

# Learning to Identify Critical States for Reinforcement Learning from Videos

Haozhe Liu<sup>1†</sup>, Mingchen Zhuge<sup>1†</sup>, Bing Li<sup>1✉</sup>, Yuhui Wang<sup>1</sup>, Francesco Faccio<sup>1,2</sup>  
Bernard Ghanem<sup>1</sup>, Jürgen Schmidhuber<sup>1,2,3</sup>

<sup>1</sup>AI Initiative, King Abdullah University of Science and Technology

<sup>2</sup>The Swiss AI Lab IDSIA/USI/SUPSI, <sup>3</sup>NNAISENSE

{haozhe.liu, mingchen.zhuge, bing.li, yuhui.wang,  
francesco.faccio, bernard.ghanem, juergen.schmidhuber}@kaust.edu.sa

## Abstract

Recent work on deep reinforcement learning (DRL) has pointed out that algorithmic information about good policies can be extracted from offline data which lack explicit information about executed actions [50, 51, 35]. For example, videos of humans or robots may convey a lot of implicit information about rewarding action sequences, but a DRL machine that wants to profit from watching such videos must first learn by itself to identify and recognize relevant states/actions/rewards. Without relying on ground-truth annotations, our new method called *Deep State Identifier* learns to predict returns from episodes encoded as videos. Then it uses a kind of mask-based sensitivity analysis to extract/identify important critical states. Extensive experiments showcase our method’s potential for understanding and improving agent behavior. The source code and the generated datasets are available at [Github](#).

## 1. Introduction

In deep reinforcement learning (DRL), the cumulative reward—also known as the return—of an episode is obtained through a long sequence of dynamic interactions between an agent (i.e., a decision-maker) and its environment. In such a setting, the rewards may be sparse and delayed, and it is often unclear which decision points were critical to achieve a specific return.

Several existing methods use the notion of localizing critical states, such as EDGE [21] and RUDDER [1]. These methods typically require explicit action information or policy parameters to localize critical states. This limits their potential applicability in settings like video-based offline RL, where an agent’s actions are often hard to measure, an-

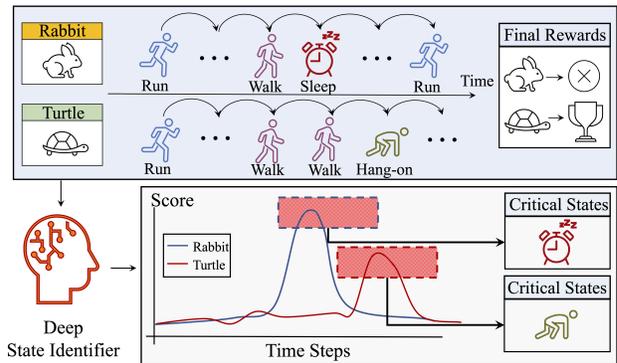


Figure 1. **Motivation of the proposed method.** In the illustrated race between a turtle and a rabbit, the *sleep* state is critical in determining the winner of the race. Our method is proposed to identify such critical states.

notate, or estimate [76, 37]. To avoid this pitfall, in this work, we explicitly study the relationship between sequential visual observations and episodic returns without accessing explicit action information.

Inspired by the existing evidence that frequently only a few decision points are important in determining the return of an episode [1, 13], and as shown in Fig. 1, we focus on identifying the state underlying these critical decision points. However, the problem of directly inferring critical visual input based on the return is nontrivial [13], and compounded by our lack of explicit access to actions or policies during inference. To overcome these problems—inspired by the success of data-driven approaches [72, 44, 27]—our method learns to infer critical states from historical visual trajectories of agents.

We propose a novel framework, namely the *Deep State Identifier*, to identify critical states in video-based environments. A principal challenge of working in such settings lies in acquiring ground-truth annotations of critical states; it is laborious to manually label in videos critical states cor-

† Equal Contribution.

✉ Corresponding Author.

Accepted to ICCV23.

responding to complex spatio-temporal patterns. The Deep State Identifier is designed to directly overcome this challenge by identifying the critical states based solely on visual inputs and rewards. Our proposed architecture comprises a return predictor and a critical state detector. The former predicts the return of an agent given a visual trajectory, while the latter learns a soft mask over the visual trajectory where the non-masked frames are sufficient for accurately predicting the return. Our training technique explicitly minimizes the number of critical states to avoid redundant information through a novel loss function. If the predictor can achieve the same performance using a small set of frames, we consider those frames critical. Using a soft mask, we obtain a rank that indicates the importance of states in a trajectory, allowing for the selection of critical states with high scores. During inference, critical states can be directly detected without relying on the existence of a return predictor. Our contributions can be summarized as follows:

- We propose a novel framework that effectively identifies critical states for reinforcement learning from videos, despite the lack of explicit action information.
- We propose new loss functions that effectively enforce compact sets of identified critical states.
- We demonstrate the utility of the learned critical states for policy improvement and comparing policies.

## 2. Related Work

In the past decade, researchers have explored the potential of combining computer vision (CV) and RL to develop more intelligent agents. A pioneering study by Koutnik et al. [32] used recurrent neural networks to tackle vision-based RL problems through an evolutionary strategy [33]. Since then, this topic has gained popularity. Mnih et al. [40, 41] trained a deep neural network using raw pixel data from Atari games to learn the Q-function for RL agents. Recently, Visual MPC [15] proposed a method using deep convolutional neural networks to predict the future states of a robot’s environment based on its current visual input. RIG [43] trains agents to achieve imagined goals in a visual environment using a combination of RL and an auxiliary visual network. Ha and Schmidhuber [22] propose a version of the world model, which employs a Variational Autoencoder (VAE) [31] to construct representations of the visual environment and help train a model using imagined future states. Robotprediction [14] designs a method for unsupervised learning of physical interactions through video prediction, achieved by an adversarial model that assists RL agents in learning to interact with the environment. More recently, researchers have explored novel CV advances, such as self-attention and self-supervised learning, applied to RL algorithms [28, 7, 73, 18, 10], leading to satisfactory

improvements. While visual input is integral to RL agents and can benefit RL in numerous ways, our paper proposes a method to assist agents in identifying the most crucial visual information for decision-making rather than solely focusing on improving visual representation.

Our method offers a novel perspective on explainable RL by identifying a small set of crucial states. Explaining the decision-making process in RL is more challenging than in CV, due to its reliance on sequential interactions and temporal dependencies. Various methods have been employed to address this challenge. Recent attention-based approaches [28, 7, 42] focus on modeling large-scale episodes offline [28, 7] to localize crucial decision-making points [42]. However, the attention structure typically operates on feature space, where the spatial correspondence is not aligned with the input space [5, 21]. Therefore, it is challenging to directly threshold attention values to identify critical temporal points. Post-training explanation is an efficient method that directly derives the explanation from an agent’s policy or value network [38, 20, 19, 16], thereby reducing memory and computation costs. Other popular explainable DRL methods include self-interpretable methods, such as Relational-Control Agent [74] and Alex [42], and model approximation methods, such as VIPER [4] and PIRL [67]. These methods are widely used in the field of DRL [38, 20, 19, 16, 74, 42, 4, 67]. For example, Alex [42] proposes using the output of the attention mechanism to enable direct observation of the information used by the agent to choose its action, making this model easier to interpret than traditional models. Tang et al. [64] use a small fraction of the available visual input and demonstrate that their policies are directly interpretable in pixel space. The PIRL method [67] produces interpretable and verifiable policies using a high-level, domain-specific language. Recent work uses policy fingerprinting [24] to build a single value function to evaluate multiple DRL policies [13, 12, 11]. The authors use only the policy parameters and the return to identify critical abstract states for predicting the return. However, policy parameters are often unavailable in practical applications, and storing them for multiple policies can require significant memory resources. We circumvent this issue by using visual states observed from the environment rather than relying on policy parameters.

Apart from the methods mentioned above, reward decomposition is also popular. Such methods [56, 29] re-engineer the agent’s reward function to make the rewards earned at each time step more meaningful and understandable. Compared to these methods, our approach evaluates the specific states. It provides a context-based framework for long-horizon trajectories in a challenging, yet practical domain, specifically learning without actions. Our method is also related to the concept of Hierarchical RL [71, 63], which aims to identify high-level subgoals [53, 47] that a

low-level policy should achieve. Using a few crucial states to explain an RL agent is closely connected to the concept of history compression [46, 48], where a neural network is trained to learn compact representations that are useful for modeling longer data sequences.

### 3. Method

#### 3.1. Problem Formulation

In Reinforcement Learning (RL) [62], an agent interacts sequentially with an environment. At each time step  $t$ , the agent observes a state  $s^{(t)}$ —in our case, the frame of a video, chooses an action  $a^{(t)}$ , obtains a scalar immediate reward  $r^{(t)} = R(s^{(t)}, a^{(t)})$ , where  $R$  is the reward function, and transitions to a new state  $s^{(t+1)}$  with probability  $P(s^{(t+1)}|s^{(t)}, a^{(t)})$ .

The behavior of an agent is expressed by its policy  $\pi(a|s)$ , which defines a probability distribution over actions given a state. The agent starts from an initial state and interacts with the environment until it reaches a specific state (a goal state or a failing state) or hits a time horizon  $T$ . Each of these interactions generates an *episode* and a *return*, i.e., the discounted cumulative reward  $\mathbf{y} = \sum_{t=0}^T \gamma^t r^{(t)}$ , where  $\gamma \in [0, 1)$  is a discount factor. Due to the general form of the return and the complex agent-environment interaction, it is generally difficult to identify which decision points—or states—are essential to achieve a specific return in an episode. In other words, it is difficult to explain the behavior of a policy.

Inspired by the success of data-driven approaches [72, 44, 27, 77], we design a learning-based method to identify a few crucial states in an episode that are critical to achieving the return  $\mathbf{y}$ . Unlike previous approaches [1, 21], we focus on identifying critical states in a video without needing an explicit representation of the policy or actions executed. More formally, let  $\{s_i, \mathbf{y}_i\}_i$  be the collected *episode-return training data*, where  $s_i = \{s_i^{(t)}\}_t$  is the  $i$ -th state trajectory,  $s_i^{(t)}$  is a state at the time step  $t$ , and  $\mathbf{y}_i$  is the return achieved in the state trajectory  $s_i$ .

To identify critical states, we suggest a novel framework, called the Deep State Identifier, consisting of the following two steps. **First**, we propose a return predictor that estimates the return  $\mathbf{y}_i$  given a state trajectory  $s_i$ . **Second**, we use the return predictor to train a critical state detector to identify critical states. The detector receives the states as input and outputs a mask over the states. It is used to measure how important each state is to the return. Fig. 2 illustrates the architecture of our method.

#### 3.2. Return Predictor

Our return predictor  $\mathcal{G}(\cdot)$  aims to predict the return of a sequence of states. We build it using a neural network and train it in a supervised manner. There are two types

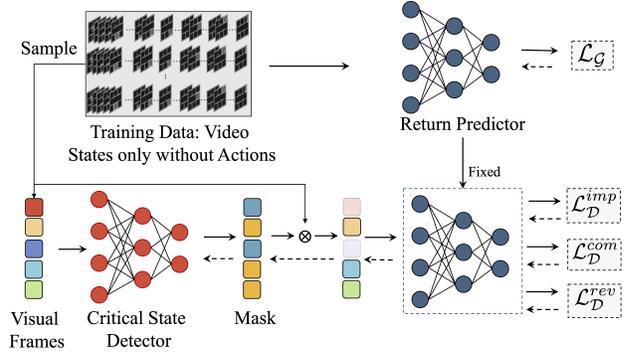


Figure 2. **Illustration of the proposed framework.** During training, our return predictor learns to predict the return of an episode from a state trajectory. Our critical state detector learns to exploit the return predictor to identify a compact set of states critical for return prediction. During testing, the critical state detector takes a state trajectory as input and automatically detects its critical states without using the return predictor.

of learning objectives depending on whether the return is discrete or continuous. For discrete return values (e.g., 1 indicates success, while 0 denotes failure), we train  $\mathcal{G}(\cdot)$  using cross-entropy loss:

$$\mathcal{L}_{\mathcal{G}}^c = \sum_i \mathcal{L}_{\mathcal{G}}^c(s_i, \mathbf{y}_i) = - \sum_i \mathbf{y}_i \log \mathcal{G}(s_i), \quad (1)$$

where  $\mathbf{y}_i$  is the category-level annotation of  $s_i$ . If the return is continuous, we employ a regression loss  $\mathcal{L}_{\mathcal{G}}^r$  to train  $\mathcal{G}(\cdot)$ ,

$$\mathcal{L}_{\mathcal{G}}^r = \sum_i \mathcal{L}_{\mathcal{G}}^r(s_i, \mathbf{y}_i) = \sum_i \|\mathcal{G}(s_i) - \mathbf{y}_i\|_2, \quad (2)$$

where  $\mathbf{y}_i \in \mathbb{R}$  is the scalar return of state trajectory  $s_i$ .

#### 3.3. Critical State Detector

In a general environment, manually labeling critical states is expensive and impractical. The unavailability of ground-truth critical states prevents our method from being fully-supervised. We hereby propose a novel way of leveraging the return predictor for training a critical state detector. Note that the critical states are elements of the state trajectory and can be discontinuous along the temporal dimension. We cast the task of identifying critical states as deriving a soft mask on a state trajectory. In particular, given a state trajectory  $s_i = \{s_i^{(t)}\}$ , the critical state detector  $\mathcal{D}$  outputs a mask on  $s_i$ , i.e.,  $\mathbf{m}_i = \mathcal{D}(s_i)$ , where  $\mathbf{m}_i = \{m_i^{(t)}\}$ ,  $m_i^{(t)} \in [0, 1]$  can be interpreted as confidence that  $s_i^{(t)}$  is a critical state. Intuitively, a high value of  $m_i^{(t)}$  indicates a higher probability that the corresponding state  $s_i^{(t)}$  is critical. To enforce  $\mathcal{D}$  to identify critical states,

we design three loss functions, namely, importance preservation loss, compactness loss, and reverse loss, for training  $\mathcal{D}$ :

$$\mathcal{L}_{\mathcal{D}} = \lambda_s \mathcal{L}_{\mathcal{D}}^{imp} + \lambda_r \mathcal{L}_{\mathcal{D}}^{com} + \lambda_v \mathcal{L}_{\mathcal{D}}^{rev}, \quad (3)$$

where  $\lambda_s$ ,  $\lambda_r$  and  $\lambda_v$  are the weights for importance preservation loss, compactness loss, and reverse loss respectively.

**Importance preservation loss.** Given a state trajectory  $\mathbf{s}_i$ , the goal of the importance preservation loss is to ensure the states discovered by the critical state detector are important to predict the return  $\mathbf{y}_i$ . Hence, the loss enforces the masked state sequence discovered by  $\mathcal{D}$  to contain a similar predictive information of the original state trajectory  $\mathbf{s}_i$ . Given the training data  $\{(\mathbf{s}_i, \mathbf{y}_i)\}$ , the importance preservation loss is defined as follows:

$$\mathcal{L}_{\mathcal{D}}^{imp} = \sum_i \mathcal{L}_{\mathcal{G}}(\mathcal{G}(\mathbf{s}_i \circ \mathcal{D}(\mathbf{s}_i)), \mathbf{y}_i), \quad (4)$$

where  $\circ$  denotes the element-wise multiplication ( $\mathbf{s}_i \circ \mathcal{D}(\mathbf{s}_i))^{(t)} \triangleq m_i^{(t)} s_i^{(t)}$ ,  $\mathcal{G}(\mathbf{s}_i \circ \mathcal{D}(\mathbf{s}_i))$  predicts the return of the masked state sequence  $\mathbf{s}_i \circ \mathcal{D}(\mathbf{s}_i)$ ,  $\mathcal{L}_{\mathcal{G}}$  stands for  $\mathcal{L}_{\mathcal{G}}^c$  or  $\mathcal{L}_{\mathcal{G}}^r$ , as defined in the previous subsection. Note that the masked state sequence can be discontinuous, and the information is dropped by skipping some redundant states. As a result, we cannot obtain a ground-truth return for a masked state sequence by running an agent in its environment. Thanks to the generalization abilities of neural networks [75, 68, 52, 49], we expect that the return predictor trained on the original state trajectories can predict well the return for masked state trajectories when critical states are not masked.

**Compactness loss.** Solely using the importance preservation loss  $\mathcal{L}_{\mathcal{D}}^{imp}$  leads to a trivial solution where the mask identifies all states in  $\mathbf{s}_i$  as critical. Critical states should instead be as compact as possible to avoid involving redundant and irrelevant states. To address this issue, we further introduce the compactness loss  $\mathcal{L}_{\mathcal{D}}^{com}$ . The compactness loss forces the discovered critical state to be as few as possible. Specifically, we employ the L1-norm to encourage the mask, *i.e.*, the output of  $\mathcal{D}$ , to be sparse given each  $\mathbf{s}_i$ :

$$\mathcal{L}_{\mathcal{D}}^{com} = \sum_i \|\mathcal{D}(\mathbf{s}_i)\|_1. \quad (5)$$

It is difficult to balance the importance preservation loss and compactness loss. The detector may ignore some critical states for compactness. We propose a reverse loss for training  $\mathcal{D}$  to mitigate this problem.

**Reverse loss.** The third loss is designed for undetected states. We remove the critical states by inverting the mask

from the original state trajectory  $\mathbf{s}_i \circ (1 - \mathcal{D}(\mathbf{s}_i))$  and process this masked sequence where the remaining states are useless for return prediction. This loss ensures that all the remaining states are not useful for estimating the return. We define the reverse loss as:

$$\mathcal{L}_{\mathcal{D}}^{rev} = - \sum_i \mathcal{L}_{\mathcal{G}}(\mathcal{G}(\mathbf{s}_i \circ (1 - \mathcal{D}(\mathbf{s}_i))), \mathbf{y}_i). \quad (6)$$

### 3.4. Iterative Training

Here we introduce the training strategy of our framework. We train the return predictor on complete and continuous state trajectories. At the same time, we use it to predict the return of masked state sequences that are incomplete and discontinuous when training the critical state detector. We iteratively train the predictor and the detector, where the learning objective of the whole framework is given by:

$$\min_{\mathcal{G}} \min_{\mathcal{D}} \mathcal{L}_{\mathcal{D}} + \mathcal{L}_{\mathcal{G}}. \quad (7)$$

After training, our critical state detector automatically detects critical states without using the return predictor. Appendix A lists the pseudo-code of the proposed method.

## 4. Experiments

### 4.1. Benchmark and Protocol Navigation

We begin this section by releasing a benchmark to test our method and facilitate the research on explainability. As shown in Table 1, we collect five datasets on three different RL environments, *i.e.*, Grid World [9, 8], Atari-Pong [6], and Atari-Seaquest [6]. We select Grid World for qualitative analysis since it is very intuitive for human understanding. We study a challenging environment with partial observation. In the context of Grid World, we define a "state" as a combination of the current visual frame and historical information. Although this surrogate representation does not equate to the full, true state of the environment, it serves as an agent's internal understanding, developed from its sequence of past observations. To elaborate, when we say that our model identifies a "state" in this context, we imply that it recognizes a specific observation or frame, based on the agent's history of previous observations. For fully observable environments like Atari, the term "state" assumes its traditional definition, providing complete information about the system at any given time. We use Atari-Pong and Atari-Seaquest environments to compare our method with similar approaches based on critical state identification, using adversarial attacks, and evaluating policy improvement. Note that evaluating critical states using adversarial attacks was first proposed by work on Edge [21]. However, Edge does not consider cross-policy attacks where the policies for training and testing the detector are different. More details can be found in the supplementary material.

Table 1. **The specification of the five collected datasets.** The datasets cover discrete and continuous returns for a comprehensive study of the proposed method.  $y$  here is the cumulative reward.

	Length	Training	Test	Total
Grid World-S (Memory: 353 MB)				
Reaching Goal	31.97	1000	200	1200
Fail	25.72	1000	200	1200
Grid World-M (Memory: 412 MB)				
Policy-1	31.97	1000	200	1200
Policy-2	38.62	995	200	1195
Atari-Pong-[S/M](Memory: 174 GB /352 GB)				
Agent Win	200	13158/17412	1213/1702	14371/19114
Agent Lose	200	8342/4088	787/298	9129/4386
Total	-	21500	2000	23500
Atari-Seaquest-S (Memory:706 GB)				
$\mathbb{E}[y]=2968.6$	2652.5	8000	2000	10000

Table 2. **Summary of improvements due to our method**, where Gain refers to improvement over the baselines. Our method improves performance across various tasks. The baselines in the 2nd-6th rows are our method using *Imp. Loss* on Grid-World-S, EDGE [21] for Atari-Pong-S, an attack with 30 randomly selected frames on Atari-Pong-M, and DQN trained with 25M time steps on Atari-Seaquest-S, respectively.

Datasets	Navigation	Task	Gain
GridWorld-S	Sec. 4.2	Critical State Identify	<b>16.38%</b>
GridWorld-S	Sec. 4.2	Sequence Reasoning	<b>Qualitative</b>
GridWorld-M	Sec. 4.3	Policy Evaluation	<b>First Study</b>
Atari-Pong-S	Sec. 4.4	In-Policy Adv. Attack	<b>18.63%</b>
Atari-Pong-M	Sec. 4.4	Robust Analysis	<b>50.35%</b>
Atari-Seaquest-S	Sec. 4.5	Policy Improvement	<b>17.65%</b>

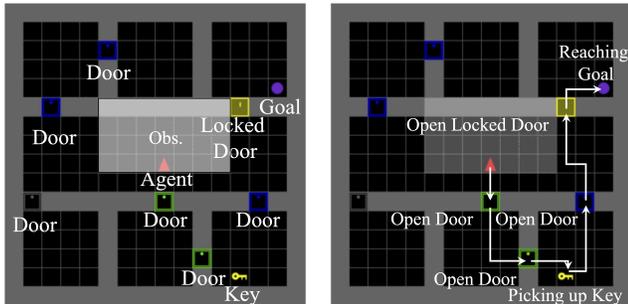


Figure 3. **Illustration of an instance of the GridWorld environment.** The environment consists of six rooms and one corridor. The agent starts from a random initial location in the corridor, and the final goal is to get the ball. Since the ball is locked in a room, the agent must pick up the key and open the yellow locked door. In a successful episode, the agent must open the unlocked doors (colored in green and blue), pick up the key, unlock the yellow door, and reach the purple ball. Note that the agent has only partial observation (colored white) of the environment at a time step.

## 4.2. Critical State Discovery

**Performance.** This section provides a qualitative analysis of the critical time point identified by our Deep State Identifier. We choose the ‘MiniGrid-KeyCorridorS6R3-v0’

Table 3. **Ablation study for the critical state detector.**

<i>Imp. Loss</i>	<i>Com. Loss</i>	<i>Rev. Loss</i>	F-1 Score (%) $\uparrow$
✓	×	×	68.98
✓	✓	×	unstable
×	✓	✓	74.42
×	×	✓	76.09
✓	✓	✓	<b>80.28</b>

task [54, 9] of the GridWorld environment, where the goal is to reach a target position in a locked room after picking up a key (see the yellow block in Fig. 3). This task is useful to visually evaluate our method since it is intuitive to identify what states are critical: top row in Fig. 4 shows that states immediately before actions such as ‘opening the door’ (S.1, S.2, S.3), ‘picking up the key’ and ‘opening the locked door’ are critical to successfully reaching the goal. Note that there is no ground truth on the critical state for a general, more complex environment.

We use a pre-defined DRL agent to collect trajectories. Since our method detects critical states by masking the trajectory, we evaluate how our critical state detector accurately assigns high scores to the states we intuitively labeled as critical. As shown in Fig. 4, our method assigns high values to human-annotated critical states and low values to remaining states, showing its effectiveness in discovering critical states.

**Ablation study.** We analyze the contribution of each component of the critical state detector loss in Tab. 3 and Fig. 5. If we remove the compactness loss and the reverse loss, our method wrongly assigns high confidence to all states in an episode, *i.e.*, all states are detected as critical ones. Similarly, if we remove the reverse loss, our method detects all states as non-critical. Finally, removing only the compactness loss, most states (including non-critical ones) are wrongly detected as critical. This ablation shows that each loss component is crucial to critical state identification.

**More Analysis.** In RL, states within an episode can be highly correlated. We show how our method can discover state dependencies essential to identifying critical states. It is challenging to capture the dependencies among states in the Gridworld since the agent can only partially observe the environment through a small local view.

Tab. 4 provides examples of states in the environment<sup>1</sup>. In Gridworld, the states that occur immediately before or after the action “opening door” are frequently observed in a trajectory. In these states, the agent can be either with or without the key. However, obtaining the key is crucial for achieving the goal of GridWorld (see Fig. 3). Without the key, the agent cannot successfully finish the task. Therefore, the states immediately before or after the action “opening door” without the key are not as critical as the

<sup>1</sup>We use a text description of states due to space constraints. We provide visual states in the supplemental material.

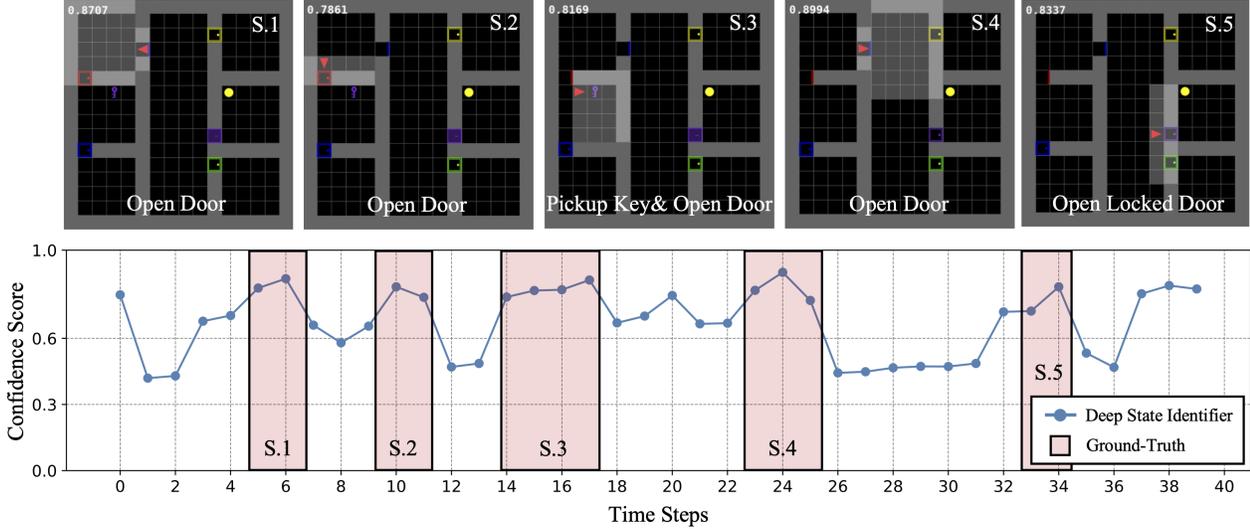


Figure 4. **The performance of our method in identifying critical states.** The top row shows human-annotated critical states (i.e., ground truth) in an episode. The bottom row shows for each time step in the environment how confident the detector is that the current state is critical. Our method assigns high scores to human-annotated critical states, demonstrating its identification abilities.

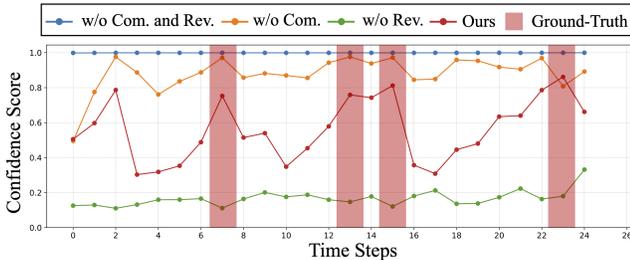


Figure 5. **Ablation study of the detector’s loss function.** For each time step and loss component, the line indicates how confident the detector is that the current input is critical. Red blocks mark the human annotation.

states immediately before or after the action “opening the door” with the key to predict the return. Tab. 4 shows how our method captures such dependencies between “opening door” and “picking up the key.” Our method successfully assigns much higher confidence to the critical states immediately before or after the action “opening door” with the key and lower confidence to the states immediately before or after the action “opening door” without the key.

### 4.3. Policy Comparison by Critical States

In general, researchers use cumulative rewards to validate policy performance. However, these metrics cannot elucidate the diverse behavioral patterns exhibited by different policies. To better distinguish and explain the behavioral differences among various policies, a return predictor is trained to recognize the distinct trajectories of each policy. Our detector then is trained to identify critical states for highlighting the contrasts between policies rather than merely focusing on returns, thus facilitating a more compre-

Table 4. **State detector’s confidence score over different states.** Our method has different confidence scores for the states immediately before and after (i.b.a.) opening a door with or without the key, which indicates that it can capture temporal dependencies among states. Normal states refer to states where the agent has a distance greater than two from positions where it can take a relevant action (pick up the key or open a door). We report the mean and standard deviation of the confidence over four random seeds.

State Description	Confidence Score
Normal States (Full)	$53.66 \pm 0.12$
Normal States Before Picking up the Key	$49.59 \pm 0.13$
State i.b.a. Opening Door (without the Key)	<b><math>67.13 \pm 0.12</math></b>
State i.b.a. Trying Locked Door (without the Key)	<b><math>50.81 \pm 0.08</math></b>
State i.b.a. Picking up the Key	<b><math>78.35 \pm 0.04</math></b>
Normal States After Picking Up the Key	$56.58 \pm 0.10$
State i.b.a. Opening Door (with the Key)	<b><math>80.65 \pm 0.06</math></b>
State i.b.a. Opening Locked Door	<b><math>87.55 \pm 0.01</math></b>

hensive comparison of their behaviors. Consequently, we can leverage the ability of the critical state detector to pinpoint the key states that discriminate between the two policies and visually represent the dissimilarities between them. As shown in Fig. 6, both policy-A and policy-B can achieve the final goal, but in policy-B, the agent always enters an invalid room after picking up the key, leading to more steps in the environment before achieving the goal. Both policies achieve a high return. However, our approach identifies the most discriminating states. Our method precisely assigns the highest confidence to the states inside the invalid room. The visualization shows that our method can explain the difference between the two policies. More details are provided in Appendix A.

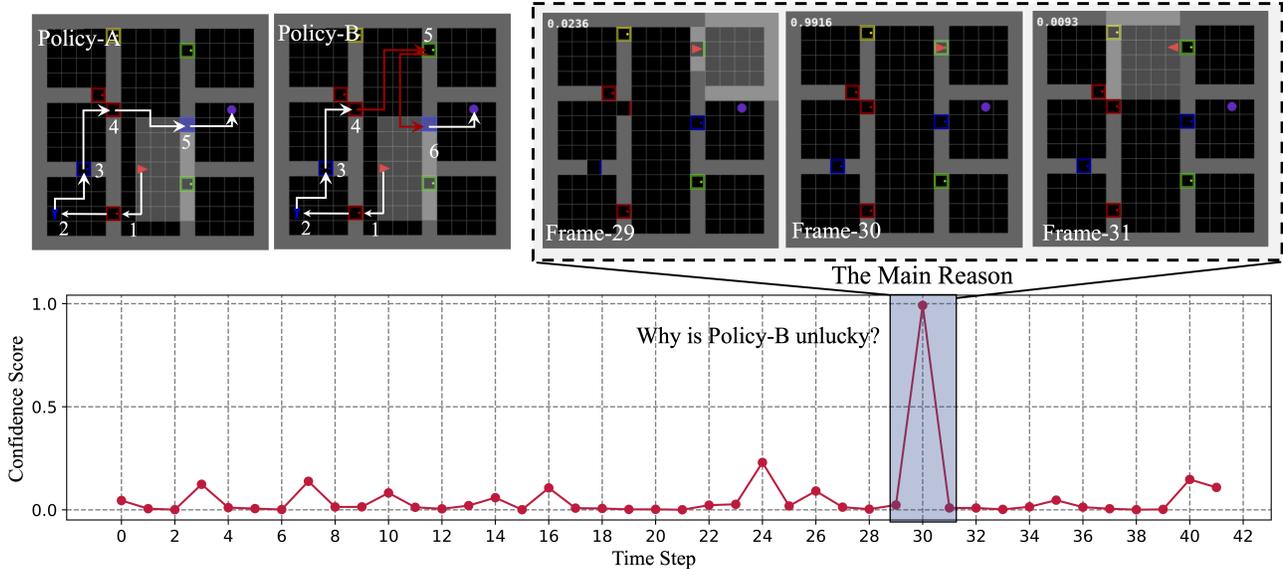


Figure 6. **Visualization of the Deep State Identifier for policy comparison.** We pre-collect policy-A and policy-B. While policy-A is optimal, policy-B first causes the agent to enter the incorrect room after picking up the key and then reach the goal. We train our method to discriminate between policy-A and policy-B, given sequences of trajectories generated by them. The critical state detector assigns high confidence to states where policy-B is suboptimal.

Table 5. **Win rate changes of the agent before/after attacks by following the protocol of EDGE [21].** We use the detected top 30 states as input to attack the policy. We report means and standard deviations over three random seeds. The reported results of all the baselines are from previous work [21].  $s, a, y, \pi$  denote the state, action, return, and policy parameters, respectively.

Method	Input	Win Rate Changes ↓
Rudder [1]	( $s, a, y$ )	$-19.93 \pm 4.43$
Saliency [57, 59, 60]	( $s, a, y$ )	$-30.33 \pm 0.47$
Attention RNN [2]	( $s, a, y, \pi$ )	$-25.27 \pm 1.79$
Rationale Net [36]	( $s, a, y, \pi$ )	$-29.20 \pm 4.24$
Edge [21]	( $s, a, y, \pi$ )	$-65.47 \pm 2.90$
Ours with single policy	( $s, y$ )	<b><math>-77.67 \pm 0.56</math></b>
Ours with multiple policies	( $s, y$ )	<b><math>-85.90 \pm 1.47</math></b>

#### 4.4. Efficient Attack using Critical States

In the previous sections, we showed that our method identifies the critical states with the highest impact on return prediction. However, for complex environments, it is difficult to evaluate the performance of this class of methods because the ground-truth critical states are not available. Following previous approaches [21], we use adversarial attacks to validate whether the identified states are critical. Intuitively, if a state is critical, introducing noise in the action that a policy would take in such a state will significantly deteriorate performance (the return will be lower). Here we follow the same protocol of previous approaches [21], and we compare the policy’s performance drop to the baseline methods when the 30 most critical states are attacked (i.e., whenever the agent reaches those states, its action is perturbed).

Table 5 shows that our method outperforms the other techniques in the Atari-Pong environment, exhibiting the most significant changes in win rates, highlighting its efficacy in localizing critical states. In particular, we achieve an 18.63% improvement over the previous SOTA method Edge[21], suggesting that the states identified by our Deep State Identifier are more crucial to achieve a high return. Note that the previous methods, such as Edge [21], are based on sequences of states and action pairs. Our method instead achieves higher performance by only observing a state sequence. In the real-world scenario, imaging systems can easily capture sequences of visual states, while actions are more difficult to collect, requiring special sensors or manual annotations. In other words, our method can work with pure visual information to achieve higher performance, resulting in flexibility toward various potential applications. Moreover, when different policies collect the training dataset, the proposed method can benefit from data diversity, inducing more satisfactory results (i.e., an 85.90 drop in winning performance).

We then analyze the attack performance across different policies to test the robustness against policy shifts. In Table 6, we set the baseline that attacks 30 states chosen randomly and attacks a policy that was never used to train our method. To ensure policy diversity for testing, we derive the policies with various random seeds, training steps, and network architectures. Compared with the baseline, our method cannot improve performance using a single policy, which indicates that a cross-policy protocol is challenging for adversarial attacks. However, when we increase the training data

Table 6. **Win rate changes of the agent before/after attacks for different policies.** We assess whether our method, trained on trajectories generated by one or multiple policies, can accurately identify critical time points within a trajectory generated by another unseen policy. We consider three kinds of unseen policies, including different random seeds (seeds), different training steps (steps), and different network architectures (Arch.), to test the performance of our method against cross-policy challenges. We report mean and standard error over three random seeds. We attack the policy perturbing its action in the top 30 states detected.

	Baseline	Ours (Single)	Ours (Multi.)
In-Policy (baseline)	54.88 ± 1.80	-77.67 ± 0.56	<b>-85.90 ± 1.47</b>
Cross-Policy (Seeds)	-63.32 ± 0.93	-30.67 ± 0.58	<b>-85.45 ± 0.86</b>
Cross-Policy (Steps)	-50.23 ± 1.21	-30.57 ± 1.01	<b>-83.72 ± 0.91</b>
Cross-Policy (Arch.)	-49.85 ± 3.50	-39.55 ± 2.38	<b>-76.50 ± 3.11</b>

Table 7. **Performance of DQN with different adaptive step strategies on Atari-Seaquest.** We base the implementation on the Tianshou Platform [70]. Our method effectively improves the performance of DQN. n-step stands for the lookahead steps.

Methods	Return ↑ ± Std.
PPO (time steps=5M) [55]	887.00 ± 4.36
SAC (time steps=5M) [23]	1395.50 ± 339.34
Rainbow (step=3,time steps=5M) [25]	2168.50 ± 332.89
DQN(time steps=10M) [40]	3094.75 ± 1022.54
DQN (n-step=random(1,5),time steps=5M) [61]	3250.25 ± 638.13
Baseline: DQN (n-step=5,time steps=5M) [61]	1987.00 ± 115.71
DQN (n-step=12,time steps=5M) [61]	1472.50 ± 407.40
DQN (n-step=grid search,time steps=5M) [61]	3936.50 ± 459.19
SAC (time steps=25M)[23]	1444.00 ± 136.86
Rainbow (time steps=25M)[25]	2151.25 ± 329.29
DQN (time steps=25M)[40]	3525.00 ± 63.87
HL based on Frequency (time steps=5M)[39, 58]	2477.00 ± 223.65
DQN + Ours (n-step≤5,time steps=5M)	<b>4147.25 ± 378.16</b>

diversity by adding policies, we achieve a higher generalization, and the model’s drop in performance improves from 49.85 to 76.50. A potential explanation is that each policy induces a specific distribution over the state space in the environment. Using different policies to collect data allows us to generalize to unseen policies and achieve more invariant representations of the policy behavior. Indeed, when the dataset can cover the distribution of states in the environment, our method generalizes to arbitrary unseen policies. We thereby achieve an environment-specific policy-agnostic solution for interoperability.

#### 4.5. Policy Improvement

We show how our method can improve DRL policies. The experimental results in the previous sections demonstrate that our Deep State Identifier can efficiently identify critical states. Here we show how one can use these states to perform rapid credit assignment for policy improvement. In particular, we combine our method with the widely-used DQN [40] for multi-step credit assignment. The objective function of traditional Multi-step DQN[25, 61] is:

$$\sum_{(s^{(j)}, a^{(j)}) \in \text{Rep.}} \left[ Q(s^{(j)}, a^{(j)}) - \left( \sum_{t=j}^{j+n-1} \gamma^{t-j} r^{(t)} + \gamma^n \max_{a^{(j+n)}} Q^T(s^{(j+n)}, a^{(j+n)}) \right) \right]^2, \quad (8)$$

where  $Q$  is the action-value function, i.e., a network predicting the expected return of the policy from a particular state-action pair, Rep. is the replay buffer,  $Q^T$  is a copy of  $Q$ , which is periodically synchronized with  $Q$  to facilitate learning,  $\gamma$  is the discount factor, and  $a$  denotes an action.

A recent study [69] highlights the importance of varying the lookahead step  $n$  in Multi-step DQN. Here we combine our method with Multi-step DQN by first identifying critical states and then dynamically setting lookahead steps to learn DQN. In other words, we set  $n$  as the number of time steps from the state to the most critical state detected within a specific range. Here, we set the maximum lookahead step to 5.

Table 7 presents preliminary results which illustrate that Multi-step DQN combined with our method improves the return of DQN from 1987.00 to 4147.25. Since our method effectively discovers states important for return prediction, our Deep State Identifier provides DQN with faster credit assignment, improving its performance. Moreover, our method performs slightly better than finely tuning the lookahead step  $n$  using grid search. Table 7 also includes improved versions of DQN [39, 58] for comparison. Our method outperforms all of them.

## 5. Conclusion

Our novel method identifies critical states from episodes encoded as videos. Its return predictor and critical state detector collaborate to achieve this. When the critical state detector is trained, it outputs a soft mask over the sequence of states. This mask can be interpreted as the detector’s belief in the importance of each state. Experimental results confirm that the generated belief distribution closely approximates the importance of each state. Our approach outperforms comparable methods for identifying critical states in the analyzed environments. It can also explain the behavioral differences between policies and improve policy performance through rapid credit assignment. Future work will focus on applying this method to hierarchical RL and exploring its potential in more complex domains.

## Acknowledgements

We thank Dylan R. Ashley for his valuable comments and help to polish the paper. This work was supported by the European Research Council (ERC, Advanced Grant

Number 742870) and the SDAIA-KAUST Center of Excellence in Data Science and Artificial Intelligence (SDAIA-KAUST AI).

## References

- [1] Jose A Arjona-Medina, Michael Gillhofer, Michael Widrich, Thomas Unterthiner, Johannes Brandstetter, and Sepp Hochreiter. Rudder: Return decomposition for delayed rewards. *NIPS*, 32, 2019.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *ICLR*, 2015.
- [3] Trapit Bansal, Jakub Pachocki, Szymon Sidor, Ilya Sutskever, and Igor Mordatch. Emergent complexity via multi-agent competition. *arXiv preprint arXiv:1710.03748*, 2017.
- [4] Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy extraction. *NIPS*, 31, 2018.
- [5] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba. Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 117(48):30071–30078, 2020.
- [6] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [7] 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. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- [8] Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of grounded language learning. *arXiv preprint arXiv:1810.08272*, 2018.
- [9] Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld environment for gymnasium. *GitHub*, 2018.
- [10] Benjamin Eysenbach, Tianjun Zhang, Ruslan Salakhutdinov, and Sergey Levine. Contrastive learning as goal-conditioned reinforcement learning. *arXiv preprint arXiv:2206.07568*, 2022.
- [11] Francesco Faccio, Vincent Herrmann, Aditya Ramesh, Louis Kirsch, and Jürgen Schmidhuber. Goal-conditioned generators of deep policies. *arXiv preprint arXiv:2207.01570*, 2022.
- [12] Francesco Faccio, Louis Kirsch, and Jürgen Schmidhuber. Parameter-based value functions. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.
- [13] Francesco Faccio, Aditya Ramesh, Vincent Herrmann, Jean Harb, and Jürgen Schmidhuber. General policy evaluation and improvement by learning to identify few but crucial states. *arXiv preprint arXiv:2207.01566*, 2022.
- [14] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. *Advances in neural information processing systems*, 29, 2016.
- [15] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2786–2793. IEEE, 2017.
- [16] Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *ICCV*, pages 3429–3437, 2017.
- [17] Kuniyuki Fukushima. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4):193–202, 1980.
- [18] Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurams, Sergey Levine, and Pieter Abbeel. Multimodal masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022.
- [19] Wenbo Guo, Sui Huang, Yunzhe Tao, Xinyu Xing, and Lin Lin. Explaining deep learning models—a bayesian non-parametric approach. *NIPS*, 31, 2018.
- [20] Wenbo Guo, Dongliang Mu, Jun Xu, Purui Su, Gang Wang, and Xinyu Xing. Lemna: Explaining deep learning based security applications. In *proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, pages 364–379, 2018.
- [21] Wenbo Guo, Xian Wu, Usman Khan, and Xinyu Xing. Edge: Explaining deep reinforcement learning policies. *NIPS*, 34:12222–12236, 2021.
- [22] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [23] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *ICML*, pages 1861–1870. PMLR, 2018.
- [24] Jean Harb, Tom Schaul, Doina Precup, and Pierre-Luc Bacon. Policy evaluation networks. *arXiv preprint arXiv:2002.11833*, 2020.
- [25] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI*, volume 32, 2018.
- [26] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, Nov. 1997.
- [27] Qibin Hou, PengTao Jiang, Yunchao Wei, and Ming-Ming Cheng. Self-erasing network for integral object attention. *Advances in Neural Information Processing Systems*, 31, 2018.
- [28] Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence modeling problem. In *Advances in Neural Information Processing Systems*, 2021.
- [29] Zoe Juozapaitis, Anurag Koul, Alan Fern, Martin Erwig, and Finale Doshi-Velez. Explainable reinforcement learning via reward decomposition. In *IJCAI/ECAI Workshop*, 2019.

- [30] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [31] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [32] Jan Koutník, Giuseppe Cuccu, Jürgen Schmidhuber, and Faustino Gomez. Evolving large-scale neural networks for vision-based reinforcement learning. In *Proceedings of the 15th annual conference on Genetic and evolutionary computation*, pages 1061–1068, 2013.
- [33] Jan Koutník, Faustino Gomez, and Jürgen Schmidhuber. Evolving neural networks in compressed weight space. In *Proceedings of the 12th annual conference on Genetic and evolutionary computation*, pages 619–626, 2010.
- [34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- [35] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62, 2022.
- [36] Tao Lei, Regina Barzilay, and Tommi Jaakkola. Rationalizing neural predictions. *EMNLP-IJCNLP*, 2017.
- [37] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.
- [38] Yang Young Lu, Wenbo Guo, Xinyu Xing, and William Stafford Noble. Dance: Enhancing saliency maps using decoys. In *ICML*, pages 7124–7133. PMLR, 2021.
- [39] Amy McGovern and Andrew G Barto. Automatic discovery of subgoals in reinforcement learning using diverse density. *Computer Science Department Faculty Publication Series*. 8., 2001.
- [40] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [41] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [42] Alexander Mott, Daniel Zoran, Mike Chrzanowski, Daan Wierstra, and Danilo Jimenez Rezende. Towards interpretable reinforcement learning using attention augmented agents. *NIPS*, 32, 2019.
- [43] Ashvin V Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforcement learning with imagined goals. *Advances in neural information processing systems*, 31, 2018.
- [44] Phuc Nguyen, Ting Liu, Gautam Prasad, and Bohyung Han. Weakly supervised action localization by sparse temporal pooling network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6752–6761, 2018.
- [45] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [46] J. Schmidhuber. Adaptive history compression for learning to divide and conquer. In *Proc. International Joint Conference on Neural Networks, Singapore*, volume 2, pages 1130–1135. IEEE, 1991.
- [47] J. Schmidhuber. Neural sequence chunkers. Technical Report FKI-148-91, Institut für Informatik, Technische Universität München, April 1991.
- [48] Jürgen Schmidhuber. Learning complex, extended sequences using the principle of history compression. *Neural Computation*, 4(2):234–242, 1992.
- [49] Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural networks*, 61:85–117, 2015.
- [50] Jürgen Schmidhuber. On learning to think: Algorithmic information theory for novel combinations of reinforcement learning controllers and recurrent neural world models. *arXiv preprint arXiv:1511.09249*, 2015.
- [51] Jürgen Schmidhuber. One big net for everything. *arXiv preprint arXiv:1802.08864*, 2018.
- [52] Juergen Schmidhuber. Annotated history of modern ai and deep learning. *arXiv preprint arXiv:2212.11279*, 2022.
- [53] Jürgen Schmidhuber and Reiner Wahnsiedler. Planning simple trajectories using neural subgoal generators. In *From Animals to Animats 2: Proceedings of the Second International Conference on Simulation of Adaptive Behavior*, volume 2, page 196. MIT Press, 1993.
- [54] Juergen Schmidhuber, Jieyu Zhao, and MA Wiering. Simple principles of metalearning. *Technical report IDSIA*, 69:1–23, 1996.
- [55] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [56] Tianmin Shu, Caiming Xiong, and Richard Socher. Hierarchical and interpretable skill acquisition in multi-task reinforcement learning. *arXiv preprint arXiv:1712.07294*, 2017.
- [57] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [58] Özgür Şimşek and Andrew Barto. Skill characterization based on betweenness. *Advances in neural information processing systems*, 21, 2008.
- [59] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825*, 2017.
- [60] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *ICML*, pages 3319–3328. PMLR, 2017.
- [61] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*, 3:9–44, 1988.
- [62] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, USA, 2018.

- [63] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211, 1999.
- [64] Yujin Tang, Duong Nguyen, and David Ha. Neuroevolution of self-interpretable agents. In *Proceedings of the 2020 Genetic and Evolutionary Computation Conference*, pages 414–424, 2020.
- [65] 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*, pages 5026–5033. IEEE, 2012.
- [66] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015.
- [67] Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaudhuri. Programmatically interpretable reinforcement learning. In *ICML*, pages 5045–5054. PMLR, 2018.
- [68] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*, 2022.
- [69] Yuhui Wang, Haozhe Liu, Miroslav Strupl, Francesco Faccio, Qingyuan Wu, Xiaoyang Tan, and Jürgen Schmidhuber. Highway reinforcement learning. *Open Review*, 2023.
- [70] Jiayi Weng, Huayu Chen, Dong Yan, Kaichao You, Alexis Duburcq, Minghao Zhang, Yi Su, Hang Su, and Jun Zhu. Tianshou: A highly modularized deep reinforcement learning library. *JMLR*, 23(267):1–6, 2022.
- [71] Marco Wiering and Jürgen Schmidhuber. Hq-learning. *Adaptive Behavior*, 6:219–246, 09 1997.
- [72] Jinheng Xie, Xianxu Hou, Kai Ye, and Linlin Shen. Clims: cross language image matching for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4483–4492, 2022.
- [73] Tao Yu, Zhizheng Zhang, Cuiling Lan, Zhibo Chen, and Yan Lu. Mask-based latent reconstruction for reinforcement learning. *NIPS*, 2022.
- [74] 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 *ICLR*, 2018.
- [75] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. Domain generalization: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [76] Deyao Zhu, Yuhui Wang, Jürgen Schmidhuber, and Mohamed Elhoseiny. Guiding online reinforcement learning with action-free offline pretraining. *arXiv preprint arXiv:2301.12876*, 2023.
- [77] Mingchen Zhuge, Deng-Ping Fan, Nian Liu, Dingwen Zhang, Dong Xu, and Ling Shao. Salient object detection

via integrity learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

This appendix provides the implementation details of our Deep State Identifier. In Section A, we provide the pseudo-code for the Deep State Identifier, its network architecture, and the hyperparameters used during training. Then, Section B discusses the datasets we collected and our experimental protocol. Finally, Section C provides additional experimental results related to the ablation study and the comparison with EDGE [21] on MuJoCo.

## A. Implementation Details

This section details our implementation of the proposed method. We implement our method and conduct our experiments using PyTorch [45]. All experiments were conducted on a cluster node equipped with 4 Nvidia Tesla A100 80GB GPUs.

The proposed method—while quite effective—is conceptually simple. The training pipeline can be written in 25 lines of pseudo-code:

```

1 import torch as T
2 def cs_detector_train(input_states, labels):
3     mask = cs_detector(input_states)
4     loss_reg = lambda_r * T.linalg.norm(mask, ord=1)
5     masked_states = mask * input_states
6     output = return_predictor(masked_states)
7     loss_sub = lambda_s * criterion(output, labels)
8     reverse_mask = torch.ones_like(mask) - mask
9     reverse_states = reverse_mask * input_states
10    output_r = return_predictor(reverse_states)
11    confused_label = torch.ones_like(output_r)
12    *0.5 #binary classification case
13    loss_vic = lambda_v * criterion(output_r,
14    confused_label)
15    loss_total = loss_reg + loss_sub + loss_vic
16    loss_total.backward()
17    optimizer_cs.step()
18 def return_predictor_train(input_states, labels):
19    output = return_predictor(input_states)
20    loss_d = criterion(output, labels)
21    loss_d.backward()
22    optimizer_return.step()
23 def main_train(input_states, labels):
24    optimizer_cs.zero_grad()
25    cs_detector_train(input_states, labels)
26    optimizer_return.zero_grad()
27    return_predictor_train(input_states, labels)

```

We use two potential network architectures in our work, 3DCNN [66], and CNN-LSTM [17, 26, 34], to implement our Deep State Identifier. Tables 8 and 9 show the specification of the corresponding architectures. We use 3DCNN architecture in Table 10 and employ LSTM structure in the other empirical studies.

To train the critical state detector and return predictor, we use the Adam optimizer [30] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . The learning rate is set as  $1 \times 10^{-4}$  and the weight decay is  $1 \times 10^{-4}$ . The input length of 3DCNN is 12 frames and is a partial observation ( $7 \times 7$  pixels) of the environment [9, 8]. The remaining hyper-parameters  $\lambda_s$ ,  $\lambda_r$ , and  $\lambda_v$  are set to 1,  $5 \times 10^{-3}$  and 2 respectively.

Table 8. **The specification of the 3DCNN-based Neural Network adopted in this paper.** In-Norm refers to the Instance Normalization, 3D Conv. is the 3D convolutional Layer, and F.C. refers to the fully connected layer. In the last layer, the [return predictor/critical state detector] has a different architecture specified in the last column.

3DCNN	Channel	Filter	Stride	In-Norm	Activation
3D Conv.	12 $\rightarrow$ 32	(1,3,3)	(1,2,2)	False	Relu
3D Conv.	32 $\rightarrow$ 64	(1,3,3)	(1,1,1)	True	Relu
3D Conv.	64 $\rightarrow$ 128	(1,3,3)	(1,2,2)	False	Relu
3D Conv.	128 $\rightarrow$ 128	(1,3,3)	(1,1,1)	True	Relu
3D Conv.	128 $\rightarrow$ 256	(3,2,2)	(1,1,1)	False	Relu
Avg Pooling	-	-	-	-	-
F.C.	256 $\rightarrow$ 512	-	-	-	-
F.C.	512 $\rightarrow$ [2/12]	-	-	-	[-/sigmoid]

Table 9. **The specification of the CNN-LSTM Neural Network in this paper.** In the last layer, the critical state detector outputs a vector with the same length as the input (i.e., 256  $\rightarrow$  1). The return predictor estimates a scalar for the whole episode (i.e., 256  $\times$  length  $\rightarrow$  2)

CNN-LSTM	Channel	Filter	Stride	In-Norm	Activation
2D Conv.	3 $\rightarrow$ 32	3	2	False	Relu
2D Conv.	32 $\rightarrow$ 64	3	1	True	Relu
2D Conv.	64 $\rightarrow$ 128	3	2	False	Relu
2D Conv.	128 $\rightarrow$ 128	3	1	True	Relu
2D Conv.	128 $\rightarrow$ 256	2	1	False	Relu
Avg Pooling	-	-	-	-	-
	Input	Hidden	Bi-Direct.		Activation
LSTM	256	128	True		-
F.C.	[length $\times$ 256] $\rightarrow$ [2/length]				[-/sigmoid]

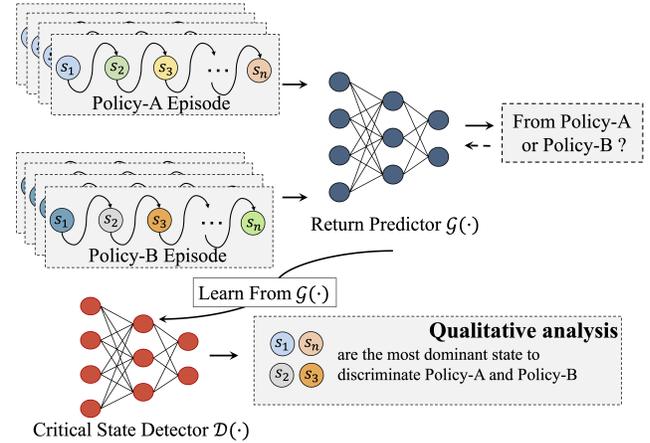


Figure 7. **Illustration of the Deep State Identifier for policy comparison.** We modify the return predictor as a binary classifier. Its training data comprises pairs  $\{s_i, c_i\}$ , where  $s_i$  represents a trajectory and  $c_i \in \mathbb{R}$  is a class label indicating whether it belongs to policy-A or policy-B. By exploiting the return predictor, the critical state detector can directly localize the states that primarily explain the difference between policy-A and policy-B.

Fig. 7 shows how we can adapt the return predictor to find the critical frame that explains the difference in behavior between the two policies. We can train the return predic-

tor to identify which of the two policies generates a specific trajectory.

## B. Experimental details

**Critical States Discovery.** We use a GridWorld environment (MiniGrid-KeyCorridorS6R3-v0) to collect a dataset (Grid-World-S) to test the accuracy of the critical state detector. Data is collected by acting in the environment using an optimal policy based on a depth-first search algorithm (DFS). Additional data is collected from a random-exploring policy. Since, in this environment, one can find critical states by visual inspection (they correspond to the states immediately before or after the action of opening doors or picking up keys), we can directly test the accuracy of the proposed method. We use the F1 score as a metric.

**Policy Comparison by Critical States.** Here, we collect a dataset, *Grid-World-M*, for our experiments on policy comparison. The labels in *Grid-World-M* are the policies that collected the corresponding episode. We use two policies to collect data: Policy-A is the optimal policy used to collect *Grid-World-S*, while Policy-B is an exploratory policy.

**Efficient Attack using Critical States.** Here we use adversarial attacks on Atari-Pong to validate whether the detected states are critical. Following the same protocol as Edge [21], we use a trained policy downloaded from <https://github.com/greydanus/baby-a3c> to collect the training data. We call the corresponding dataset *Atari-Pong-S*. In particular, we collect 21500 episodes for training and 2000 for testing, and we fix the length of each episode as 200. We augment the input by randomly increasing or decreasing the length within 50 frames, and the padding value is set as 0. To validate the generalization of the proposed method for unseen policies, we then collect another dataset, denoted *Atari-Pong-M*. We train policies with different seeds using the same implementation as Edge [21] from <https://github.com/greydanus/baby-a3c>. In particular, we use ten different policies to collect training data. In cross-policy (seeds), we use the trained policy on different random seeds to test the performance. In cross-policy (steps), we use the policy trained with 80M and 40M steps for training and testing our method, respectively. In cross-policy (Arch.), we change the architecture to make the setting more challenging. In particular, we train our method using a policy with 32 channels but test it by attacking a policy trained using 64 channels. The result in each case is collected by attacking the agent for 1500 episodes using three random seeds.

**Policy Improvement.** We test the potential of our method to improve policy performance in the Atari-Seaquest environment. We first train the policies based on DQN following the implementation of Tianshou [70]. Then we use the trained policies to collect a dataset called *Atari-Seaquest-S*, consisting of 8000 trajectories for training and 2000 trajec-

ories for testing. The average length of the trajectories is 2652.5, and the average return is 2968.6. We cut the trajectory into subsequences with 200 states for training. To stabilize the training, we equip an orthogonal regularization for our method. Considering the output of the predictor is a matrix,  $\mathcal{M} \in \mathbb{R}^{b \times l}$  where  $n$  refers to the batch size and  $l$  is the length of  $\mathbf{m}$ , we drive the model to minimize the accumulation of  $\mathcal{M}\mathcal{M}^T$ . As the critical states of each trajectory are generally with different time steps, this regularization can benefit our approach. We train our Deep State Identifier on this video dataset and then test its effectiveness by re-training a new adaptive multi-step DQN from scratch, where the critical state detector adaptively determines the lookahead step. We use our trained critical state detector to determine the lookahead steps for rapid credit assignment during re-training.

## C. Experimental Results

To justify the effectiveness of the proposed method, we carry out some additional visualization and analysis. Table 10 shows some statistics of the output of the critical state detector and return predictor. We observe that the identified states are few (the L1 Norm is low), and the output of the return predictor does not change when it ignores non-critical states. If instead, the return predictor observes only states identified as non-critical, then the performance is much lower. These results further validate the effectiveness of the proposed method. We provide additional visualization of the performance of our method when using different losses for the critical state detector. The results are consistent with our empirical studies. In particular, Fig. 8(a) shows that when using only the importance preservation loss, all the states are considered critical. When adding only the compactness loss (see Fig. 8(b)) or the reverse loss (see Fig. 8(c)), the performance is still not satisfactory. The proposed method can precisely detect the critical states only when using all three losses. Indeed, as shown in Fig. 9, our method correctly outputs high confidence when the agent observes critical states (0.73, 0.94, and 0.92) and low confidence (0.6) otherwise.

### C.1. Non-Vision Environment

We also tested the performance in non-vision environments [3] and compared our method with the same methods in Table 5. As shown in Table 12, our method achieves a win rate change of **-45.10** on the MuJoCo [65] environment You-Should-Not-Pass game, surpassing the performance of EDGE (-35.13) by 28.38%. In the Kick-and-Defense environment, our method achieves a win rate change of **-48.03**, outperforming EDGE (-43.47) by 10.49%. The consistent improvement indicates that our method exhibits strong robustness in non-vision environments.

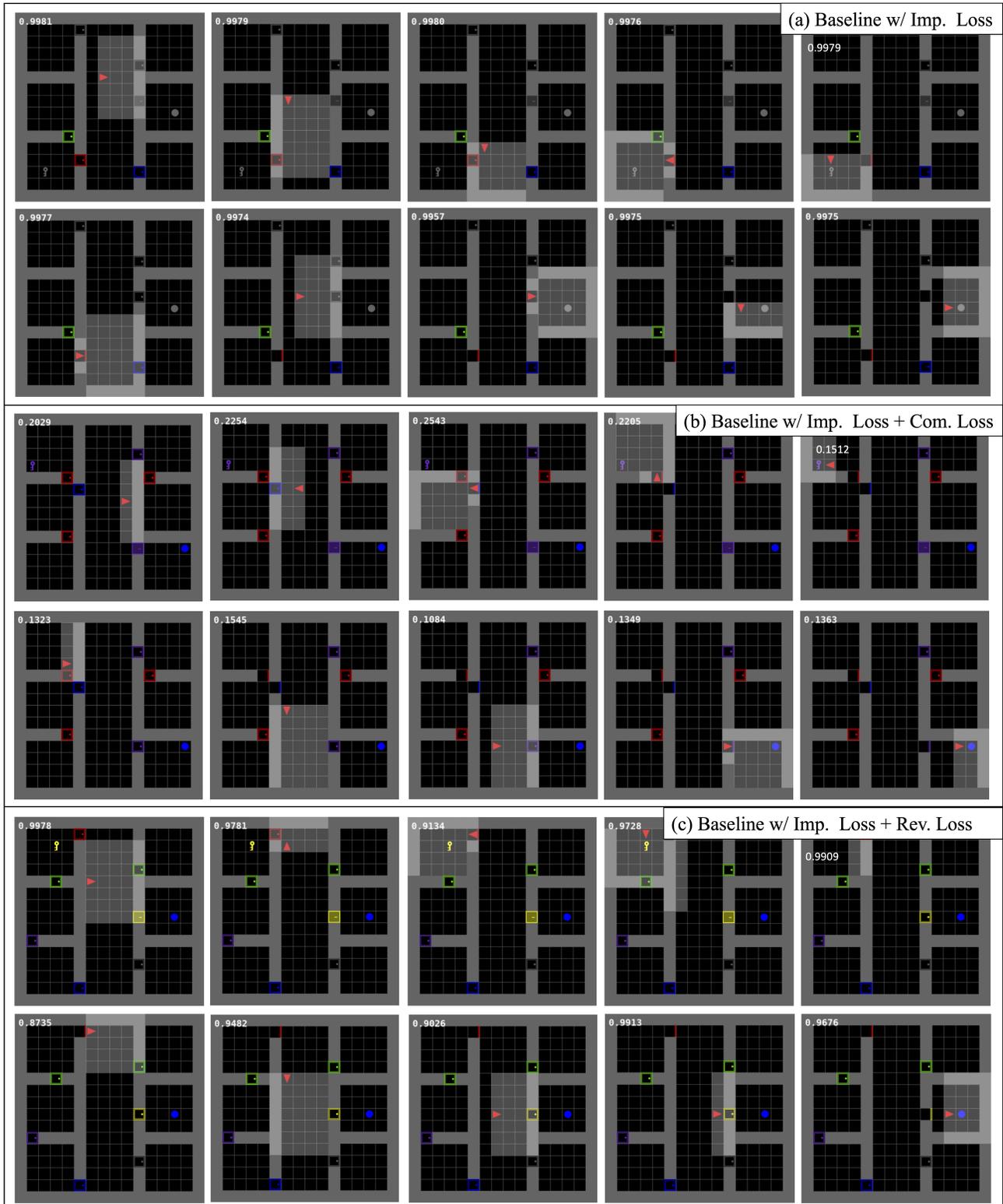


Figure 8. **Visualization of our method with different losses.** The number at the top-left corner indicates the confidence score predicted by the critical state detector, indicating whether the corresponding state is important. (a) Baseline trained with Importance Preservation loss; (b) Baseline with Importance Preservation loss and Compactness loss. (c) Baseline with Importance Preservation loss and Reverse loss. None of them can detect critical states effectively.

Table 10. **Ablation study for the Deep State Identifier.** Clean Acc. refers to the accuracy of the return predictor in the test set; Masked Acc. is the accuracy of the return predictor with the input (critical states) detected by the critical state detector; R-Masked Acc. is the accuracy of the return predictor where the masked is inverted (non-critical states are treated as critical and vice versa); L1(Mask) and Var(Mask) are the L1 norm and the average variance of the output of the critical state detector respectively.

Imp. Loss	Com. Loss	Rev. Loss	3DCNN	CNN-LSTM	Clean Acc. (%) ↑	Masked Acc.(%) ↑	R-Masked Acc.(%) ↓	L1(Mask) ↓	Var(Mask) ↑
✓	×	×	✓	×	90.07	90.07	<b>44.57</b>	63.74	$2 \times 10^{-6}$
✓	✓	×	✓	×	91.45	87.71	89.28	<b>8.79</b>	0.01
✓	×	✓	✓	×	91.45	91.39	76.12	63.73	0.03
✓	✓	✓	✓	×	90.78	89.45	64.55	57.35	0.04
✓	✓	✓	×	✓	<b>98.66</b>	<b>98.44</b>	55.58	41.05	<b>0.12</b>

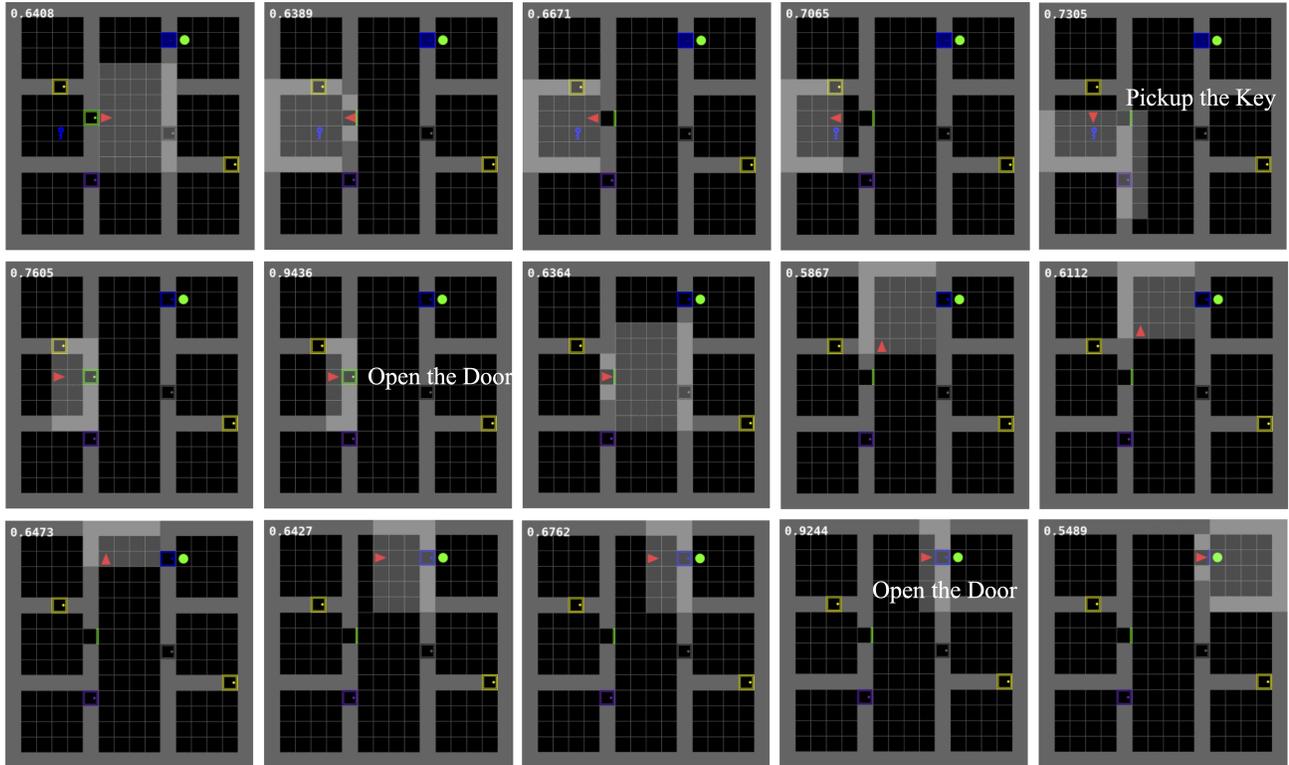


Figure 9. **Sampled observations from an episode collected by our method.** The number at the top-left corner indicates the confidence score predicted by the critical state detector, indicating whether the corresponding state is critical. Our method can localize the critical states effectively.

Table 11. **Sensitivity Analysis of the Deep State Identifier.** We show the F1 score ↑ of our method using different hyperparameters on GridWorld-S datasets.

$\lambda_r (\times 10^{-3})$	1	2.5	5	7.5	10	Variance
F1 Score	76.69	78.44	80.28	78.26	76.44	<b>1.39</b>
$\lambda_s$	0.5	0.75	1	1.25	1.5	Variance
F1 Score	76.68	77.50	80.28	78.18	78.83	<b>1.22</b>
$\lambda_v$	1.5	1.75	2	2.25	2.5	Variance
F1 Score	77.77	77.04	80.28	78.76	<b>83.78</b>	<b>2.39</b>

## C.2. Sensitivity analysis

We evaluate the performance of our method using different values of hyperparameters  $\lambda_r$ ,  $\lambda_s$ , and  $\lambda_v$ . Table 11 shows that our algorithm has moderate sensitivity to hyperparameters when their values are within a specific range.

Table 12. **Win rate changes of the agent before/after attacks by following the protocol of EDGE [21]** We compare the methods on two MuJoCo environments: You-Should-Not-Pass game [3] (MuJoCo-Y) and Kick-And-Defend game [3] (MuJoCo-K).

Method	MuJoCo-Y	MuJoCo-K
Rudder [1]	-32.53	-21.80
Saliency [57, 59, 60]	-29.33	-37.87
Attention RNN [2]	-33.93	-41.20
Rationale Net [36]	-30.00	-7.13
Edge [21]	-35.13	-43.47
Ours	<b>-45.10</b>	<b>-48.03</b>

For example, given  $\lambda_s \in [0.5, 1.5]$ , the performance variance is only 1.22, indicating stable performance. To determine the optimal hyperparameters, we searched a range

of values. The best hyperparameters found in GridWorld-S were then used in all other environments.