

# SYMMETRIC REINFORCEMENT LEARNING LOSS FOR ROBUST LEARNING ON DIVERSE TASKS AND MODEL SCALES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reinforcement learning (RL) training is inherently unstable due to factors such as moving targets and high gradient variance. Reinforcement Learning from Human Feedback (RLHF) and Reinforcement Learning from AI Feedback (RLAIF) introduce additional challenges. For instance, diverse preferences complicate the alignment process, and prediction errors in a trained reward model can become more severe as the LLM generates unseen outputs. These RL challenges create confusion about whether the probability of an action for a given state should be increased or decreased, similar to the noise in labels for classification tasks. In this work, we enhance the stability of the RL training procedure by adapting reverse cross-entropy (RCE) from supervised learning for noisy data to define a symmetric RL loss. We demonstrate performance improvements across various tasks and scales. We conduct experiments in discrete action tasks (Atari games) and continuous action space tasks (MuJoCo benchmark and Box2D) using Symmetric A2C (SA2C) and Symmetric PPO (SPPO), with and without added noise. Notably, SPPO shows strong performance across different hyperparameters. Furthermore, we validate the benefits of the symmetric RL loss in the RLHF framework using PPO for natural language processing tasks, demonstrating improved performance in tasks such as IMDB positive sentiment and TL;DR summarization.

## 1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) have shown impressive performance across various natural language processing tasks (Chung et al., 2022; Wei et al., 2023), robot control (Huang et al., 2022; Driess et al., 2023), and healthcare (Lee et al., 2023c; Huang et al., 2020). However, as these LLMs are typically trained to predict the next word in a provided dataset, they require post-training processing to make them useful for particular tasks. Reinforcement Learning from Human Feedback (RLHF) trains LLMs to generate responses aligned with user preferences through human feedback. Additionally, Reinforcement Learning from AI Feedback (RLAIF), which leverages feedback from well-trained AI models, has also been employed (Lee et al., 2023a; Bai et al., 2022). Thus, adapting fundamental Reinforcement Learning (RL) algorithms such as REINFORCE (Williams, 1992), A2C (Mnih et al., 2016), and PPO (Schulman et al., 2017) to suit the fine-tuning of LLMs for LLM tasks is an area of active interest (Ahmadian et al., 2024; Ouyang et al., 2022; Rafailov et al., 2023).

RL methods (Sutton et al., 2000; Sutton & Barto, 2018a) have led to substantial breakthroughs in tasks such as robot control and game playing. Still, they entail learning instability compared to supervised learning due to factors such as moving targets, high-gradient variance, and training value functions. The RL literature has proposed various methods to make the RL learning process more robust, such as preventing overestimation with Double DQN (van Hasselt et al., 2015), reducing variance with Generalized Advantage Estimation (GAE) (Schulman et al., 2018), updates within the trust region (Schulman et al., 2015; 2017), and encouraging diverse behavior with Soft Actor-Critic (SAC) (Haarnoja et al., 2018). In addition to the methods devised specifically for RL problems, RL literature has also adopted supervised learning techniques to make the learning process more robust. For example, ensembles have been used for more accurate value function prediction, while Layer

054 Normalization and Batch Normalization have been employed to constrain predictions for out-of-  
055 distribution samples, thereby mitigating the overestimation and extrapolation.

056 RLHF (Ouyang et al., 2022; Lee et al., 2023b) and RLAIIF Lee et al. (2023a); Bai et al. (2022); Byun  
057 et al. (2024) potentially introduce additional training challenges. For example, these algorithms of-  
058 ten receive feedback from multiple sources (human or AI models) to align LLMs, and each feedback  
059 provider may have different preferences, meaning a sample considered preferable by one provider  
060 could be deemed undesirable by another (Ethayarajh et al., 2024; Chakraborty et al., 2024). In ad-  
061 dition, RLHF and RLAIIF often leverage a trained reward model to provide feedback on samples  
062 generated by the LLM. This indirection raises the question: *does the learned reward model provide*  
063 *the correct reward?* The reward model has prediction errors itself (See Figure 1), but as the LLM  
064 is trained with RL, its outputs deviate from the reward model’s training dataset, introducing more  
065 error in the reward model’s predictions for out-of-distribution samples.

066 The challenges associated with RL, RLHF, and RLAIIF, as mentioned above, can introduce confusion  
067 when calculating advantage values in RL algorithms like A2C and PPO. Specifically, an action that  
068 should have a positive advantage value may have a negative sign in the next update, depending  
069 on which samples (states, actions) are generated and how the batch is composed during advantage  
070 normalization. The sign of the advantage determines whether the probability of a corresponding  
071 action for a given state increases or decreases in policy gradient algorithms. If the advantages are  
072 predicted incorrectly, this can lead to learning in the opposite direction. We hypothesize that these  
073 difficulties are similar to noisy classification tasks in supervised learning, where some labels are  
074 incorrect.

075 In this paper, we leverage a technique developed for classification tasks with noisy labels, employing  
076 a robust loss function to enhance the learning procedures of A2C and PPO. We define a symmetric  
077 RL loss, whose fundamental mechanism aligns with the robust loss function used in supervised  
078 learning (Wang et al., 2019), to improve the robustness of RL procedure for A2C and PPO (See  
079 Section 4.3). We apply this symmetric RL loss to A2C and PPO, naming them Symmetric A2C  
080 (SA2C) and Symmetric PPO (SPPO), and evaluate their performance across various tasks and model  
081 scales.

082 First, we assess the performance gains of SA2C and SPPO on Atari games (Mnih et al., 2016),  
083 which have discrete action spaces, as well as on the MuJoCo benchmark (Todorov et al., 2012) and  
084 Box2D (Catto, 2011) environments, which have continuous action spaces. For these control tasks,  
085 we introduce a noisy reward variant, hypothesizing that it will increase confusion in advantage  
086 prediction to better evaluate our method. Additionally, we test our method on RLHF tasks using  
087 LLMs, such as IMDB positive sentiment analysis (Maas et al., 2011) and TL;DR summarization  
088 (Völske et al., 2017). The IMDB task involves generating positive sentiment for a given context and  
089 TL;DR is a summarization task where an LLM is required to summarize content.

090 SA2C and SPPO demonstrate better performance improvements across diverse control tasks com-  
091 pared to A2C and PPO. Notably, both SA2C and SPPO perform well in settings with added noise  
092 to the reward. Additionally, SPPO shows consistent performance improvements across various hy-  
093 perparameters (Table 4). We analyze why SPPO exhibits more robust improvements than SA2C  
094 in Section 5.4. Furthermore, SPPO shows superior performance to PPO in RLHF tasks, such as  
095 IMDB positive sentiment and TL;DR summarization. We demonstrate SPPO outperforming PPO  
096 on reward in both tasks, and SPPO’s summarization is significantly better, as measured by win-rate  
097 against PPO, judged by GPT-4 Turbo (gpt-4-turbo-2024-04-09).

098 In summary, our key contributions are:

- 099 • We propose the symmetric RL loss for A2C and PPO, along with the gradient analysis that  
100 aligns with the gradient behavior of robust loss functions used in noisy classification tasks  
101 in Section 4.3.
- 102 • We conduct experiments across various environments and model scales, demonstrating per-  
103 formance improvements to validate the effectiveness of the symmetric RL loss for general  
104 control tasks and RLHF tasks in Section 5.
- 105 • We analyze how PPO can introduce additional confusion in advantage estimates, which  
106 justifies using symmetric RL loss (See Section 5.4). This shows that SPPO demonstrates  
107 consistent improvement across a range of hyperparameters.

## 2 RELATED WORK

We briefly introduce robust loss functions studied in the context of noise in supervised learning classification tasks. Ghosh et al. (2017) prove that, in the presence of a noisy dataset, the mean absolute error (MAE) has a slower learning speed compared to cross-entropy loss (CE), but the model learns more robustly. Zhang & Sabuncu (2018) propose a generalized cross entropy loss  $L_q$ , which becomes CE when  $q \rightarrow 0$ , and becomes MAE when  $q \rightarrow 1$ . By adjusting this parameter  $0 \leq q \leq 1$ , robust learning is achieved in noisy datasets. The symmetric cross entropy (SCE) (Wang et al., 2019) that we mainly refer to suggests a symmetric cross-entropy loss. This loss not only considers the flow of information from the true distribution to the model’s predictions but also incorporates information flowing in the reverse direction. SCE works better than GCE in general, especially for data with high noise rates. Ma et al. (2020) introduce various loss functions and classify them into types: *Active Loss* and *Passive Loss* functions. They demonstrate that normalizing the loss can help improve robustness. They use a combination of one active loss and one passive loss like SCE. We define a loss function that considers reverse information to match the RL version and use it to improve the RL procedure.

In the RL literature, Wang et al. (2018) proposes using a confusion matrix to handle perturbed rewards, predicting surrogate rewards for robust policy updates. While this method appears effective for Atari games, later research (Chen et al., 2024) shows that it does not outperform corresponding baselines in continuous tasks. Additionally, introducing noise in RL has demonstrated performance benefits. For instance, Obando-Ceron et al. (2023) show that smaller batch sizes improve performance, and Schaul et al. (2022) present that policy churn aids exploration. These studies primarily conduct experiments on Atari games, which require navigating many novel states. However, whether noise is beneficial or not in continuous action spaces remains debatable (Mai et al., 2022; Byun & Perrault, 2024). Our work proposes a robust loss function designed to handle noise (confusion in advantage prediction) without judging whether the noise is beneficial.

Reinforcement Learning from Human Feedback (RLHF) Ouyang et al. (2022); Lee et al. (2023b) and Reinforcement Learning from AI Feedback (RLAIF) (Lee et al., 2023a; Bai et al., 2022) have contributed to the success of large language models (LLMs) by aligning them with user preferences. However, these methods require training a reward model and a value function. Each of these components has prediction errors, and finding appropriate hyperparameters for training requires significant effort. Direct Preference Optimization (DPO) (Rafailov et al., 2023) eliminates the cost associated with the reward model by rearranging PPO loss for ranking-based feedback (e.g., sample A is preferred over sample B). Ethayarajh et al. (2024) remove the requirement ranking-based feedback by modifying DPO loss further, allowing a model to be trained with bad or good labels. Additionally, Chakraborty et al. (2024) demonstrate that feedback from diverse people, each with different preferences, makes a single reward model difficult to reflect preferences correctly. Recent studies focus on sentence-level feedback (Lightman et al., 2023; Wang et al., 2024), but DPO and KTO cannot utilize sentence-level feedback. Therefore, we propose the reverse RL loss term, which can make PPO in existing RLHF methods more robust.

## 3 PRELIMINARIES

### 3.1 REINFORCEMENT LEARNING

Reinforcement Learning (RL) formulates a Markov decision process (MDP) (Puterman, 2014; Sutton & Barto, 2018b) defined by the tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma, \mu)$ . At each timestep  $t$ , an action  $a_t \in \mathcal{A}$  is sampled from an agent’s policy  $\pi_\theta(\cdot | s_t)$  for a given state  $s_t \in \mathcal{S}$ . For the taken action  $a_t$ , the reward function returns a reward  $\mathcal{R}(s_t, a_t)$  where  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ , and the transition probability  $\mathcal{P}(\cdot | s_t, a_t)$  determines the next state  $s_{t+1}$ .  $\gamma$  is the discount factor, and  $\mu$  represents the initial state distribution for  $s_0$ . The RL objective is to find the optimal  $\theta$  that maximizes the expected discounted sum of rewards:

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}_{\substack{s_0 \sim \mu \\ a_t \sim \pi_\theta(\cdot | s_t) \\ s_{t+1} \sim \mathcal{P}(\cdot | s_t, a_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]. \quad (1)$$

### 3.2 A2C AND PPO ALGORITHMS

The Advantage Actor-Critic (A2C) algorithm (Mnih et al., 2016) is an actor-critic method that combines value-based and policy-based approaches. A2C uses the advantage function  $A$  to reduce the variance in policy updates. The policy  $\pi_\theta$  is updated by following the gradient of the objective function to maximize the sum of rewards as defined in 1:

$$\nabla_\theta J(\pi_\theta) = \sum_{t=0} \nabla_\theta \log \pi_\theta(a_t | s_t) A(s_t, a_t) \quad (2)$$

Proximal Policy Optimization (PPO) (Schulman et al., 2017) aims to update the policy within a trust region. This is achieved through a clipped loss function to ensure that the new policy does not deviate too much from the old policy. The PPO loss function can be written as:

$$L_{\text{ppo}}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)] \quad (3)$$

where  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$  is the probability ratio, and  $\epsilon$  is a small hyperparameter that controls the range of the clipping. The advantage function estimates how much better an action  $a$  is compared to the other actions for at a given state  $s$ . Both algorithms increase the probability of  $a$  for  $s$  if the corresponding advantage  $A(a, s) > 0$  and decrease it if  $A(a, s) < 0$ . In the approach section, we introduce the connection between A2C and PPO with the cross-entropy loss for classification and define the symmetric RL loss.

### 3.3 SYMMETRIC CROSS ENTROPY

Symmetric Cross Entropy (SCE) (Wang et al., 2019) is designed for noisy classification datasets. Cross Entropy (CE) loss (Equation 4) performs effectively when the data is clean; however, it encounters challenges in the presence of noise. Given a true distribution  $q$  and a predicted distribution  $p$ ,  $p$  is learned based on the information derived from  $q$  according to information theory. However, when  $q$  is noisy,  $p$  can only approximate the true distribution to a limited extent. To address this issue, Symmetric Cross Entropy (SCE) also consider incorporates information in the opposite direction through Reverse Cross Entropy (RCE) (Equation 5).

$$L_{\text{ce}} = - \sum_{k=1}^K q(k|\mathbf{x}) \log p(k|\mathbf{x}) \quad (4)$$

$$L_{\text{rce}} = - \sum_{k=1}^K p(k|\mathbf{x}) \log q(k|\mathbf{x}) \quad (5)$$

where  $k \in \{1, \dots, K\}$  is a class and  $\mathbf{x}$  is an input. RCE loss has been proven to be robust to a certain amount of noise, but the learning speed is too slow. Therefore, SCE combines CE and RCE losses (Equation 6),

$$L_{\text{sce}} = \alpha L_{\text{ce}} + \beta L_{\text{rce}} \quad (6)$$

where  $\alpha$  and  $\beta$  are constants determining the contribution of each part. SCE demonstrates performance improvement across various noisy ratios and types. As mentioned in the introduction section, the RL training process can lead to noisy advantage predictions, so we propose a symmetric RL loss in the next approach section.

## 4 APPROACH

This section introduces the reverse RL loss and proposes the symmetric RL loss for A2C (Mnih et al., 2016) and PPO (Schulman et al., 2017), an RL version of Symmetric Cross Entropy (SCE) (Wang et al., 2019). A2C and PPO training procedures basically increase or decrease the probability of an action depending on the advantage sign, but the advantage prediction involves noise due to several factors. A highly engineered reward function is required to eliminate errors, and the trained reward model has a prediction error in RLHF (Ouyang et al., 2022) and RLAIIF (Lee et al., 2023a; Bai et al., 2022). Receiving feedback from multiple sources further complicates the training of the reward model (Chakraborty et al., 2024). Additionally, the value function also has model errors,

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

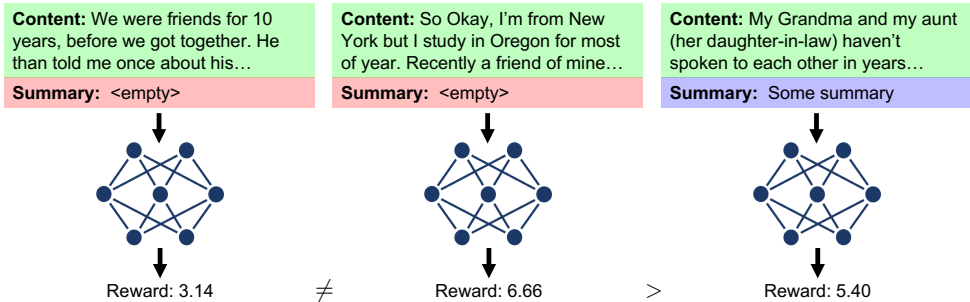


Figure 1: Example of reward prediction errors in a trained reward model for TL;DR summarization. The generated summary samples (left and middle) are both *empty*, yet they receive significantly different rewards. The middle sample is higher than some summarization text (right) and even scores higher (6.66) than the average reward score of SPPO (6.13). The full text for these samples can be found in Appendix 15.

and the sign of the advantage in advantage normalization depends on how the batch is composed. PPO increases sample efficiency compared to A2C, but the off-policy part can introduce confusion in advantage predictions (See Section 5.4). Similar to SCE, which is robust to noisy data, the symmetric RL loss contributes to robust learning in an RL environment that can introduce noisy predictions.

#### 4.1 REVERSE REINFORCEMENT LEARNING LOSS

Given a true (target) distribution  $q$  and a predicted distribution  $p$ , if  $q$  is noisy, training  $p$  can be challenging and  $p$  cannot accurately reflect the true distribution. Reverse Cross Entropy (RCE) considers the reverse information from  $p$ . We propose that the reverse RL losses for A2C and PPO also incorporate reverse information to address noisy factors in the RL training procedure. The RCE loss (Equation 5) defines  $\log 0 = Z$  where  $Z < 0$  is some constant for  $q(k|\mathbf{x}) = 0$ . We also use this definition for the negative advantage and this also is useful to prove the robustness of the reverse RL losses. For all tasks conducted in this paper, we use  $Z = -1$ . Note that the constant terms  $Z$  and  $\beta$  in Equation 4 and 9 are multiplied together, so we control the impact of the reverse RL loss solely by adjusting  $\beta$ . For example,  $(\beta = 1.0, Z = -1.0)$  and  $(\beta = 10.0, Z = -0.1)$  yield the exact same results. Suppose there exist  $k$  actions and  $a^{(i)}$  indicates  $i^{\text{th}}$  action.  $\pi_{\theta}^{(i)} = \pi_{\theta}(a^{(i)}|s)$  for a state  $s$ . Let’s denote the possible action probabilities set  $s$  as  $\pi_{\theta}(s) = \{\pi_{\theta}^{(1)}, \pi_{\theta}^{(2)}, \dots, \pi_{\theta}^{(k)}\}$ . Note that we discretize the continuous action space for continuous action tasks (Tang & Agrawal, 2020). One thing we need to note is that when updating a policy, we use advantages instead of label sets in RL. Advantages can have negative values (negative labels) unlike ordinary labels. We only consider the sign of the advantage<sup>1</sup> because this advantage is the role of the label in supervised learning. For a sampled action probability  $\pi_{\theta}^{(i)}$  and the corresponding advantage  $A(s, a^{(i)}) = A^{(i)}$ , the *sample-wise reverse A2C (RA2C) loss* is:

$$L_{\text{ra2c}}(\pi_{\theta}(s), A^{(i)}) = \begin{cases} \sum_{j \in [k] \setminus \{i\}} -\pi_{\theta}^{(j)} A^{(i)} Z, & \text{if } A^{(i)} > 0 \\ \sum_{j \in [k] \setminus \{i\}} \pi_{\theta}^{(j)} A^{(i)} Z, & \text{if } A^{(i)} < 0 \end{cases} \quad (7)$$

For a positive advantage  $A$ , the difference between A2C’s loss  $A \log \pi$  and CE loss  $1 \log p$  is that A2C can be considered as CE multiplied by the advantage. In terms of gradients,  $A$  is a constant, so A2C reflects the information  $A$  times more strongly than the CE loss. Thus, we also reflect the reverse direction  $A$  times more strongly. Similarly, since PPO has  $\pi_{\text{old}}^{(i)}$  term in the loss, the sample-wise reverse PPO (RPPO) loss just introduces the additional constant  $\pi_{\text{old}}^{(i)}$  for a sampled action

<sup>1</sup>We do not consider when the advantages are zero because those are not considered when updating a policy.

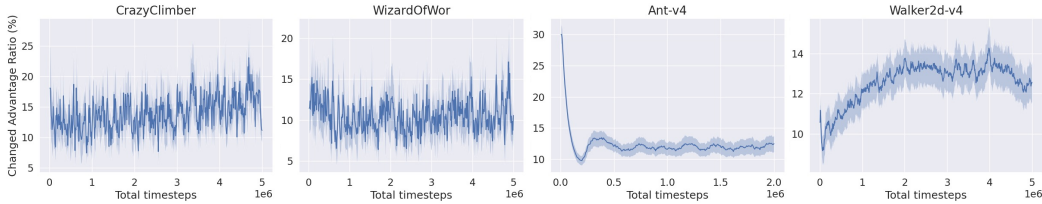


Figure 2: Change of advantage rate (%): The graphs show how often the advantage signs flip in various environments as training progresses. In Atari games, often over 5% of samples change signs, while in MuJoCo tasks, usually over 10% of samples change signs after the advantage normalization. We use 5 different random seeds for CrazyClimber and WizardOfWor, and 30 different random seeds for Ant-v4 and Walker2d-v4. The line is the mean of the change ratio across the seeds, and the shaded area represents standard errors.

probability  $\pi_{\theta}^{(i)}$  to consider the same amount of reverse information:

$$L_{\text{rppo}}(\pi_{\theta}(s), A^{(i)}, \pi_{\text{old}}^{(i)}) = \begin{cases} \sum_{j \in [k] \setminus \{i\}} -\frac{\pi_{\theta}^{(j)} A^{(i)} Z}{\pi_{\text{old}}^{(i)}}, & \text{if } A^{(i)} > 0 \\ \sum_{j \in [k] \setminus \{i\}} \frac{\pi_{\theta}^{(j)} A^{(i)} Z}{\pi_{\text{old}}^{(i)}}, & \text{if } A^{(i)} < 0 \end{cases} \quad (8)$$

We define the symmetric RL loss, which is composed of the original RL loss (A2C or PPO) and the corresponding reverse RL loss, in Section 4.2. We then analyze why these reverse RL losses contribute to the RL learning procedure in Section 4.3.

#### 4.2 SYMMETRIC REINFORCEMENT LEARNING LOSS

The *Symmetric Reinforcement Learning (SRL) loss*  $L_{\text{srl}}$  consists of two parts like SCE (Equation 6): the original actor loss  $L_{\text{rl}}$  (A2C or PPO) and the corresponding reverse RL loss  $L_{\text{rev}}$  (RA2C or RPPO).  $L_{\text{srl}}$  flexibly adjust the symmetric learning framework with two additional hyperparameters ( $\alpha > 0$  and  $\beta > 0$ ) as follows:

$$L_{\text{srl}} = \alpha L_{\text{rl}} + \beta L_{\text{rev}} \quad (9)$$

We name A2C and PPO using the symmetric RL loss as *Symmetric A2C (SA2C)* and *Symmetric PPO (SPPO)*, respectively. The meanings of  $\alpha$  and  $\beta$  align with SCE, where  $\alpha$  represents the degree of actively training a policy, and  $\beta$  serves as auxiliary support to stabilize the entire learning process. In the following section, we analyze the gradient of the two types of losses.

#### 4.3 GRADIENT ANALYSIS

For an input  $\mathbf{x}$  and the corresponding correct label  $k$ , the cross entropy (CE) loss gradient is  $-\frac{1}{p_{\theta}(k|\mathbf{x})} \nabla_{\theta} p_{\theta}(k|\mathbf{x})$ . Smaller  $p_{\theta}$  values aggressively increase the magnitude of the gradient. CE loss rapidly increases uncertain predictions. If there is no noise, this method is correct, but it may lead to incorrect predictions on noisy datasets and excessive overfitting (Zhang & Sabuncu, 2018). A2C and PPO losses also have the same issue. For A2C, the gradient is simply multiplied by an advantage  $A$ , i.e.,  $-\frac{A(s,a)}{\pi_{\theta}(a|s)} \nabla_{\theta} \pi_{\theta}(a|s)$ . In the case of PPO, the magnitude of the gradient tends to increase as the probability of an action decreases. Consider a sample that passes the clipping function: the difference between  $\pi_{\text{old}}$  and  $\pi$  is within the  $\epsilon$  bound. As the denominator  $\pi_{\text{old}}$  gets smaller, the magnitude of the gradient increases.

**Detailed analysis:** The symmetric RL loss gradient analysis aligns with the analysis of SCE. For simplicity, we set  $\alpha$  and  $\beta$  to 1 and examine the gradient direction for two types of A2C loss (RL and reverse RL) with respect to the action logits  $z$ . We use the notation defined in Section 4.1 and introduce the case when  $A^{(i)} > 0$ . For the full derivation including SPPO and  $A^{(i)} < 0$ , please refer to Appendix A. The sample-wise SA2C loss is:

$$L_{\text{sa2c}} = L_{\text{a2c}} + L_{\text{ra2c}} \quad (10)$$

Table 1: Mean final scores and standard errors (over the last 10 episodes) of PPO and SPPO on Atari games, without and with binary symmetric channel (BSC) noise with a crossover probability of 0.1 across 5 seeds. Full results can be found in Table 11.

	Without Noise		$\epsilon \sim \text{BSC}(0.1)$	
	PPO	SPPO	PPO	SPPO
Alien	<b>1128 ± 105</b>	1081 ± 79	525 ± 26	<b>713 ± 26</b>
Centipede	2961 ± 379	<b>3694 ± 224</b>	4759 ± 257	<b>7525 ± 769</b>
CrazyClimber	86764 ± 3568	<b>103588 ± 2871</b>	71144 ± 11060	<b>99810 ± 2487</b>
Gravitar	371 ± 47	<b>442 ± 67</b>	269 ± 39	<b>332 ± 61</b>
Qbert	4352 ± 128	<b>4412 ± 282</b>	2827 ± 1927	<b>4020 ± 2415</b>
MsPacman	837 ± 62	<b>1204 ± 86</b>	704 ± 41	<b>1011 ± 52</b>
NameThisGame	<b>5665 ± 280</b>	5423 ± 63	2681 ± 143	<b>5187 ± 247</b>
UpNDown	58289 ± 21226	<b>126830 ± 27534</b>	8815 ± 1395	<b>73490 ± 33553</b>

The gradients for each part are:

$$\frac{\partial L_{a2c}(\pi^{(i)}, A^{(i)})}{\partial z_y} = \begin{cases} A^{(i)}(\pi^{(i)} - 1), & \text{if } i = y \\ A^{(i)}\pi^{(y)}, & \text{if } i \neq y \end{cases} \quad (11)$$

$$\frac{\partial L_{ra2c}(\pi^{(i)}, A^{(i)})}{\partial z_y} = \begin{cases} -A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1), & \text{if } i = y \text{ and } A^{(i)} > 0 \\ -A^{(i)}Z\pi^{(y)}\pi^{(i)}, & \text{if } i \neq y \text{ and } A^{(i)} > 0 \end{cases} \quad (12)$$

Thus, the SA2C loss gradient is:

$$\frac{\partial L_{sa2c}}{\partial z_y} = \begin{cases} \underbrace{A^{(i)}(\pi^{(i)} - 1) - A^{(i)}Z\pi^{(i)}(\pi^{(i)} - 1)}_{\substack{\nabla L_{a2c} < 0 & \nabla L_{ra2c} < 0}}, & \text{if } i = y \text{ and } A^{(i)} > 0 \\ \underbrace{A^{(i)}\pi^{(y)} - A^{(i)}Z\pi^{(y)}\pi^{(i)}}_{\substack{\nabla L_{a2c} > 0 & \nabla L_{ra2c} > 0}}, & \text{if } i \neq y \text{ and } A^{(i)} > 0 \end{cases} \quad (13)$$

For both cases, the gradient directions of the RL (A2C) loss and the reverse RL (RA2C) loss are aligned. When  $i = y$  and  $A^{(i)} > 0$ , the gradient of the RA2C loss is  $-A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1)$ , reaching its maximum magnitude at  $\pi^{(y)} = 0.5$  as a parabolic function. This means that the accelerator helps the probability  $\pi^{(i)}$  increase most rapidly when the action to take is ambiguous. When  $i \neq y$  and  $A^{(i)} > 0$ , the probability of actions other than  $a^{(i)}$  is reduced, and this reduction is influenced by the confidence of both  $\pi^{(i)}$  and  $\pi^{(y)}$ . Specifically, the gradient of the RA2C loss is  $-A^{(i)}Z\pi^{(y)}\pi^{(i)}$ . When both  $\pi^{(i)}$  and  $\pi^{(y)}$  are 0.5, representing the most ambiguous predictions, the accelerator aids the A2C loss in reducing  $\pi^{(y)}$  most effectively. Thus, the RA2C loss helps deviate from ambiguous predictions as an accelerator. SPPO’s loss gradients are also aligned like SA2C and follow the same mechanism (See Appendix B.2).

## 5 EXPERIMENTS

To validate the effectiveness of our algorithm, we conduct experiments on various tasks and models of different scales. First, we experiment on Atari games (Mnih et al., 2013) featuring discrete action spaces (Section 5.1), as well as MuJoCo benchmark tasks (Todorov et al., 2012) and Box2D tasks (Catto, 2011) (Section 5.2) with continuous action spaces using Stable-Baselines3 (Raffin et al., 2021). In these control tasks, we also create a variant of each that introduces reward noise, hypothesizing that it will create more confusion in advantage prediction. SPPO performs better than SA2C for various reverse RL loss hyperparameters  $\beta$ . We also evaluate our method on IMDB and TL;DR datasets using TRIL Chang et al. (2023) to determine whether our approach is practical for LLM tasks. We primarily present the experimental results for PPO in the main paper. In the latter part of this section, we analyze why our method works better with PPO than A2C (Section 5.4), conduct hyperparameter sensitivity tests, and examine the training cost (Section 5.5).

Table 2: Mean final scores and standard errors (over the last 10 episodes) of PPO and SPPO on Atari games, without and with binary symmetric channel (BSC) noise with a crossover probability of 0.1 across 5 seeds. To leverage the reverse RL loss, we discretize the continuous action space. DPPO is added as another baseline ( $\alpha = 1.0, \beta = 0.0$ ), and DSPPO is our proposed method. Full results can be found in Table 12 and 13.

$\epsilon \sim \mathcal{N}(0, 0.05^2)$	<b>Ant</b>	<b>Hopper</b>	<b>HalfCheetah</b>	<b>HumanoidStandup</b>
PPO	601 $\pm$ 47	1936 $\pm$ 147	2068 $\pm$ 208	80945 $\pm$ 2130
DPPO	1897 $\pm$ 86	2153 $\pm$ 106	2722 $\pm$ 188	<b>146038 <math>\pm</math> 1841</b>
<b>DSPPO</b>	<b>2095 <math>\pm</math> 102</b>	<b>2333 <math>\pm</math> 109</b>	<b>3118 <math>\pm</math> 195</b>	145974 $\pm$ 2520
	<b>Walker2d</b>	<b>Swimmer</b>	<b>BipedalWalker</b>	<b>LunarLanderContinuous</b>
PPO	1270 $\pm$ 107	44 $\pm$ 3.0	158 $\pm$ 15.2	181 $\pm$ 13.8
DPPO	3419 $\pm$ 100	57 $\pm$ 3.6	<b>274 <math>\pm</math> 7.1</b>	281 $\pm$ 5.7
<b>DSPPO</b>	<b>3523 <math>\pm</math> 129</b>	<b>72 <math>\pm</math> 5.1</b>	267 $\pm$ 8.8	<b>294 <math>\pm</math> 3.3</b>

## 5.1 DISCRETE ACTION SPACE TASKS

We first conduct experiments on Atari games (Mnih et al., 2016) that the action spaces are discrete to evaluate SPPO and SA2C. We primarily select 22 games based on the reported score for A2C in Schulman et al. (2017), focusing on games where the A2C scores are not close to 0, as this allows us to demonstrate meaningful score changes.

The reward functions for Atari games only output 0 or 1 and are well-defined. To introduce some reward noise, we flip the reward from 0 to 1 or from 1 to 0 with a probability of 10%. We denote this noise setting as a Binary Symmetric Channel (BSC). This setting is analogous to a potential problem in ranking-based feedback (Ouyang et al., 2022) from humans or AI, where evaluators may have different preferences, resulting in reversed scores. We observe that SA2C shows marginal improvements (Table 8), with a narrow range of effective hyperparameter  $\beta$  values. In contrast, SPPO performs well in both noise-free and noisy environments (See Section 5.4 for discussion). Table 1 presents partial results, while the complete results for SPPO, including training curves (Figure C.3), can be found in Table 11. SPPO achieves 16 out of 22 wins in noise-free settings and 19 out of 22 wins in noisy settings.

## 5.2 CONTINUOUS ACTION SPACE TASKS

Next, we perform experiments on MuJoCo benchmark (Todorov et al., 2012) and Box2D (Catto, 2011) continuous action space environments. To utilize the reverse RL loss, we need other action probabilities for a sampled action probability. However, conventional RL uses a multivariate Gaussian distribution as a policy, so it cannot provide the other action probabilities. Thus, we discretize the continuous action space (Tang & Agrawal, 2020), naming these methods DA2C and DPPO, and add them as additional baseline comparisons.

Note that discretizing the continuous action space generally works better than the original RL methods like A2C and PPO for these tasks if the continuous action space is discretized with a sufficient number of bins. This discretized distribution can represent more complex distributions than a diagonal Gaussian distribution (where the covariance is diagonal). We apply the reverse RL loss to both DA2C and DSPPO.

Since the reward functions in these environments are highly engineered, we perturb the reward function with Gaussian noise with a mean of 0 and a standard deviation of 0.05. Table 2 shows partial results for SPPO under noise settings. The full experiment results are in Table 12 and 13. Similar to the Atari game results, SA2C without noise shows tied performance in the noiseless setting, and improvements when the reward noise is introduced. SPPO consistently shows robust performance gains across a wide range of  $\beta$  values for both settings.



Table 3: RM Score indicates the reward model score, Perplexity measures the uncertainty of the model, and Win Rate is judged by GPT-4 Turbo by comparing the generated output and reference text. We use 4 different random seeds for each task.

	IMDB Sentiment		TL;DR Summarization		
	RM Score ( $\uparrow$ )	Perplexity ( $\downarrow$ )	RM Score ( $\uparrow$ )	Perplexity ( $\downarrow$ )	Win Rate ( $\uparrow$ )
SFT	0.54 $\pm$ 0.00	33.02 $\pm$ 0.09	5.83 $\pm$ 0.02	18.35 $\pm$ 0.02	42.00 $\pm$ 2.58
PPO	0.89 $\pm$ 0.02	41.09 $\pm$ 0.43	5.94 $\pm$ 0.08	19.08 $\pm$ 0.17	43.25 $\pm$ 3.82
<b>SPPO</b>	<b>0.92 <math>\pm</math> 0.01</b>	40.60 $\pm$ 0.44	<b>6.13 <math>\pm</math> 0.02</b>	19.27 $\pm$ 0.21	<b>52.50 <math>\pm</math> 2.40</b>

### 5.3 RLHF TASKS

The final tasks are RLHF tasks to determine if our method is applicable to large language models. The first task is IMDB positive sentiment. The objective of the IMDB task is to generate positive sentiment continuations for movie reviews (Maas et al., 2011). The sentiment classifier (Sanh et al., 2019) is used as a reward model to evaluate how positive a provided text is. The base policy is GPT-2 (Radford et al., 2019), which we fine-tune using PPO or SPPO. We evaluate this model based on the reward score and perplexity. SPPO shows improvement in both reward score and perplexity compared to PPO.

The second RLHF task is TL;DR summarization (Völske et al., 2017). The objective is to summarize a post from Reddit. The reward model is a fine-tuned GPT-J (Wang & Komatsuzaki, 2021) with LoRA adapters (Hu et al., 2021) by Chang et al. (2023). The training dataset for this reward model is the filtered dataset with additional human preference data used in Stiennon et al. (2020). The base policy model is an open-source GPT-J model (CarperAI/openai\_summarize\_tldr\_sft) with added LoRA adapters. Note that the open-source GPT-J model is often outputs empty summarizations for most evaluation data. Therefore, we report results after 10 epochs of RL updates as an alternative to SFT, as it begins to consistently summarize posts. We evaluate SPPO based on the reward score, perplexity, and win rate. This win rate is judged by GPT-4 Turbo (OpenAI, 2024) (gpt-4-turbo-2024-04-09) by comparing the generated output and reference text. Even though the perplexity of SPPO is slightly higher than that of PPO, there is an improvement in the reward score and a significantly increased win rate.

In the introduction section, we mention that RLHF or RLAIFF have additional errors due to a trained reward model. We check whether the trained reward model used in TL;DR has reward prediction errors. Figure 4.1 shows a dramatic example: the generated summary sample (left) and the middle sample were both *empty*, but the two rewards show a huge gap. The middle sample even scores (6.66) better than those learned with an SPPO score (6.13). Wrong summaries, like *empty*, can score higher than a summarized text (right). These cases are observed very often. This makes the RL training procedure more noisy and means that the sign of advantage changes depending on how the batch is composed. The full text for these samples can be found in Appendix D.

### 5.4 WHY SPPO WORKS BETTER THAN SA2C

The motivation for using the reverse RL loss is to address the issue of ambiguity in advantage predictions (Section 4.3). We hypothesize that the PPO advantage prediction (sign) is less consistent than in A2C during policy updates, but this does not mean that PPO is worse than A2C. There are two main reasons why consistency is not maintained. First, PPO has improved sample efficiency compared to A2C, but after the first epoch, subsequent updates become off-policy, affecting advantage estimates. Second, PPO often uses advantage normalization to restrict large advantage values from being involved with policy updates to stabilize the learning process. In addition, PPO often uses smaller mini-batch sizes (e.g., 64), whereas A2C uses the entire dataset for policy updates. Many popular RL code baselines, such as Stable Baselines3 (Raffin et al., 2021), RL4LMs (Ramamurthy et al., 2023), TRL (von Werra et al., 2020), and TRLX (Havrilla et al., 2023) use PPO advantage normalization by default, whereas A2C does not. Our experiments on the usefulness of advantage normalization also show that the performance increase in IMDB is greater than the performance decrease in TL;DR (Appendix 14).

Table 4: Percentage improvement of SPPO over PPO. The percentage improvements are computed across 22 Atari games. We simply fix  $\alpha = 0.5$  to reduce the total loss magnitude and vary  $\beta$  to control the impact of the reverse RL loss. We exclude very large improvements (e.g., 2000%) from calculating the average. This large improvements result from PPO’s significant learning failures.

$\alpha = 0.5$ is fixed	$\beta = 0.5$	$\beta = 1.0$	$\beta = 5.0$	$\beta = 10.0$	$\beta = 25.0$
SPPO under 0% noise	7.83%	10.15%	24.98%	21.52%	18.92%
SPPO under 10% noise	1.74%	21.89%	148.46%	166.73%	136.50%

We examine the ratio of advantage sign changes before and after normalization for PPO in Atari games and MuJoCo tasks (Figure 2). This ratio varies across different environments. The advantage sign changes usually exceed 5% for Atari games and 10% for MuJoCo and Box2D environments. These changes introduce the confusion, which makes the reverse RL loss more effective for PPO. This observation aligns with our motivation for using symmetric RL loss to handle noisy data, similar to how it is addressed in noisy classification tasks in supervised learning.

Additionally, since A2C uses the entire dataset (rather than using advantage normalization with small batches) for the policy updates, it introduces less confusion in advantage prediction. As a result, SA2C demonstrates performance comparable to A2C in settings without reward noise (Table 8 and 9), and improvements in settings with reward noise (Table 8 and 10), where advantage estimation is more likely to be confused.

## 5.5 HYPERPARAMETERS AND TRAINING COST

Although the symmetric RL loss introduces three additional hyperparameters (Equation 4.2):  $\alpha$ ,  $\beta$ , and  $Z$ , we simply fix  $\alpha = 0.5$  in all experiments to reduce the overall magnitude of the symmetric loss. Additionally, since  $\beta$  and  $Z$  are constants that are multiplied together, we can fix one and adjust the other. For example,  $(\beta = 1.0, Z = -1.0)$  and  $(\beta = 10.0, Z = -0.1)$  yield the same results. In our experiments, we fix  $Z = -1$  and adjust  $\beta$  to determine the influence of the reverse RL loss.

We test the sensitivity of  $\beta$  for SPPO on Atari games with and without noise in the rewards. Table 5 presents the percentage improvements compared to PPO. We exclude excessively large improvements (e.g., 2000%) to avoid skewing the average. These significant improvements typically result from PPO’s training failure, while SPPO remains stable (Gopher and WizardOfWor in Figure C.3). Fixing  $\alpha = 0.5$  and  $Z = -1$ , we vary  $\beta$  and observe consistent improvements, demonstrating SPPO’s robustness across hyperparameters. Also, we use the default values of Stable Baselines3 (Raffin et al., 2021) for the other RL hyperparameters; more details can be found in Appendix C.1.

The symmetric RL loss introduces the reverse RL loss term, which is essentially another form of cross-entropy that does not significantly increase training time. In practice, there is no increase in training time for the continuous tasks discussed in Section 5.2 and the LLM tasks in Section 5.3, and a 10–20% increase for the Atari games in Section 5.1.

## 6 CONCLUSION AND FUTURE WORK

We present Symmetric RL loss, inspired by Symmetric Cross Entropy (SCE) (Wang et al., 2019) from supervised learning, to enhance the robustness of RL. SCE leverages reverse information to handle noisy data, which we adapt to RL algorithms like A2C and PPO, resulting in SA2C and SPPO. We test SA2C and SPPO on various discrete and continuous action space tasks and further evaluate SPPO on RLHF tasks, including IMDB positive sentiment and TL;DR summarization. Our results show that SPPO consistently outperforms PPO. We argue this is mainly due to PPO’s off-policy parts and advantage normalization with small batch sizes, which lead to advantage sign changes (confusion). The Symmetric RL loss for SPPO alleviates this training difficulty.

While we only propose the Symmetric RL loss specifically for A2C and PPO, exploring its integration into other RL algorithms, such as DQN or SAC, is an intriguing future direction. Additionally, developing more diverse reverse RL loss functions, like the Normalized Loss Functions (Ma et al., 2020) proposed after SCE, is also an interesting direction for future work.

## REFERENCES

- 540  
541  
542 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,  
543 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learn-  
544 ing from human feedback in llms, 2024.
- 545  
546 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones,  
547 Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Ols-  
548 son, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-  
549 Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse,  
550 Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mer-  
551 cado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna  
552 Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Con-  
553 erly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario  
554 Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai:  
555 Harmlessness from ai feedback, 2022.
- 556  
557 Ju-Seung Byun and Andrew Perrault. Normality-guided distributional reinforcement learning for  
558 continuous control, 2024. URL <https://arxiv.org/abs/2208.13125>.
- 559  
560 Ju-Seung Byun, Jiyun Chun, Jihyung Kil, and Andrew Perrault. Ares: Alternating reinforcement  
561 learning and supervised fine-tuning for enhanced multi-modal chain-of-thought reasoning through  
562 diverse ai feedback, 2024. URL <https://arxiv.org/abs/2407.00087>.
- 563  
564 Erin Catto. Box2d, a 2d physics engine for games, 2011. URL <http://box2d.org>.
- 565  
566 Souradip Chakraborty, Jiahao Qiu, Hui Yuan, Alec Koppel, Furong Huang, Dinesh Manocha, Am-  
567 rit Singh Bedi, and Mengdi Wang. Maxmin-rlhf: Towards equitable alignment of large language  
568 models with diverse human preferences, 2024.
- 569  
570 Jonathan D Chang, Kianté Brantley, Rajkumar Ramamurthy, Dipendra Misra, and Wen Sun.  
571 Tril: Transformers reinforcement and imitation learning library. [https://github.com/](https://github.com/Cornell-RL/tril)  
572 Cornell-RL/tril, 2023.
- 573  
574 Xi Chen, Zhihui Zhu, and Andrew Perrault. The distributional reward critic architecture for  
575 perturbed-reward reinforcement learning, 2024. URL [https://arxiv.org/abs/2401.](https://arxiv.org/abs/2401.05710)  
576 05710.
- 577  
578 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan  
579 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu,  
580 Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pel-  
581 lat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao,  
582 Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin,  
583 Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language  
584 models, 2022.
- 585  
586 Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter,  
587 Ayzan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar,  
588 Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc  
589 Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied  
590 multimodal language model, 2023.
- 591  
592 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model  
593 alignment as prospect theoretic optimization, 2024.
- 594  
595 Aritra Ghosh, Himanshu Kumar, and P. S. Sastry. Robust loss functions under label noise for deep  
596 neural networks, 2017.
- 597  
598 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
599 maximum entropy deep reinforcement learning with a stochastic actor, 2018.

- 594 Alexander Havrilla, Maksym Zhuravinskyi, Duy Phung, Aman Tiwari, Jonathan Tow, Stella Bi-  
595 derman, Quentin Anthony, and Louis Castricato. trIX: A framework for large scale reinforce-  
596 ment learning from human feedback. In *Proceedings of the 2023 Conference on Empirical*  
597 *Methods in Natural Language Processing*, pp. 8578–8595, Singapore, December 2023. Asso-  
598 ciation for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.530. URL <https://aclanthology.org/2023.emnlp-main.530>.  
599
- 600 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and  
601 Weizhu Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685,  
602 2021. URL <https://arxiv.org/abs/2106.09685>.  
603
- 604 Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and  
605 predicting hospital readmission, 2020.
- 606 Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot  
607 planners: Extracting actionable knowledge for embodied agents, 2022.  
608
- 609 Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton  
610 Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif: Scaling rein-  
611 forcement learning from human feedback with ai feedback, 2023a.
- 612 Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel,  
613 Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human  
614 feedback, 2023b.
- 615 Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an  
616 ai chatbot for medicine. *New England Journal of Medicine*, 388(13):1233–1239, 2023c.  
617 doi: 10.1056/NEJMSr2214184. URL [https://www.nejm.org/doi/full/10.1056/](https://www.nejm.org/doi/full/10.1056/NEJMSr2214184)  
618 [NEJMSr2214184](https://www.nejm.org/doi/full/10.1056/NEJMSr2214184).  
619
- 620 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan  
621 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step, 2023.
- 622 Xingjun Ma, Hanxun Huang, Yisen Wang, Simone Romano, Sarah M. Erfani, and James Bailey.  
623 Normalized loss functions for deep learning with noisy labels. *CoRR*, abs/2006.13554, 2020.  
624 URL <https://arxiv.org/abs/2006.13554>.  
625
- 626 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher  
627 Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and  
628 Rada Mihalcea (eds.), *Proceedings of the 49th Annual Meeting of the Association for Compu-*  
629 *tational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June  
630 2011. Association for Computational Linguistics. URL [https://aclanthology.org/](https://aclanthology.org/P11-1015)  
631 [P11-1015](https://aclanthology.org/P11-1015).  
632
- 633 Vincent Mai, Kaustubh Mani, and Liam Paull. Sample efficient deep reinforcement learning via  
634 uncertainty estimation, 2022. URL <https://arxiv.org/abs/2201.01666>.
- 635 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
636 Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*,  
637 abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>.  
638
- 639 Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim  
640 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement  
641 learning. *CoRR*, abs/1602.01783, 2016. URL <http://arxiv.org/abs/1602.01783>.
- 642 Johan Obando-Ceron, Marc G. Bellemare, and Pablo Samuel Castro. Small batch deep reinforce-  
643 ment learning, 2023. URL <https://arxiv.org/abs/2310.03882>.  
644
- 645 OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 646 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong  
647 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-  
ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,  
and Ryan Lowe. Training language models to follow instructions with human feedback, 2022.

- 648 Martin L. Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John  
649 Wiley & Sons, 2014.
- 650
- 651 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language  
652 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 653
- 654 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and  
655 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model,  
656 2023.
- 657 Antonin Raffin. `RL baselines3 zoo`. [https://github.com/DLR-RM/  
658 rl-baselines3-zoo](https://github.com/DLR-RM/rl-baselines3-zoo), 2020.
- 659
- 660 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah  
661 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of  
662 Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/  
663 20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 664 Rajkumar Ramamurthy, Prithviraj Ammanabrolu, Kianté Brantley, Jack Hessel, Rafet Sifa, Chris-  
665 tian Bauckhage, Hannaneh Hajishirzi, and Yejin Choi. Is reinforcement learning (not) for natural  
666 language processing: Benchmarks, baselines, and building blocks for natural language policy  
667 optimization, 2023.
- 668
- 669 Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version  
670 of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108, 2019. URL [http://  
671 arxiv.org/abs/1910.01108](http://arxiv.org/abs/1910.01108).
- 672 Tom Schaul, André Barreto, John Quan, and Georg Ostrovski. The phenomenon of policy churn,  
673 2022. URL <https://arxiv.org/abs/2206.00730>.
- 674
- 675 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region  
676 policy optimization. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International  
677 Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp.  
678 1889–1897, Lille, France, 07–09 Jul 2015. PMLR. URL [https://proceedings.mlr.  
679 press/v37/schulman15.html](https://proceedings.mlr.press/v37/schulman15.html).
- 680
- 681 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
682 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL [http://arxiv.org/abs/  
683 1707.06347](http://arxiv.org/abs/1707.06347).
- 684
- 685 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-  
686 dimensional continuous control using generalized advantage estimation, 2018.
- 687
- 688 Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,  
689 Dario Amodei, and Paul F. Christiano. Learning to summarize from human feedback. *CoRR*,  
690 abs/2009.01325, 2020. URL <https://arxiv.org/abs/2009.01325>.
- 691
- 692 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press,  
693 Cambridge, MA, USA, 2018a. ISBN 978-0262039246.
- 694
- 695 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018b.
- 696
- 697 Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient  
698 methods for reinforcement learning with function approximation. In *Advances in neural informa-  
699 tion processing systems*, pp. 1057–1063, 2000.
- 700
- 701 Yunhao Tang and Shipra Agrawal. Discretizing continuous action space for on-policy optimization,  
2020.
- 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,  
2012. doi: 10.1109/IROS.2012.6386109.

- 702 Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-  
703 learning, 2015.  
704
- 705 Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. TL;DR: Mining Reddit to  
706 learn automatic summarization. In Lu Wang, Jackie Chi Kit Cheung, Giuseppe Carenini,  
707 and Fei Liu (eds.), *Proceedings of the Workshop on New Frontiers in Summarization*, pp. 59–  
708 63, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi:  
709 10.18653/v1/W17-4508. URL <https://aclanthology.org/W17-4508>.
- 710 Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan  
711 Lambert, and Shengyi Huang. Trl: Transformer reinforcement learning. <https://github.com/huggingface/trl>, 2020.  
712
- 713 Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language  
714 Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.  
715
- 716 Jingkan Wang, Yang Liu, and Bo Li. Reinforcement learning with perturbed rewards. *CoRR*,  
717 abs/1810.01032, 2018. URL <http://arxiv.org/abs/1810.01032>.  
718
- 719 Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang  
720 Sui. Math-shepherd: Verify and reinforce llms step-by-step without human annotations, 2024.  
721
- 722 Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross  
723 entropy for robust learning with noisy labels. *CoRR*, abs/1908.06112, 2019. URL <http://arxiv.org/abs/1908.06112>.  
724
- 725 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc  
726 Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models,  
727 2023.  
728
- 729 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement  
730 learning. *Machine learning*, 8(3-4):229–256, 1992.
- 731 Zhilu Zhang and Mert R. Sabuncu. Generalized cross entropy loss for training deep neural networks  
732 with noisy labels. *CoRR*, abs/1805.07836, 2018. URL <http://arxiv.org/abs/1805.07836>.  
733

## 734 A GRADIENT OF RL LOSS AND REVERSE RL LOSS

735  
736  
737 Suppose there exist  $k$  actions, and  $a^{(i)}$  indicates the  $i^{\text{th}}$  action. Let  $\pi_{\theta}^{(i)} = \pi_{\theta}(a^{(i)}|s)$  denote the  
738 policy for a state  $s$ . The set  $\pi_{\theta}(s) = \{\pi_{\theta}^{(1)}, \pi_{\theta}^{(2)}, \dots, \pi_{\theta}^{(k)}\}$  represents the possible action probabil-  
739 ities set for  $s$ .  $A^{(i)}$  indicates the corresponding advantage of the sampled action  $a^{(i)}$  for  $s$ .  $Z < 0$   
740 is a constant used in the reverse RL loss to handle the computational issue where  $\log 0 = -\infty$ . For  
741 simplicity of notation, we drop  $\theta$ ,  $s$ , and  $a$  from the policy  $\pi$ . Note that  $A^{(i)}$  and  $Z$  are not involved  
742 with the gradient as they are constants with respect to  $\theta$ .  
743  
744

### 745 A.1 A2C LOSS

746  
747 The derivation of the A2C loss  $L_{\text{a2c}}$  with respect to logits  $z$  is presented as follows:

748 For  $i = y$ ,

$$\begin{aligned}
 \frac{\partial \pi^{(i)}}{\partial z_y} &= \frac{\partial}{\partial z_y} \frac{e^{z_i}}{\sum_{w=1}^k e^{z_w}} \\
 &= \frac{e^{z_i} \sum_{w=1}^k e^{z_w} - e^{z_i} e^{z_i}}{(\sum_{w=1}^k e^{z_w})^2} \\
 &= \pi^{(i)}(1 - \pi^{(i)})
 \end{aligned}
 \tag{14}$$

756 For  $i \neq y$ ,

$$\begin{aligned}
 757 \frac{\partial \pi^{(i)}}{\partial z_y} &= \frac{\partial}{\partial z_y} \frac{e^{z_i}}{\sum_{w=1}^k e^{z_w}} \\
 758 &= -\frac{e^{z_i} e^{z_y}}{(\sum_{w=1}^k e^{z_w})^2} \\
 759 &= -\pi^{(i)} \pi^{(y)}
 \end{aligned} \tag{15}$$

764 The sample-wise A2C loss is:

$$765 L_{\text{a2c}}(\pi^{(i)}, A^{(i)}) = -A^{(i)} \log \pi^{(i)} \tag{16}$$

767 For  $i = y$ ,

$$\begin{aligned}
 768 \frac{\partial L_{\text{a2c}}}{\partial z_y} &= \frac{\partial}{\partial z_y} -A^{(i)} \log \pi^{(i)} \\
 769 &= -A^{(i)} \frac{\partial}{\partial z_y} \log \pi^{(i)} \\
 770 &= -\frac{A^{(i)} \partial \pi^{(i)}}{\pi^{(i)}} \\
 771 &= A^{(i)} (\pi^{(i)} - 1) \text{ by (14)}
 \end{aligned} \tag{17}$$

777 For  $i \neq y$ ,

$$\begin{aligned}
 778 \frac{\partial L_{\text{a2c}}}{\partial z_y} &= \frac{\partial}{\partial z_y} -A^{(i)} \log \pi^{(i)} \\
 779 &= -A^{(i)} \frac{\partial}{\partial z_y} \log \pi^{(i)} \\
 780 &= -\frac{A^{(i)} \partial \pi^{(i)}}{\pi^{(i)}} \\
 781 &= A^{(i)} \pi^{(y)} \text{ by (15)}
 \end{aligned} \tag{18}$$

787 In summary, we have the following form for  $L_{\text{a2c}}(\pi^{(i)}, A^{(i)})$ :

$$788 \frac{\partial L_{\text{a2c}}(\pi^{(i)}, A^{(i)})}{\partial z_y} = \begin{cases} A^{(i)} (\pi^{(i)} - 1), & \text{if } i = y \\ A^{(i)} \pi^{(y)}, & \text{if } i \neq y \end{cases} \tag{19}$$

## 791 A.2 REVERSE A2C LOSS

792 The derivation of the reverse A2C loss  $L_{\text{ra2c}}$  with respect to logits  $z$  is presented as follows:

$$793 L_{\text{ra2c}}(\pi^{(i)}, A^{(i)}) = \begin{cases} \sum_{j \in [k] \setminus \{i\}} -\pi^{(j)} A^{(i)} Z, & \text{if } A^{(i)} > 0 \\ \sum_{j \in [k] \setminus \{i\}} \pi^{(j)} A^{(i)} Z, & \text{if } A^{(i)} < 0 \end{cases} \tag{20}$$

798 For  $i = y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned}
 800 \frac{\partial L_{\text{ra2c}}}{\partial z_y} &= \frac{\partial}{\partial z_y} \sum_{j \in [k] \setminus \{i\}} -\pi^{(j)} A^{(i)} Z \\
 801 &= -A^{(i)} Z \sum_{j \in [k] \setminus \{i\}} \frac{\partial \pi^{(j)}}{\partial z_y} \\
 802 &= -A^{(i)} Z \sum_{j \in [k] \setminus \{i\}} -\pi^{(j)} \pi^{(y)} \text{ by (15)} \\
 803 &= A^{(i)} Z \pi^{(y)} (1 - \pi^{(i)}) \\
 804 &= -A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1)
 \end{aligned} \tag{21}$$

810 For  $i \neq y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned}
811 \quad \frac{\partial L_{\text{ra2c}}}{\partial z_y} &= \frac{\partial}{\partial z_y} \sum_{j \in [k] \setminus \{i\}} -\pi^{(j)} A^{(i)} Z \\
812 \quad &= -A^{(i)} Z \sum_{j \in [k] \setminus \{i\}} \frac{\partial \pi^{(j)}}{\partial z_y} \\
813 \quad &= -A^{(i)} Z \left( \sum_{j \in [k]} \frac{\partial \pi^{(j)}}{\partial z_y} - \frac{\partial \pi^{(i)}}{\partial z_y} \right) \\
814 \quad &= A^{(i)} Z \left( \sum_{j \in [k]} -\pi^{(j)} \pi^{(y)} + \pi^{(y)} \pi^{(y)} + \pi^{(y)} (1 - \pi^{(y)}) - \pi^{(i)} \pi^{(y)} \right) \text{ by (14) and (15)} \\
815 \quad &= -A^{(i)} Z \pi^{(y)} \pi^{(i)} \\
816 \quad & \\
817 \quad & \\
818 \quad & \\
819 \quad & \\
820 \quad & \\
821 \quad & \\
822 \quad & \\
823 \quad & \\
824 \quad & \\
825 \quad & \\
826 \quad & \\
827 \quad & \tag{22}
\end{aligned}$$

828 For  $i = y$  and  $A^{(i)} < 0$ , the only difference from Equation 21 is the negative sign, thus:

$$\begin{aligned}
829 \quad \frac{\partial L_{\text{ra2c}}}{\partial z_y} &= A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1) \text{ by (21)} \\
830 \quad & \\
831 \quad & \tag{23} \\
832 \quad &
\end{aligned}$$

833 For  $i \neq y$  and  $A^{(i)} < 0$ , the only difference from Equation 22 is the negative sign, thus:

$$\begin{aligned}
834 \quad \frac{\partial L_{\text{ra2c}}}{\partial z_y} &= A^{(i)} Z \pi^{(y)} \pi^{(i)} \text{ by (22)} \\
835 \quad & \\
836 \quad & \tag{24} \\
837 \quad &
\end{aligned}$$

838 In summary, we have the following form for  $L_{\text{a2c}}(\pi^{(i)}, A^{(i)})$ :

$$\begin{aligned}
839 \quad \frac{\partial L_{\text{ra2c}}(\pi^{(i)}, A^{(i)})}{\partial z_y} &= \begin{cases} -A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1), & \text{if } i = y \text{ and } A^{(i)} > 0 \\ -A^{(i)} Z \pi^{(y)} \pi^{(i)}, & \text{if } i \neq y \text{ and } A^{(i)} > 0 \\ A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1), & \text{if } i = y \text{ and } A^{(i)} < 0 \\ A^{(i)} Z \pi^{(y)} \pi^{(i)}, & \text{if } i \neq y \text{ and } A^{(i)} < 0 \end{cases} \\
840 \quad & \\
841 \quad & \tag{25} \\
842 \quad & \\
843 \quad & \\
844 \quad &
\end{aligned}$$

### 845 A.3 PPO LOSS

846 The derivation of the PPO loss  $L_{\text{ppo}}$  with respect to the logits  $z$  is presented as follows. The PPO loss includes a clipping function and a minimum operation. When these conditions are not satisfied, there is no gradient.

847 The sample-wise PPO loss is:

$$\begin{aligned}
848 \quad L_{\text{ppo}}(\pi^{(i)}, A^{(i)}, \pi_{\text{old}}^{(i)}) &= -\frac{\pi^{(i)}}{\pi_{\text{old}}^{(i)}} A^{(i)} \\
849 \quad & \\
850 \quad & \tag{26} \\
851 \quad &
\end{aligned}$$

852 For  $i = y$ ,

$$\begin{aligned}
853 \quad \frac{\partial L_{\text{ppo}}}{\partial z_y} &= \frac{\partial}{\partial z_y} -\frac{\pi^{(i)}}{\pi_{\text{old}}^{(i)}} A^{(i)} \\
854 \quad &= -\frac{A^{(i)}}{\pi_{\text{old}}^{(i)}} \frac{\partial \pi^{(i)}}{\partial z_y} \\
855 \quad &= -\frac{A^{(i)} \pi^{(i)} (\pi^{(i)} - 1)}{\pi_{\text{old}}^{(i)}} \text{ by (14)} \\
856 \quad & \\
857 \quad & \tag{27} \\
858 \quad & \\
859 \quad & \\
860 \quad & \\
861 \quad & \\
862 \quad & \\
863 \quad &
\end{aligned}$$



864 For  $i \neq y$ ,

$$\begin{aligned}
866 \frac{\partial L_{\text{ppo}}}{\partial z_y} &= \frac{\partial}{\partial z_y} - \frac{\pi^{(i)}}{\pi_{\text{old}}^{(i)}} A^{(i)} \\
867 &= \frac{A^{(i)}}{\pi_{\text{old}}^{(i)}} \frac{\partial \pi^{(i)}}{\partial z_y} \\
868 &= \frac{A^{(i)} \pi^{(i)} \pi^{(y)}}{\pi_{\text{old}}^{(i)}} \text{ by (15)}
\end{aligned} \tag{28}$$

#### 875 A.4 REVERSE PPO LOSS

877 The derivation of the reverse PPO loss  $L_{\text{rppo}}$  with respect to logits  $z$  is presented as follows. As with  
878 PPO, the reverse PPO loss only considers samples that pass the clipping function and the minimum  
879 operation.

880 From Section A.2, we have the following form for  $L_{\text{rppo}}(\pi^{(i)}, A^{(i)}, \pi_{\text{old}}^{(i)})$ :

$$\frac{\partial L_{\text{rppo}}(\pi^{(i)}, A^{(i)}, \pi_{\text{old}}^{(i)})}{\partial z_y} = \begin{cases} -\frac{A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1)}{\pi_{\text{old}}^{(i)}}, & \text{if } i = y \text{ and } A^{(i)} > 0 \\ -\frac{A^{(i)} Z \pi^{(y)} \pi^{(i)}}{\pi_{\text{old}}^{(i)}}, & \text{if } i \neq y \text{ and } A^{(i)} > 0 \\ \frac{A^{(i)} Z \pi^{(y)} (\pi^{(y)} - 1)}{\pi_{\text{old}}^{(i)}}, & \text{if } i = y \text{ and } A^{(i)} < 0 \\ \frac{A^{(i)} Z \pi^{(y)} \pi^{(i)}}{\pi_{\text{old}}^{(i)}}, & \text{if } i \neq y \text{ and } A^{(i)} < 0 \end{cases} \tag{29}$$

## 890 B GRADIENT ANALYSIS OF RL LOSS AND REVERSE RL LOSS

### 892 B.1 SYMMETRIC A2C GRADIENT ANALYSIS

895 The gradient analysis of the symmetric RL loss follows the SCE analysis. We adopt their analysis  
896 and extend it to cover the RL loss analysis. We set  $\alpha$  and  $\beta$  to 1 for simplicity and evaluate the  
897 gradient direction of both RL and reverse RL losses with respect to the logits  $z$ . We show that the  
898 gradient directions for both types are the same and that the reverse RL loss helps deviate ambiguous  
899 predictions where the probability is around 0.5. We first show how the symmetric A2C (SA2C) loss  
900 behaves. Note that  $Z < 0$  is a constant used in the reverse RL loss to handle  $\log 0 = -\infty$ .

$$902 L_{\text{sa2c}} = L_{\text{a2c}} + L_{\text{ra2c}} \tag{30}$$

904 For  $i = y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned}
906 \frac{\partial L_{\text{sa2c}}}{\partial z_y} &= \frac{\partial L_{\text{a2c}}}{\partial z_y} + \frac{\partial L_{\text{ra2c}}}{\partial z_y} \\
907 &= \underbrace{A^{(i)} (\pi^{(i)} - 1)}_{\nabla L_{\text{a2c}} < 0} - \underbrace{A^{(i)} Z \pi^{(i)} (\pi^{(i)} - 1)}_{\nabla L_{\text{ra2c}} < 0} \text{ by (17) and (21)}
\end{aligned} \tag{31}$$

912 For  $i \neq y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned}
914 \frac{\partial L_{\text{sa2c}}}{\partial z_y} &= \frac{\partial L_{\text{a2c}}}{\partial z_y} + \frac{\partial L_{\text{ra2c}}}{\partial z_y} \\
915 &= \underbrace{A^{(i)} \pi^{(y)}}_{\nabla L_{\text{a2c}} > 0} - \underbrace{A^{(i)} Z \pi^{(y)} \pi^{(i)}}_{\nabla L_{\text{ra2c}} > 0} \text{ by (18) and (22)}
\end{aligned} \tag{32}$$

For  $i = y$  and  $A^{(i)} < 0$ ,

$$\begin{aligned} \frac{\partial L_{\text{sa2c}}}{\partial z_y} &= \frac{\partial L_{\text{a2c}}}{\partial z_y} + \frac{\partial L_{\text{ra2c}}}{\partial z_y} \\ &= \underbrace{A^{(i)}(\pi^{(i)} - 1)}_{\nabla L_{\text{a2c}} > 0} - \underbrace{A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1)}_{\nabla L_{\text{ra2c}} > 0} \quad \text{by (17) and (21)} \end{aligned} \quad (33)$$

For  $i \neq y$  and  $A^{(i)} < 0$ ,

$$\begin{aligned} \frac{\partial L_{\text{sa2c}}}{\partial z_y} &= \frac{\partial L_{\text{a2c}}}{\partial z_y} + \frac{\partial L_{\text{ra2c}}}{\partial z_y} \\ &= \underbrace{A^{(i)}\pi^{(y)}}_{\nabla L_{\text{a2c}} < 0} - \underbrace{A^{(i)}Z\pi^{(y)}\pi^{(i)}}_{\nabla L_{\text{ra2c}} < 0} \quad \text{by (18) and (22)} \end{aligned} \quad (34)$$

For the above cases, the gradient directions of the RL (A2C) loss and the reverse RL (RA2C) loss are the same as SCE gradients. Essentially, the RA2C loss acts as an accelerator. In the case of  $i = y$  and  $A^{(i)} > 0$ , the gradient of the RA2C loss is  $-A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1)$ , with the largest gradient magnitude at  $\pi^{(y)} = 0.5$  as a parabolic function. In other words, the accelerator helps the probability  $\pi^{(i)}$  increase most quickly when it is ambiguous which action to take. In the case of  $i \neq y$  and  $A^{(i)} > 0$ , the probability of other actions except  $a^{(i)}$  is reduced, and this reduction is influenced by the confidence of both  $\pi^{(i)}$  and  $\pi^{(y)}$ . Specifically, the gradient of the RA2C loss is  $-A^{(i)}Z\pi^{(y)}\pi^{(i)}$ . When both  $\pi^{(i)}$  and  $\pi^{(y)}$  are 0.5, indicating the most ambiguous predictions, the accelerator helps the A2C loss reduce  $\pi^{(y)}$  most aggressively.

When  $A^{(i)} < 0$ , the gradient direction is simply reversed. The behavior of the gradient itself remains the same as when  $A^{(i)} > 0$ . In the case of  $i = y$ , RA2C decreases the probability  $\pi^{(y)}$  more when  $\pi^{(y)}$  is around 0.5. For  $i \neq y$ , RA2C helps increase  $\pi^{(y)}$  more when both  $\pi^{(i)}$  and  $\pi^{(y)}$  are ambiguous (both around 0.5).

## B.2 SYMMETRIC PPO GRADIENT ANALYSIS

$$L_{\text{sppo}} = L_{\text{ppo}} + L_{\text{rppo}} \quad (35)$$

For  $i = y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned} \frac{\partial L_{\text{sppo}}}{\partial z_y} &= \frac{\partial L_{\text{ppo}}}{\partial z_y} + \frac{\partial L_{\text{rppo}}}{\partial z_y} \\ &= \underbrace{\frac{A^{(i)}\pi^{(i)}(\pi^{(i)} - 1)}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{ppo}} < 0} - \underbrace{\frac{A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1)}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{rppo}} < 0} \quad \text{by (27) and (29)} \end{aligned} \quad (36)$$

For  $i \neq y$  and  $A^{(i)} > 0$ ,

$$\begin{aligned} \frac{\partial L_{\text{sppo}}}{\partial z_y} &= \frac{\partial L_{\text{ppo}}}{\partial z_y} + \frac{\partial L_{\text{rppo}}}{\partial z_y} \\ &= \underbrace{\frac{A^{(i)}\pi^{(i)}\pi^{(y)}}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{ppo}} > 0} - \underbrace{\frac{A^{(i)}Z\pi^{(y)}\pi^{(i)}}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{rppo}} > 0} \quad \text{by (28) and (29)} \end{aligned} \quad (37)$$

For  $i = y$  and  $A^{(i)} < 0$ ,

$$\begin{aligned} \frac{\partial L_{\text{sppo}}}{\partial z_y} &= \frac{\partial L_{\text{ppo}}}{\partial z_y} + \frac{\partial L_{\text{rppo}}}{\partial z_y} \\ &= \underbrace{\frac{A^{(i)}\pi^{(i)}(\pi^{(i)} - 1)}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{ppo}} < 0} + \underbrace{\frac{A^{(i)}Z\pi^{(y)}(\pi^{(y)} - 1)}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{rppo}} < 0} \quad \text{by (27) and (29)} \end{aligned} \quad (38)$$

For  $i \neq y$  and  $A^{(i)} < 0$ ,

$$\begin{aligned}
 \frac{\partial L_{\text{sppo}}}{\partial z_y} &= \frac{\partial L_{\text{ppo}}}{\partial z_y} + \frac{\partial L_{\text{rppo}}}{\partial z_y} \\
 &= \underbrace{\frac{A^{(i)} \pi^{(i)} \pi^{(y)}}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{ppo}} > 0} + \underbrace{\frac{A^{(i)} Z \pi^{(y)} \pi^{(i)}}{\pi_{\text{old}}^{(i)}}}_{\nabla L_{\text{rppo}} > 0} \quad \text{by (28) and (29)} \tag{39}
 \end{aligned}$$

Basically, the mechanism of RPPO is the same as RA2C, except for  $\pi_{\text{old}}^{(i)}$ , which does not change the gradient sign. Therefore, RPPO also helps PPO deviate from ambiguous predictions, acting as an accelerator.

## C EXPERIMENTAL SETUPS AND RESULTS

### C.1 HYPERPARAMETERS

**Atari games:** We primarily follow the hyperparameter settings of RL Baselines3 Zoo (Raffin, 2020). Most hyperparameter values remain unchanged across environments. Only  $\alpha$  and  $\beta$  are adjusted for the reverse RL loss. For SA2C without noise, we use ( $\alpha = 0.5, \beta = 5.0$ ) for all environments. For SA2C with noise, we use ( $\alpha = 0.5, \beta = 1.0$ ) for (Alien, MsPacman, Qbert, TimePilot, VideoPinball, Assault, Gravitar, StarGunner, UpNDown), and ( $\alpha = 0.5, \beta = 1.0$ ) for others. For SPPO without noise, we use ( $\alpha = 0.5, \beta = 1.0$ ) for all environments. For SPPO with noise, we use ( $\alpha = 0.5, \beta = 10.0$ ) for all environments. We do not use any GPU for Atari games.

Table 5: Hyperparameters for Atari games

	Without Noise	$\epsilon \sim \mathbf{BSC}(0.1)$
<b>SA2C</b>		
- ( $\alpha = 0.5, \beta = 1.0$ )	-	(Alien, Assault, Gravitar, MsPacman, Qbert, StarGunner, TimePilot, UpNDown, VideoPinball)
- ( $\alpha = 0.5, \beta = 5.0$ )	All environments	All others except those mentioned above
<b>SPPO</b>		
- ( $\alpha = 0.5, \beta = 1.0$ )	All environments	-
- ( $\alpha = 0.5, \beta = 10.0$ )	-	All environments

**MuJoCo and Box2D:** We use `n.envs = 4` and `n.steps = 8` for A2C and SA2C. We follow Stable-Baselines3’s default hyperparameters (Raffin et al., 2021) for other settings. Only  $\alpha$  and  $\beta$  are adjusted for the reverse RL loss. For table visibility, let  $\{\text{Ant} = 1, \text{BipedalWalker} = 2, \text{HalfCheetah} = 3, \text{Hopper} = 4, \text{HumanoidStandup} = 5, \text{InvertedDoublePendulum} = 6, \text{LunarLanderContinuous} = 7, \text{Swimmer} = 8, \text{Walker2d} = 9\}$ . We do not use any GPU for these tasks.

Table 6: Hyperparameters for MuJoCO and Box2D environments

	Without Noise	$\epsilon \sim \mathcal{N}(0, 0.05^2)$
<b>SA2C</b>		
- ( $\alpha = 0.5, \beta = 0.2$ )	(1, 4, 8, 9)	(1)
- ( $\alpha = 0.5, \beta = 0.5$ )	(2, 5)	-
- ( $\alpha = 0.5, \beta = 5.0$ )	(3, 6, 7)	(7)
- ( $\alpha = 0.5, \beta = 10.0$ )	-	(2, 3, 4, 5, 6, 8, 9)
- $Z = -1$	All environments	All environments
- (timesteps= $2e6$ )	All environments	All environments
- (Number of bins= 11)	All environments	All environments
<b>SPPO</b>		
- ( $\alpha = 0.5, \beta = 20.0$ )	All environments	(1, 7)
- ( $\alpha = 0.5, \beta = 25.0$ )	-	(2, 3, 5, 6, 9)
- ( $\alpha = 0.5, \beta = 50.0$ )	-	(4, 8)
- $Z = -1$	All environments	All environments
- (timesteps= $1e6$ )	(2, 6)	(2, 6)
- (timesteps= $2e6$ )	(1, 3, 4, 8, 9)	(1, 3, 4, 8, 9)
- (timesteps= $5e6$ )	(9)	(9)
- (timesteps= $1e7$ )	(5)	(5)
- (Number of bins= 11)	All environments	All environments

**IMDB and TL;DR:** We basically use the provided implementation (Chang et al., 2023) and follow their hyperparameters, with the addition of the advantage normalization step for PPO. The scripts used in our experiments are available in the code repository for further detail. We use a single Nvidia A100 (80GB) for our experiments.

Table 7: Hyperparameters for IMDB positive sentiment and TL;DR summarization

	IMDB	Box2D
PPO		
- model:	GPT-2	GPT-J
- updates:	60	100
- trajectories per update:	112	64
- epochs per update	5	4
- batch size	28	32
- learning rate	5e-6	5e-6
- discount factor	0.99	1.0
- GAE lambda	0.95	0.95
- clip range	0.2	0.2
SPPO	$(\alpha = 0.5, \beta = 0.4)$	$(\alpha = 0.5, \beta = 0.2)$

## C.2 EXPERIMENTAL RESULTS: A2C AND SA2C

Table 8: Mean final scores and standard errors (over the last 10 episodes) of A2C and SA2C on Atari games, without and with binary symmetric channel (BSC) noise with a crossover probability of 0.1 across 5 seeds.

	Without Noise		$\epsilon \sim \text{BSC}(0.1)$	
	A2C	SA2C	A2C	SA2C
Alien	<b>913 ± 100</b>	771 ± 51	481 ± 72	<b>496 ± 37</b>
Assault	<b>1538 ± 199</b>	1061 ± 41	287 ± 226	<b>399 ± 133</b>
Asterix	2308 ± 86	<b>2377 ± 164</b>	1403 ± 305	<b>1430 ± 208</b>
BeamRider	1121 ± 61	<b>1335 ± 43</b>	<b>1087 ± 339</b>	902 ± 196
Centipede	<b>3588 ± 430</b>	3574 ± 295	3108 ± 243	<b>3540 ± 194</b>
CrazyClimber	98774 ± 2516	<b>99330 ± 4371</b>	93042 ± 8711	<b>97058 ± 6251</b>
DemonAttack	4309 ± 325	<b>5017 ± 625</b>	<b>30 ± 21</b>	19 ± 3
Frostbite	255 ± 2	<b>257 ± 3</b>	241 ± 9	<b>286 ± 48</b>
Gopher	960 ± 80	<b>1036 ± 138</b>	947 ± 91	<b>996 ± 114</b>
Gravitar	143 ± 18	<b>201 ± 16</b>	<b>279 ± 48</b>	183 ± 36
Krull	6387 ± 267	<b>7672 ± 819</b>	<b>7564 ± 486</b>	6337 ± 754
MsPacman	1175 ± 43	<b>1495 ± 104</b>	<b>926 ± 44</b>	916 ± 100
NameThisGame	<b>5945 ± 102</b>	5614 ± 166	2280 ± 257	<b>2372 ± 141</b>
Qbert	1646 ± 240	<b>2103 ± 261</b>	620 ± 96	<b>641 ± 77</b>
Riverraid	4368 ± 582	<b>5461 ± 456</b>	1609 ± 65	<b>2511 ± 190</b>
RoadRunner	14971 ± 1396	<b>18624 ± 1812</b>	<b>5606 ± 1788</b>	3830 ± 1517
Seaquest	836 ± 7	<b>988 ± 92</b>	650 ± 22	<b>653 ± 22</b>
StarGunner	<b>2222 ± 114</b>	1766 ± 120	<b>1194 ± 645</b>	622 ± 54
TimePilot	<b>3992 ± 198</b>	3116 ± 137	2232 ± 259	<b>3288 ± 106</b>
UpNDown	<b>8313 ± 1544</b>	1638 ± 761	4228 ± 1187	<b>7093 ± 2772</b>
VideoPinball	<b>24948 ± 3038</b>	19618 ± 1888	20319 ± 2157	<b>25035 ± 3914</b>
WizardOfWor	<b>824 ± 136</b>	674 ± 125	496 ± 87	<b>752 ± 156</b>
<b>Wins (SA2C)</b>	<b>12 / 22</b>		<b>15 / 22</b>	

Table 9: Mean final scores and standard errors (over the last 10 episodes) of A2C and SA2C on MuJoCo benchmark tasks and Box2D environments without Gaussian noise across 30 seeds.

Without Noise	Ant	Hopper	HalfCheetah	HumanoidStandup
A2C	757 ± 116	1410 ± 112	1393 ± 163	121850 ± 4264
DA2C	2220 ± 96	<b>1944 ± 116</b>	<b>2325 ± 209</b>	152135 ± 3937
<b>DSA2C</b>	<b>2287 ± 94</b>	1797 ± 139	2266 ± 203	<b>159142 ± 129</b>
Walker2d	Swimmer	BipedalWalker	LunarLanderContinuous	
A2C	1348 ± 130	95.8 ± 19.0	124 ± 23	79.0 ± 20.2
DA2C	<b>2131 ± 154</b>	<b>142.4 ± 17.0</b>	234 ± 22	176.7 ± 20.9
<b>DSA2C</b>	1662 ± 164	128.5 ± 16.2	<b>274 ± 16</b>	<b>221.2 ± 10.7</b>
InvertedDoublePendulum				
A2C	1670 ± 500			
DA2C	9139 ± 94			
<b>DSA2C</b>	<b>9145 ± 93</b>			

Table 10: Mean final scores and standard errors (over the last 10 episodes) of A2C and SA2C on MuJoCo benchmark tasks and Box2D environments with Gaussian noise (mean 0 and standard deviation 0.05) across 30 seeds.

$\epsilon \sim \mathcal{N}(0, 0.05^2)$	Ant	Hopper	HalfCheetah	HumanoidStandup
A2C	673 ± 108	1083 ± 92	1610 ± 163	101064 ± 4933
DA2C	1296 ± 80	<b>1323 ± 87</b>	1510 ± 126	126241 ± 3973
<b>DSA2C</b>	<b>1520 ± 83</b>	1307 ± 102	<b>1696 ± 163</b>	<b>128064 ± 4391</b>
Walker2d	Swimmer	BipedalWalker	LunarLanderContinuous	
A2C	786 ± 86	28.9 ± 4.4	158 ± 20	-3.7 ± 15.9
DA2C	<b>1599 ± 138</b>	36.8 ± 4.6	210 ± 21	106 ± 20
<b>DSA2C</b>	1423 ± 129	<b>53.1 ± 7.0</b>	<b>222 ± 20</b>	<b>179 ± 12</b>
InvertedDoublePendulum				
A2C	3852 ± 634			
DA2C	7900 ± 364			
<b>DSA2C</b>	<b>8323 ± 217</b>			

## C.3 EXPERIMENTAL RESULTS: PPO AND SPPO

Table 11: Mean final scores and standard errors (over the last 10 episodes) of PPO and SPPO on Atari games, without and with binary symmetric channel (BSC) noise with a crossover probability of 0.1 across 5 seeds.

	Without Noise		$\epsilon \sim \text{BSC}(0.1)$	
	PPO	SPPO	PPO	SPPO
Alien	<b>1128 ± 105</b>	1081 ± 79	525 ± 26	<b>713 ± 26</b>
Assault	3134 ± 193	<b>3385 ± 214</b>	2327 ± 401	<b>3698 ± 363</b>
Asterix	2599 ± 101	<b>2976 ± 150</b>	1272 ± 106	<b>1739 ± 329</b>
BeamRider	<b>2176 ± 251</b>	1635 ± 404	<b>1828 ± 130</b>	1580 ± 96
Centipede	2961 ± 379	<b>3694 ± 224</b>	4759 ± 257	<b>7525 ± 769</b>
CrazyClimber	86764 ± 3568	<b>103588 ± 2871</b>	71144 ± 11060	<b>99810 ± 2487</b>
DemonAttack	7872 ± 302	<b>7901 ± 455</b>	<b>161 ± 24</b>	132 ± 13
Frostbite	268 ± 5	<b>286 ± 6</b>	<b>509 ± 108</b>	23 ± 16
Gopher	787 ± 48	<b>875 ± 78</b>	478 ± 38	<b>7765 ± 3366</b>
Gravitar	371 ± 47	<b>442 ± 67</b>	269 ± 39	<b>332 ± 61</b>
Krull	6628 ± 417	<b>7578 ± 588</b>	5602 ± 481	<b>9015 ± 381</b>
MsPacman	837 ± 62	<b>1204 ± 86</b>	704 ± 41	<b>1011 ± 52</b>
NameThisGame	<b>5665 ± 280</b>	5423 ± 63	2681 ± 143	<b>5187 ± 247</b>
Qbert	4352 ± 128	<b>4412 ± 282</b>	2827 ± 1927	<b>4020 ± 2415</b>
Riverraid	6128 ± 272	<b>6343 ± 219</b>	2460 ± 127	<b>3998 ± 248</b>
RoadRunner	<b>28382 ± 2254</b>	22562 ± 2875	1204 ± 157	<b>3830 ± 1230</b>
Seaquest	<b>902 ± 2</b>	888 ± 6	652 ± 16	<b>814 ± 15</b>
StarGunner	11848 ± 722	<b>14746 ± 1876</b>	1514 ± 110	<b>23250 ± 6292</b>
TimePilot	<b>3850 ± 151</b>	3548 ± 220	3506 ± 318	<b>3936 ± 420</b>
UpNDown	58289 ± 21226	<b>126830 ± 27534</b>	8815 ± 1395	<b>73490 ± 33553</b>
VideoPinball	22408 ± 4292	<b>29485 ± 2851</b>	31680 ± 2318	<b>37048 ± 6989</b>
WizardOfWor	3186 ± 256	<b>3762 ± 387</b>	940 ± 158	<b>4442 ± 1332</b>
<b>Wins (SPPO)</b>	<b>16 / 22</b>		<b>19 / 22</b>	

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

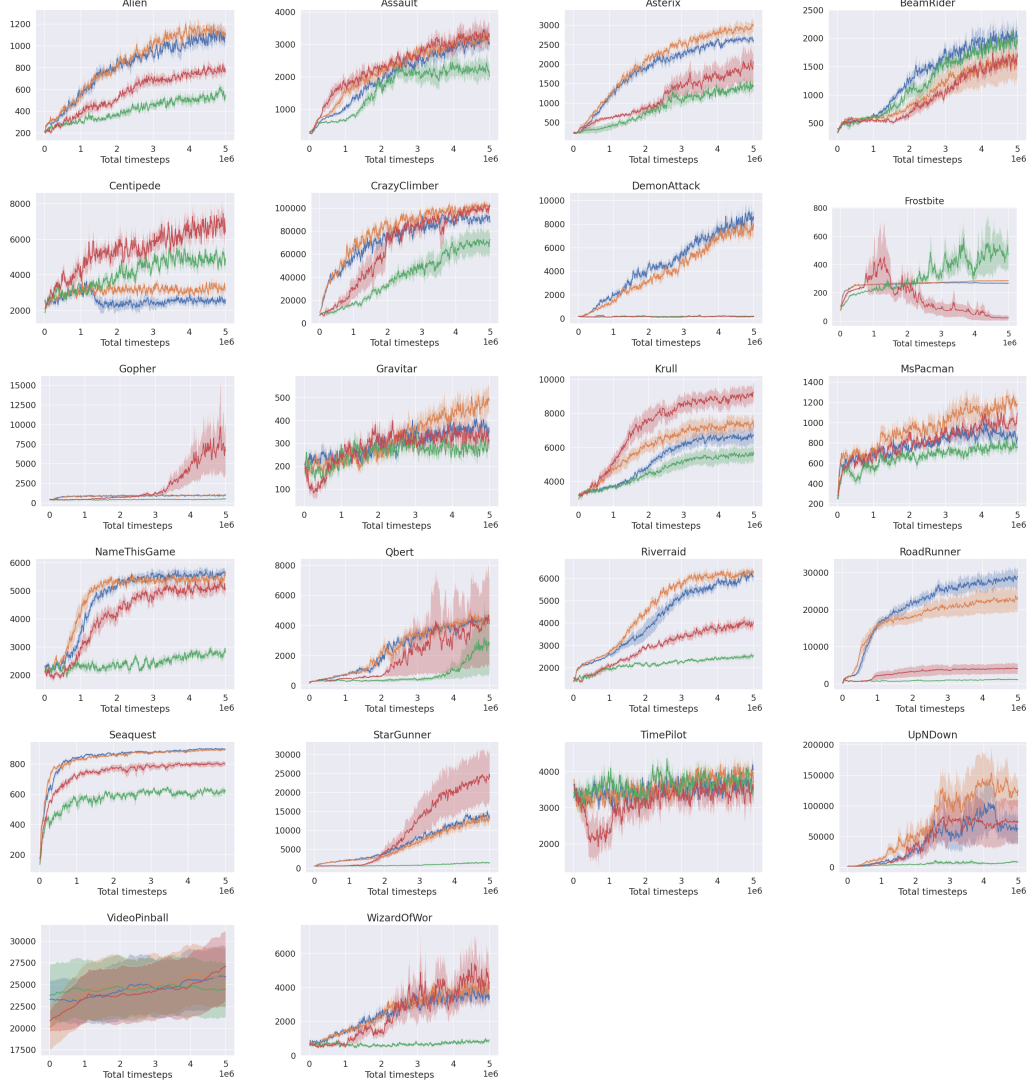


Figure 3: Result of training plots for SPPO and PPO for Atari games. The blue line indicates the original PPO without any added noise, while the orange line represents SPPO without added noise. The green line indicates PPO with 10% noise, and the red line represents SPPO with 10% noise. We fix  $\alpha = 0.5$  for all environments, with  $\beta = 1.0$  for the experiments without noise and  $\beta = 10.0$  for the noise environments.



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

Table 12: Mean final scores and standard errors (over the last 10 episodes) of PPO and SPPO on MuJoCo benchmark tasks and Box2D environments without Gaussian noise across 30 seeds.

Without Noise	Ant	Hopper	HalfCheetah	HumanoidStandup
PPO	2068 ± 166	<b>2875 ± 137</b>	2282 ± 191	93763 ± 3402
DPPO	2735 ± 109	2154 ± 119	3478 ± 279	176320 ± 6538
<b>DSPPO</b>	<b>2885 ± 100</b>	2299 ± 115	<b>4104 ± 258</b>	<b>189301 ± 5915</b>
	Walker2d	Swimmer	BipedalWalker	LunarLanderContinuous
PPO	2793 ± 199	112 ± 5.0	247 ± 8.3	134 ± 10.9
DPPO	4443 ± 119	<b>131 ± 0.3</b>	265 ± 15.5	241 ± 7.7
<b>DSPPO</b>	<b>4587 ± 154</b>	130 ± 0.6	<b>274 ± 6.2</b>	<b>250 ± 6.9</b>
InvertedDoublePendulum				
PPO	7454 ± 394			
DPPO	8928 ± 136			
<b>DSPPO</b>	<b>9015 ± 101</b>			

Table 13: Mean final scores and standard errors (over the last 10 episodes) of PPO and SPPO on MuJoCo benchmark tasks and Box2D environments with Gaussian noise (mean 0 and standard deviation 0.05) across 30 seeds.

$\epsilon \sim \mathcal{N}(0, 0.05^2)$	Ant	Hopper	HalfCheetah	HumanoidStandup
PPO	601 ± 47	1936 ± 147	2068 ± 208	80945 ± 2130
DPPO	1897 ± 86	2153 ± 106	2722 ± 188	<b>146038 ± 1841</b>
<b>DSPPO</b>	<b>2095 ± 102</b>	<b>2333 ± 109</b>	<b>3118 ± 195</b>	145974 ± 2520
	Walker2d	Swimmer	BipedalWalker	LunarLanderContinuous
PPO	1270 ± 107	44 ± 3.0	158 ± 15.2	181 ± 13.8
DPPO	3419 ± 100	57 ± 3.6	<b>274 ± 7.1</b>	281 ± 5.7
<b>DSPPO</b>	<b>3523 ± 129</b>	<b>72 ± 5.1</b>	267 ± 8.8	<b>294 ± 3.3</b>
InvertedDoublePendulum				
PPO	8050 ± 244			
DPPO	8963 ± 100			
<b>DSPPO</b>	<b>9147 ± 61</b>			

## C.4 ON AND OFF ADVANTAGE NORMALIZATION

Table 14: Comparison with and without advantage normalization over 4 different random seeds.

PPO	IMDB	TL;DR
Without $A$ Normalization	$0.77 \pm 0.01$	$6.06 \pm 0.02$
With $A$ Normalization	$0.89 \pm 0.02$	$5.94 \pm 0.08$

## D EXAMPLES OF REWARD MODEL ERRORS

**Warning: This section contains harmful language.**

Table 15: Example showing a trained reward model with errors that are not consistent for empty outputs, and the reward for an empty output is greater than that for a non-empty summarization. [...] indicates omitted content for brevity.

<p><b>Subreddit:</b> r/relationships (Sample ID: 37)</p> <p><b>TITLE:</b> I'm a dumb [21] male and so I'm having a lot of trouble interpreting the signals that this [21] girl may or may not be sending me. A little help please?</p> <p><b>Post:</b> So okay, I'm from New York but I study in Oregon for most of the year. Recently a friend of mine who I was not really close started facebook messaging me, that was about 3 months ago, since then we've talked almost everyday. [...] Is she trying to play hard to get? Am I looking way too into this and maybe she was just occupied that weekend? I really have no idea how to evaluate this. Do any of you guys have any suggestions/ideas?</p> <p><b>Generated Summary:</b> &lt;empty&gt;</p> <p><b>Reward Model Output:</b> 6.66</p>
<p><b>Subreddit:</b> r/relationship_advice (Sample ID: 60)</p> <p><b>TITLE:</b> My bf [23] doesn't speak of his childhood, but I[f22] know he's traumatized.</p> <p><b>Post:</b> We were friends for 10 years, before we got together. He than told me once about his terrible childhood. (He told only 3 of his friends his story) Now we're a couple for quite a few months and well, sometimes there's stuff I know that reminds him of his childhood, but it's like he's forgotten that he had told me. [...] (But striking, the things he thinks are important are always the things his parents should have done, to save him from the traumatizing stuff.)I know he likes to put his problems far away. But on the other hand, I'm his girlfriend now and we're pretty serious, isn't it good to speak about it maybe just once, so he knows I know his secret/won't tell, and most of all, I'm always there for him? What do you think?</p> <p><b>Generated Summary:</b> &lt;empty&gt;</p> <p><b>Reward Model Output:</b> 3.14</p>
<p><b>Subreddit:</b> r/AskReddit (Sample ID: 27)</p> <p><b>TITLE:</b> Dear Reddit, What silly/irrelevant/ridiculous family miscommunications have lead to feuds lasting years?</p> <p><b>Post:</b> My Grandma and my aunt (her daughter-in-aw) haven't spoken to each other in years over a phone that didn't get hung up. My aunt and uncle screen their calls and frequently do not return them-- one time, my grandma called and left a message then thought she hung up the phone. [...] Why continue to hold a silly grudge? To complicate matters further, my grandma has a daughter who lives with her and likes to be in other peoples business-- I think she is also part of the problem here as she won't drop it either. Grandma is innocent but has a daughter and daughter-in-law who won't grow up and drop it</p> <p><b>Generated Summary:</b> Grandma and Aunt haven't spoken in years over a phone that didn't get hung up. Grandma wants to reconcile and clear the air, but Aunt won't go near her, won't let her husband and kids go there, and avoids.</p> <p><b>Reward Model Output:</b> 5.40</p>