine</i> , 34(6):26–38, 2017.					
601 602 603	Fahiem Bacchus, Craig Boutilier, and Adam Grove. Rewarding behaviors. In <i>Proceedings of the 13th National Conference on Artificial Intelligence</i> , pp. 1160–1167. AAAI Press, 1996.					
604 605 606	Fahiem Bacchus, Craig Boutilier, and Adam Grove. Structured solution methods for non-markovian decision processes. In <i>Proceedings of the 14th National Conference on Artificial Intelligence</i> , pp. 112–117. AAAI Press, 1997.					
607 608	Richard Bellman. Dynamic programming. science, 153(3731):34-37, 1966.					
609 610	Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.					
611 612 613 614	Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In <i>The Thirty-fourth Annual Conference on Advances in Neural Information Processing Systems</i> , 2021.					
616 617 618	Antonio Coronato, Muddasar Naeem, Giuseppe De Pietro, and Giovanni Paragliola. Reinforcement learning for intelligent healthcare applications: A survey. <i>Artificial intelligence in medicine</i> , 109: 101964, 2020.					
619 620 621	Rati Devidze, Parameswaran Kamalaruban, and Adish Singla. Exploration-guided reward shap- ing for reinforcement learning under sparse rewards. In <i>The Thirty-fifth Annual Conference on</i> <i>Advances in Neural Information Processing Systems</i> , 2022.					
622 623 624	Joseph Early, Tom Bewley, Christine Evers, and Sarvapali Ramchurn. Non-markovian reward mod- elling from trajectory labels via interpretable multiple instance learning. In <i>The Thirty-fifth Annual</i> <i>Conference on Advances in Neural Information Processing Systems</i> , 2022.					
625 626 627	Yonathan Efroni, Nadav Merlis, and Shie Mannor. Reinforcement learning with trajectory feedback. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 35, pp. 7288–7295, 2021.					
628 629 630	Roy Frostig, Matthew James Johnson, and Chris Leary. Compiling machine learning programs via high-level tracing. <i>Systems for Machine Learning</i> , 4(9), 2018.					
631 632 633	Tanmay Gangwani, Yuan Zhou, and Jian Peng. Learning guidance rewards with trajectory-space smoothing. In <i>The Thirty-third Annual Conference on Advances in Neural Information Processing Systems</i> , 2020.					
634 635 636	Qitong Gao, Ge Gao, Juncheng Dong, Vahid Tarokh, Min Chi, and Miroslav Pajic. Off-policy evaluation for human feedback. <i>The Thirty-sixth Annual Conference on Advances in Neural In- formation Processing Systems</i> , 2024.					
637 638 639 640	Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In <i>The Thirty-fifth International Conference on Machine Learning</i> . PMLR, 2018.					
641 642	Ben Hambly, Renyuan Xu, and Huining Yang. Recent advances in reinforcement learning in finance. <i>Mathematical Finance</i> , 33(3):437–503, 2023.					
643 644 645	Beining Han, Zhizhou Ren, Zuofan Wu, Yuan Zhou, and Jian Peng. Off-policy reinforcement learn- ing with delayed rewards. In <i>The Thirty-ninth International Conference on Machine Learning</i> . PMLR, 2022.					
647	Michael Janner, Qiyang Li, and Igor Mordatch. Reinforcement learning as one big sequence modeling problem, 2021.					

648 649 650	Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey. <i>Journal of artificial intelligence research</i> , 4:237–285, 1996.						
651 652	Daniel Kahneman. Evaluation by moments: Past and future. In Daniel Kahneman and Amos Tversky (eds.), <i>Choices, Values, and Frames</i> , pp. 693–708. Cambridge University Press, 2000.						
653 654 655 656	Changyeon Kim, Jongjin Park, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Preference transformer: Modeling human preferences using transformers for RL. In <i>The Eleventh International Conference on Learning Representations</i> , 2023.						
657 658 659	B Ravi Kiran, Ibrahim Sobh, Vincent Talpaert, Patrick Mannion, Ahmad Al Sallab, Senthil Yoga- mani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. <i>IEEE Transactions on Intelligent Transportation Systems</i> , 22(6):3239–3258, 2021.						
660 661	Gary Klein. Naturalistic decision making. Human factors, 50(3):456-460, 2008.						
662 663	Ilya Kostrikov. JAXRL: Implementations of Reinforcement Learning algorithms in JAX, 10 2021. URL https://github.com/ikostrikov/jaxrl.						
664 665 666 667 668	Sanjay Krishnan, Animesh Garg, Richard Liaw, Brijen Thananjeyan, Lauren Miller, Florian T Poko- rny, and Ken Goldberg. Swirl: A sequential windowed inverse reinforcement learning algorithm for robot tasks with delayed rewards. <i>The international journal of robotics research</i> , 38(2-3): 126–145, 2019.						
669 670	Yang Liu, Yunan Luo, Yuanyi Zhong, Xi Chen, Qiang Liu, and Jian Peng. Sequence modeling of temporal credit assignment for episodic reinforcement learning, 2019.						
671 672 673	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In <i>The Sixth International Conference on Learning Representations</i> , 2018.						
674 675 676	Wentao Luo, Jianfu Zhang, P. Feng, D. Yu, and Zhijun Wu. A Deep Transfer-Learning-Based Dynamic Reinforcement Learning for Intelligent Tightening System . <i>International Journal of Intelligent Systems</i> , 2021.						
677 678 679	Ben R Newell, David A Lagnado, and David R Shanks. <i>Straight choices: The psychology of decision making</i> . Psychology Press, 2022.						
680 681 682	Emilio Parisotto, H Francis Song, Jack W Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant M Jayakumar, Max Jaderberg, Raphael Kaufman, Aidan Clark, Seb Noury, et al. Stabilizing transformers for reinforcement learning, 2020.						
683 684 685 686 687	Vihang Patil, Markus Hofmarcher, Marius-Constantin Dinu, Matthias Dorfer, Patrick M Blies, Johannes Brandstetter, José Arjona-Medina, and Sepp Hochreiter. Align-rudder: Learning from few demonstrations by reward redistribution. In <i>The Thirty-ninth International Conference on Machine Learning</i> . PMLR, 2022.						
688 689	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language under- standing by generative pre-training, 2018.						
690 691 692 693	David Raposo, Sam Ritter, Adam Santoro, Greg Wayne, Theophane Weber, Matt Botvinick, Hado van Hasselt, and Francis Song. Synthetic returns for long-term credit assignment. <i>arXiv preprint arXiv:2102.12425</i> , 2021.						
694 695 696	Zhizhou Ren, Ruihan Guo, Yuan Zhou, and Jian Peng. Learning long-term reward redistribution via randomized return decomposition. In <i>The Ninth International Conference on Learning Representations</i> , 2021.						
697 698 699	Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement learning framework for autonomous driving. <i>arXiv preprint arXiv:1704.02532</i> , 2017.						
700 701	Shitian Shen and Min Chi. Reinforcement learning: the sooner the better, or the later the better? In <i>Proceedings of the 2016 conference on user modeling adaptation and personalization</i> , pp. 37–44, 2016.						

Yuting Tang, Xin-Qiang Cai, Yao-Xiang Ding, Qiyu Wu, Guoqing Liu, and Masashi Sugiyama. Reinforcement learning from bagged reward. In ICML 2024 Workshop: Aligning Reinforcement Learning Experimentalists and Theorists, 2024. Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud-den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite, 2018. Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. IEEE, 2012. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In The Thirtieth Annual Con-ference on Advances in Neural Information Processing Systems, 2017. Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Juny-oung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. Nature, 2019. Michael Widrich, Markus Hofmarcher, Vihang Prakash Patil, Angela Bitto-Nemling, and Sepp Hochreiter. Modern hopfield networks for return decomposition for delayed rewards. In Deep RL Workshop NeurIPS 2021, 2021. Hongyang Yang, Xiao-Yang Liu, Shan Zhong, and Anwar Walid. Deep reinforcement learning for automated stock trading: An ensemble strategy. In Proceedings of the first ACM international conference on AI in finance, 2020. Chao Yu, Jinpeng Liu, and Shamim Nemati. Reinforcement learning in healthcare: A survey. ACM Computing Surveys (CSUR), 54(5):1–36, 2021. Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Deep reinforcement learning with relational inductive biases. In The Seventh International Conference on Learning Represen-tations, 2019. Hao Zhang, Hao Wang, and Zhen Kan. Exploiting transformer in sparse reward reinforcement learning for interpretable temporal logic motion planning. IEEE Robotics and Automation Letters, 2023. Zeyu Zheng, Junhyuk Oh, and Satinder Singh. On learning intrinsic rewards for policy gradient methods. In The Thirty-first Annual Conference on Advances in Neural Information Processing Systems, 2018.

756 IMPLEMENTATION DETAILS А

758 A.1 BENCHMARKS WITH COMPOSITE DELAYED REWARD 759

760 In this work, we introduced a new problem setting, called composite delayed rewards, within the 761 suite of locomotion benchmark tasks in both the MuJoCo environment and the DeepMind Control Suite. Our experiments were conducted using the OpenAI Gym platform (Brockman et al., 2016) 762 and the DeepMind Control Suite (Tassa et al., 2018), focusing on tasks with extended horizons and a fixed maximum trajectory length of T = 1000. We utilized MuJoCo version 2.0 for our 764 simulations, which can be accessed at http://www.mujoco.org/. MuJoCo operates under 765 a commercial license, and we ensured full compliance with its licensing terms. Additionally, the 766 DeepMind Control Suite, distributed under the Apache License 2.0, was used in accordance with its 767 licensing requirements. 768

Experiments involving composite delayed rewards with varying delay steps (5, 25, 50, 100, 200, and 769 500) and different composite types (SumSquare, SquareSum, Max, and the traditional Sum) were 770 conducted to validate the effectiveness of the proposed method. In the Max experiment, the scaling 771 parameter β is set to 3. To evaluate its performance, commonly used delayed reward algorithms 772 were adapted to fit within the composite delayed reward framework, acting as baselines. In these 773 experiments, each segment with composite delayed rewards was treated as an independent trajectory, 774 and the modified algorithms were applied accordingly.

775 776 777

787 788

789

801 802 A.2 IMPLEMENTATION DETAILS AND HYPER-PARAMETER CONFIGURATION

778 In our experiments, the policy optimization module was implemented based on soft actor-critic 779 (SAC) (Haarnoja et al., 2018). We evaluated the performance of our proposed methods with the same configuration of hyper-parameters in all environments. The back-end SAC followed the JaxRL 781 implementation (Kostrikov, 2021), which is available under the MIT License.

782 The proposed CoDeTr was built upon the GPT implementation in JAX (Frostig et al., 2018), avail-783 able under the Apache License 2.0. Our experiments employed a Causal Transformer with three 784 layers and four self-attention heads, followed by an in-sequence bidirectional attention layer with 785 one self-attention head. For a comprehensive overview of the CoDeTr's hyper-parameter settings, 786 please refer to Table 1.

Table 1: Hyper-parameters of CoDeTr.

790		
791	Hyper-parameter	Value
700	Number of Causal Transformer layers	3
792	Number of in-sequence attention layers	1
793	Number of attention heads	4
794	Embedding dimension	256
795	Batch size	64
796	Dropout rate	0.1
797	Learning rate	0.00005
798	Optimizer	AdamW (Loshchilov & Hutter, 2018)
799	Weight decay	0.0001
800	Warmup steps	100
801	Total gradient steps	10000

For the baseline methods, the IRCR (Gangwani et al., 2020) method was implemented following the 803 descriptions provided in the original paper. Both the RRD (Ren et al., 2021) and LIRPG (Zheng 804 et al., 2018) methods are distributed under the MIT License. The code for HC (Han et al., 805 2022) is available in the supplementary material at https://openreview.net/forum?id= 806 nsjkNB20KsQ, while the code for RBT (Tang et al., 2024) can be found in the original paper. 807

To maintain consistency in the policy optimization process across all methods, each was subjected 808 to 1,000,000 training iterations. For the proposed method, a dataset of 10,000 time steps was first 809 gathered to pre-train the reward model. This model underwent 100 pre-training iterations, which was deemed necessary to properly initialize the reward model before commencing the main policy learning phase. After this warm-up period, the reward model was updated for 10 iterations following the addition of each new trajectory. Furthermore, to monitor performance systematically, evaluations were conducted every 5,000 time steps. During prediction, the sequence length used for prediction is set to H = 100. All computations were performed on NVIDIA GeForce A100 GPUs with 40GB of memory, which were dedicated to both training and evaluation tasks.

817 A.3 DATA NORMALIZATION PROCEDURES

The normalization process for our data varies depending on the type of composite delayed reward. Specifically:

• For **SumSquare**, the normalization is calculated as:

$$\frac{\sum \hat{R}_{\rm co}}{T \cdot \sum (r_{\rm max})^2}.$$

• For **SquareSum**, the normalization is given by:

$$\frac{\sum \hat{R}_{\rm co}}{\sum \left(\frac{T}{n} \cdot (r_{\rm max} \cdot n)^2\right)}.$$

• For Max, the normalization is computed as:

$$\frac{\sum \hat{R}_{\rm co}}{\sum r_{\rm max}}.$$

Here, \hat{R}_{co} represents the predicted composite delayed reward, T is the total number of time steps in a trajectory, r_{max} is the maximum possible reward in the environment, and n is the number of delayed steps in each segment.

In essence, the normalization process involves scaling $\sum R_{co}$ by the maximum achievable reward in the given environment. This approach ensures that results from experiments with different delayed steps are on the same scale, enabling meaningful comparisons across varying delayed steps and their impact on the learning process.

842 843

844 845

816

818

820 821

823 824 825

835

836

837

B ALGORITHM

Algorithm 1 Policy Optimization with CoDeTr 846 1: Initialize: replay buffer \mathcal{D} , CoDeTr parameters ψ , and policy π . 847 2: while training is not complete do 848 3: Collect a trajectory \mathcal{T} by interacting with the environment using the current policy π . 849 Store trajectory \mathcal{T} with composite delayed reward information based on sequences 4: 850 $\{(\tau, R_{\rm co}(\tau))\}$ in replay buffer $\mathcal{D}_{...}$ 851 5: Sample batches from replay buffer \mathcal{D} . 852 Compute the mean squared error loss $(R_{\rm co}(\tau) - \hat{R}_{\rm co}(\tau))^2$ for CoDeTr using the sampled 6: 853 sequences from the replay buffer. 854 7: Update CoDeTr parameters ψ based on the computed loss. 855 Relabel instance-level rewards in replay buffer \mathcal{D} using the updated CoDeTr. 8: 9: Optimize policy π using the relabeled data with an off-the-shelf RL algorithm (e.g., 856 SAC (Haarnoja et al., 2018)). 10: end while 858 859

The training process involves alternating between updating the reward model and optimizing the policy, which creates a continuous loop of mutual improvement. First, the agent collects trajectories by interacting with the environment according to the current policy. These trajectories are then used to train the CoDeTr, which learns to predict instance-level rewards and composite delayed rewards for sequences. The training is done by minimizing the mean squared error (MSE) loss between the predicted composite reward $R_{co}(\tau)$ and the observed composite delayed reward $\hat{R}_{co}(\tau)$. This loss function allows CoDeTr to accurately capture the relationships and dependencies within each sequence, ensuring that both individual and sequence-level contributions are effectively represented. Using the updated CoDeTr model, the rewards for state-action pairs in replay buffer are relabeled, providing more accurate feedback for policy optimization. The updated rewards are used to further refine the policy using reinforcement learning algorithms like SAC, enabling the agent to learn effective strategies even in environments with delayed rewards. This iterative procedure enhances both the reward model and the policy through each training cycle.

C ADDITIONAL RESULT

872 873

874 875

876

877

878

879

880

882 883

893

894

895

896

897

898

899 900 901

902

Table 2: Performance comparison across different settings, utilizing various delayed steps ranging from 25 to 200 in the Ant-v2 environment, evaluated over 3 independent trials. The scores presented are normalized to ensure comparability across different configurations. The methods that demonstrated the best performance, along with those that were statistically comparable based on a paired t-test at a significance level of 5%, are highlighted in boldface for emphasis.

Delayed Type	SAC	LIRPG	HC	IRCR	RRD	RBT	CoDeTr(ours)
Sum	$ \begin{array}{c c} 0.0004 \\ (0.0002) \end{array} $	-0.1759 (0.0631)	$0.0025 \\ (0.0058)$	$\begin{array}{c} 0.03364 \\ (0.0281) \end{array}$	$\begin{array}{c} 0.3327 \\ (0.2095) \end{array}$	0.5699 (0.0162)	0.5493 (0.0187)
SumSquare	$\begin{array}{ c c c } -0.0067 \\ (0.0022) \end{array}$	-0.0159 (0.0027)	$0.0198 \\ (0.0005)$	-0.0617 (0.0273)	$0.0308 \\ (0.0207)$	$\begin{array}{c} 0.2821 \\ (0.0703) \end{array}$	0.3910 (0.0431)
SquareSum	$\begin{array}{c c} -0.0902 \\ (0.1031) \end{array}$	-0.0012 (0.0001)	-0.0280 (0.0275)	-0.1060 (0.0303)	$\begin{array}{c} 0.0575 \\ (0.0173) \end{array}$	$\begin{array}{c} 0.0890 \\ (0.0240) \end{array}$	0.1992 (0.0110)
Max	$\begin{array}{c} 0.04093 \\ (0.0065) \end{array}$	-0.1982 (0.0373)	$0.0108 \\ (0.0103)$	-0.0093 (0.0524)	$\begin{array}{c} 0.2193 \\ (0.0416) \end{array}$	$\begin{array}{c} 0.4669 \\ (0.1078) \end{array}$	0.5318 (0.0821)

In Table 2, the performance of different methods is compared across various composite delayed reward types: Sum, SumSquare, SquareSum, and Max, on the Ant-v2 environment. Our proposed method, CoDeTr, consistently performed well across all composite delayed reward configurations, either achieving the best results or performing comparably to the top baseline methods. In particular, CoDeTr showed strong performance under different composite delayed reward settings, demonstrating its ability to handle complex reward structures effectively. These results indicate that our approach is robust and adaptable, providing high-quality performance across a range of composite delayed reward scenarios.

D DISCUSSION

Limitation. Our experimental results demonstrate that the proposed approach performs effectively 903 across various types of composite delayed rewards, showing notable improvements over baseline 904 methods. However, we also observed increasing difficulty in efficiently learning the policy as the 905 delay length grew longer. This challenge arises because longer delays weaken the temporal connec-906 tion between specific actions and their resulting outcomes, increasing uncertainty when attempting 907 to determine which actions contributed to the observed reward. Consequently, the diminished ability 908 to effectively assign credit to individual actions complicates the policy training process, leading to 909 slower convergence and reduced overall performance in scenarios with extended delay lengths, as 910 evidenced in our experiments.

911

Future Direction. A key area for future work lies in addressing the challenges posed by longer
 reward delays. Our experiments have shown that increasing the delay length significantly compli cates credit assignment to individual actions. To better capture long-range dependencies in such
 settings, future research could focus on developing advanced temporal credit assignment methods,
 such as improved attention mechanisms or memory-augmented neural networks. These techniques
 may enhance the model's ability to trace rewards back to responsible actions, even in situations with extended delays.

Expanding the use of composite delayed rewards to broader application scenarios represents another
 promising direction. Domains such as healthcare, autonomous vehicles, and industrial robotics often
 involve delayed and complex feedback that makes instance-level rewards impractical. Investigating
 how our proposed approach can generalize to these real-world applications would demonstrate its
 practical utility and robustness. Moreover, such exploration could help identify potential modifica tions required to adapt the framework to specific challenges, such as safety and real-time require ments inherent in these domains.

In addition, integrating human-in-the-loop feedback with composite delayed rewards could significantly enhance the learning process. Human evaluators often assign feedback based on pivotal events and use non-linear reasoning, which traditional reward models may fail to capture. Incorporating human feedback more directly, possibly through preference learning models aligned with composite delayed rewards, could improve the agent's ability to learn behaviors that align with human expectations. This approach would be particularly valuable in interactive environments where understanding human intent is crucial for the agent's success.