
Shielding Regular Safety Properties in Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To deploy reinforcement learning (RL) systems in real-world scenarios we need
2 to consider requirements such as safety and constraint compliance, rather than
3 blindly maximizing for reward. In this paper we study RL with regular safety
4 properties. We present a constrained problem based on the satisfaction of regular
5 safety properties with high probability and we compare our setup to the some
6 common constrained Markov decision processes (CMDP) settings. We also present
7 a meta-algorithm with provable safety-guarantees, that can be used to shield the
8 agent from violating the regular safety property during training and deployment.
9 We demonstrate the effectiveness and scalability of our framework by evaluating
10 our meta-algorithm in both the tabular and deep RL setting.

11 1 Introduction

12 The field of safe reinforcement learning (RL) [6, 28] has gained in-
13 creasing interest, as practitioners begin to understand the challenges
14 of applying RL in the real world [26]. There exist several distinct
15 paradigms in the literature, including constrained optimization
16 [2, 20, 49, 58, 62, 74], logical constraint satisfaction [17, 24, 36–
17 38, 66], safety-critical control [15, 19, 53], all of which are unified
18 by prioritizing safety- and risk-awareness during the decision making
19 process.

20 Constrained Markov decision processes (CMDP) [4] have emerged
21 as a popular framework for modelling safe RL, or RL with con-
22 straints. Typically, the goal is to obtain a policy that maximizes
23 reward while simultaneously ensuring that the expected cumulative cost remains below a pre-defined
24 threshold. A key limitation of this setting is that constraint violations are enforced in expectation
25 rather than with high probability, the constraint thresholds also have limited semantic meaning, can
26 be very challenging to tune and in some cases inappropriate for highly safety-critical scenarios
27 [66]. Furthermore, the cost function in the CMDP is typically Markovian and thus fails to capture a
28 significantly expressive class of safety properties and constraints.

29 Regular safety properties [9] are interesting because for all but the simplest properties the correspond-
30 ing cost function is non-Markovian. Our problem setup consists of the standard RL objective with
31 regular safety properties as constraints, we note that there has been a significant body of work that
32 combines temporal logic constraints with RL [17, 24, 36–38, 66], although many of these do not
33 explicitly separate reward and safety in the same way that we do.

34 Our approach relies on shielding [3], which is a safe exploration strategy that ensures the satisfaction
35 of temporal logic constraints by deploying the learned policy in conjunction with a reactive system

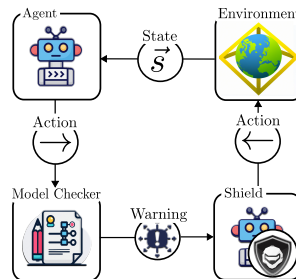


Figure 1: Diagrammatic representation of runtime verification and shielding.

36 that overrides any *unsafe* actions. Most shielding approaches typically make highly restrictive
37 assumptions, such as full knowledge of the environment dynamics [3], or access to a simulator [29],
38 although there has been recent work to deal with these restrictions [30, 39, 73]. In this paper, we
39 opt for the most permissive setting, where the dynamics of the environment are unknown, runtime
40 verification of the agent is realized by finite horizon model checking with a learned approximation of
41 the environment dynamics. However, in principle our framework is flexible enough to accommodate
42 more standard model checking procedures as long as certain assumptions are met.

43 Our approach can be summarised as an online shielding approach (see Fig. 1), that dynamically
44 identifies unsafe actions during training and deployment, and deploys a safe ‘backup policy’ when
45 necessary. We summarise the main contributions of our paper as follows:

46 (1) We state a constrained RL problem based on the satisfaction of regular safety properties with high
47 probability, and we identify the conditions whereby our setup generalizes several CMDP settings,
48 including *expected* and *probabilistic cumulative cost* constraints.

49 (2) We present several model checking algorithms that can verify the finite-horizon satisfaction
50 probability of regular safety properties, this includes statistical model checking procedures that can
51 be used if either the transition probabilities are unavailable or if the state space is too large.

52 (3) We develop a set of sample complexity results for the statistical model checking procedures
53 introduced in point (2), which are then used to develop a shielding meta-algorithm with provable
54 safety guarantees, even in the most permissive setting (i.e., no access to the transition probabilities).

55 (4) We empirically demonstrate the effectiveness of our framework on a variety of regular safety
56 properties in both a tabular and deep RL settings.

57 2 Related Work

58 **Safety Paradigms in Reinforcement Learning.** There exist many safety paradigms in RL, the most
59 popular being constrained MDPs. For CMDPs several constrained optimization algorithms have
60 been developed, most are gradient-based methods built upon Lagrange relaxations of the constrained
61 problem [20, 49, 58, 62] or projection-based local policy search [2, 74]. Model-based approaches to
62 CMDP [7, 11, 41, 64] have also gathered recent interest as they enjoy better sample complexity than
63 their model-free counterparts, which can be imperative for safe learning [44].

64 Linear Temporal Logic (LTL) constraints [17, 24, 36–38, 66] for RL have been developed as an
65 alternative to CMDPs to specify stricter and more expressive constraints. The LTL formula is typically
66 treated as the entire task specification, although some works have aimed to separate LTL satisfaction
67 and reward into two distinct objectives [66]. The typical procedure in this setting is to identify end
68 components of the MDP that satisfy the LTL constraint and construct a corresponding reward function
69 such that the optimal policy satisfies the LTL constraint with maximal probability. Formal PAC-style
70 guarantees have been developed for this setting [27, 36, 66, 71] although they typically rely on
71 non-trivial assumptions. We note that LTL constraints can capture regular safety properties, although
72 we explicitly separate reward and safety, making the work in this paper distinct from previous work.

73 More rigorous safety-guarantees can be obtained by using *safety filters* [3], *control barrier functions*
74 (CBF) [5], and *model predictive safety certification* (MPSC) [67, 68]. To achieve zero-violation
75 training these methods typically assume that the dynamics of the system are known and thus they
76 are typically restricted to low-dimensional systems. While these methods come from safety-critical
77 control, they are closely related to safe reinforcement learning [15].

78 **Learning Over Regular Structures.** RL and regular properties have been studied in conjunction
79 before, perhaps most famously as ‘Reward Machines’ [42, 43] – a type of finite state automaton that
80 specifies a different reward function at each automaton state. Reward machines do not explicitly
81 deal with safety, rather non-Markovian reward functions that depend on histories distinguished by
82 regular languages. Several methods have been developed to exploit the structure of these automata
83 and dramatically speed up learning [42, 43, 55, 61], e.g., *counterfactual experiences*.

84 Regular decision processes (RDP) [13] are a specific class non-Markovian DPs [8] that have also
85 been studied in several works [13, 22, 51, 59, 65]. Most of these works are theoretical and slightly
86 out-of-scope for this paper, as the RDP setting does not explicitly handle safety and encompasses
87 both non-Markovian rewards and transition probabilities.

88 **Shielding.** From formal methods, shielding for safe RL [3] forces hard constraints on policies, using
 89 a reactive system that ‘shields’ the agent from taking unsafe actions. Synthesising a *correct-by-*
 90 *construction* reactive ‘shield’ typically requires access to the environment dynamics and can be
 91 computationally demanding when the state or action space is large. Several recent works have aimed
 92 to scale the concept of shielding to more general settings, relaxing the prerequisite assumptions for
 93 shielding, by either only assuming access to a ‘black box’ model for planning [29], or learning a world
 94 model from scratch [30, 39, 73]. Other notable works that can be viewed as shielding include, MASE
 95 [69] – a safe exploration algorithm with access to an ‘emergency reset button’, and Recovery-RL
 96 [63] – which has access to a ‘recovery policy’ that is activated when the probability of reaching an
 97 unsafe state is too high. A simple form of shielding with LTL specifications has also been considered
 98 [37, 54], but experimentally these methods have only been tested in quite simple settings.

99 3 Preliminaries

100 For a finite set \mathcal{S} , let $Pow(\mathcal{S})$ denote the power set of \mathcal{S} . Also, let $Dist(\mathcal{S})$ denote the set of
 101 distributions over \mathcal{S} , where a distribution $\mu : \mathcal{S} \rightarrow [0, 1]$ is a function such that $\sum_{s \in \mathcal{S}} \mu(s) = 1$. Let
 102 \mathcal{S}^* and \mathcal{S}^ω denote the set of finite and infinite sequences over \mathcal{S} respectively. The set of all finite and
 103 infinite sequences is denoted $\mathcal{S}^\infty = \mathcal{S}^* \cup \mathcal{S}^\omega$. We denote as $|\rho|$ the length of a sequence $\rho \in \mathcal{S}^\infty$,
 104 where $|\rho| = \infty$ if $\rho \in \mathcal{S}^\omega$. We also denote as $\rho[i]$ the $i + 1$ -th element of a sequence, when $i < |\rho|$,
 105 and we denote as $\rho \downarrow = \rho[|\rho| - 1]$ the last element of a sequence, when $\rho \in \mathcal{S}^*$. A sequence ρ_1 is a
 106 prefix of ρ_2 , denoted $\rho_1 \preceq \rho_2$, if $|\rho_1| \leq |\rho_2|$ and $\rho_1[i] = \rho_2[i]$ for all $0 \leq i \leq |\rho_1|$. A sequence ρ_1 is
 107 a proper prefix of ρ_2 , denoted $\rho_1 \prec \rho_2$, if $\rho_1 \preceq \rho_2$ and $\rho_1 \neq \rho_2$.

108 **Labelled MDPs and Markov Chains.** An MDP is a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{P}_0, \mathcal{R}, AP, L)$, where
 109 \mathcal{S} and \mathcal{A} are finite sets of states and actions resp.; $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow Dist(\mathcal{S})$ is the *transition*
 110 *function*; $\mathcal{P}_0 \in Dist(\mathcal{S})$ is the *initial state distribution*; $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the *reward function*;
 111 AP is a set of *atomic propositions*, where $\Sigma = Pow(AP)$ is the *alphabet* over AP ; and $L : \mathcal{S} \rightarrow \Sigma$
 112 is a *labelling function*, where $L(s)$ denotes the set of atoms that hold in a given state
 113 $s \in \mathcal{S}$. A memory-less (stochastic) *policy* is a function $\pi : \mathcal{S} \rightarrow Dist(\mathcal{A})$ and its *value function*,
 114 denoted $V_\pi : \mathcal{S} \rightarrow \mathbb{R}$ is defined as the *expected reward* from a given state under policy π , i.e.,
 115 $V_\pi(s) = \mathbb{E}_\pi[\sum_{t=0}^T \mathcal{R}(s_t, a_t) | s_0 = s]$, where T is a fixed episode length. Furthermore, denote as
 116 $\mathcal{M}_\pi = (\mathcal{S}, \mathcal{P}_\pi, \mathcal{P}_0, AP, L)$ the *Markov chain* induced by a fixed policy π , where the transition
 117 function is such that $\mathcal{P}_\pi(s' | s) = \sum_{a \in \mathcal{A}} \mathcal{P}(s' | s, a) \pi(a | s)$. A path $\rho \in \mathcal{S}^\infty$ through \mathcal{M}_π is a finite (or
 118 infinite) sequence of states. Using standard results from measure theory it can be shown that the set
 119 of all paths $\{\rho \in \mathcal{S}^\omega \mid \rho_{pref} \preceq \rho\}$ with a common prefix ρ_{pref} is measurable [9].

120 **Probabilistic CTL.** (PCTL) [9] is a branching-time temporal logic for specifying properties of
 121 stochastic systems. A well-formed PCTL property can be constructed with the following grammar,

$$\begin{aligned} \Phi &::= \text{true} \mid a \mid \neg\Phi \mid \Phi \wedge \Phi \mid \mathbb{P}_{\bowtie p}[\varphi] \\ \varphi &::= X\Phi \mid \Phi U \Phi \mid \Phi U^{\leq n} \Phi \end{aligned}$$

122 where $a \in AP$, $\bowtie \in \{<, >, \leq, \geq\}$ is a binary comparison operator, and $p \in [0, 1]$ is a probability.
 123 Negation \neg and conjunction \wedge are the familiar logical operators from propositional logic, and next X ,
 124 until U and bounded until $U^{\leq n}$ are the temporal operators from CTL [9]. We make the distinction
 125 here between state formula Φ and path formula φ . The satisfaction relation for state formula Φ is
 126 defined in the standard way for Boolean connectives. For probabilistic quantification we say that
 127 $s \models \mathbb{P}_{\bowtie p}[\varphi]$ iff $\Pr(s \models \varphi) := \Pr(\rho \in \mathcal{S}^\omega \mid \rho[0] = s, \rho \models \varphi) \bowtie p$. Let $\Pr^{\mathcal{M}}(s \models \varphi)$ be the
 128 probability w.r.t. the Markov chain \mathcal{M} . For path formula φ the satisfaction relation is as follows,

$$\begin{aligned} \rho \models X\Phi & \quad \text{iff} \quad \rho[1] \models \Phi \\ \rho \models \Phi_1 U \Phi_2 & \quad \text{iff} \quad \exists j \geq 0 \text{ s.t. } (\rho[j] \models \Phi_2 \wedge \forall 0 \leq i < j, \rho[i] \models \Phi_1) \\ \rho \models \Phi_1 U^{\leq n} \Phi_2 & \quad \text{iff} \quad \exists 0 \leq j \leq n \text{ s.t. } (\rho[j] \models \Phi_2 \wedge \forall 0 \leq i < j, \rho[i] \models \Phi_1) \end{aligned}$$

129 From the standard operators of propositional logic we may derive disjunction \vee , implication \rightarrow and
 130 coimplication \leftrightarrow . We also note that the common temporal operators ‘eventually’ \diamond and ‘always’ \square ,
 131 and their bounded counterparts $\diamond^{\leq n}$ and $\square^{\leq n}$ can be derived in a familiar way, i.e., $\diamond \Phi ::= \text{true} U \Phi$,
 132 $\square \Phi ::= \neg \diamond \neg \Phi$, resp. $\diamond^{\leq n} \Phi ::= \text{true} U^{\leq n} \Phi$, $\square^{\leq n} \Phi ::= \neg \diamond^{\leq n} \neg \Phi$.

133 **Regular Safety Property.** A linear time property $P_{safe} \subseteq \Sigma^\omega$ over the alphabet Σ is a safety property
 134 if for all words $w \in \Sigma^\omega \setminus P_{safe}$, there exists a finite prefix w_{pref} of w such that $P_{safe} \cap \{w' \in \Sigma^\omega \mid$

135 $w_{pref} \preceq w'\} = \emptyset$. Any such sequence w_{pref} is called a *bad prefix* for P_{safe} , a bad prefix w_{pref}
 136 is called *minimal* iff there does not exist $w'' \prec w_{pref}$ such that w'' is a bad prefix for P_{safe} . Let
 137 $BadPref(P_{safe})$ and $MinBadPref(P_{safe})$ denote the set of bad and minimal bad prefixes resp.

138 A safety property $P_{safe} \in \Sigma^\omega$ is *regular* if the set $BadPref(P_{safe})$ constitutes a regular language. That
 139 is, there exists some *deterministic finite automata* (DFA) that accepts the bad prefixes for P_{safe} [9],
 140 that is, a path $\rho \in \mathcal{S}^\omega$ is ‘unsafe’ if the trace $trace(\rho) = L(\rho[0]), L(\rho[1]), \dots \in \Sigma^\omega$ is accepted by
 141 the corresponding DFA.

142 **Definition 3.1** (DFA). *A deterministic finite automata is a tuple $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$, where \mathcal{Q}
 143 is a finite set of states, Σ is a finite alphabet, $\Delta : \mathcal{Q} \times \Sigma \rightarrow \mathcal{Q}$ is the transition function, \mathcal{Q}_0 is the
 144 initial state, and $\mathcal{F} \subseteq \mathcal{Q}$ is the set of accepting states. The extended transition function Δ^* is the
 145 total function $\Delta^* : \mathcal{Q} \times \Sigma^* \rightarrow \mathcal{Q}$ defined recursively as $\Delta^*(q, w) = \Delta(\Delta^*(q, w \setminus w\downarrow), w\downarrow)$. The
 146 language accepted by DFA \mathcal{D} is denoted $\mathcal{L}(\mathcal{D}) = \{w \in \Sigma^* \mid \Delta^*(\mathcal{Q}_0, w) \in \mathcal{F}\}$.*

147 Furthermore, we denote as $P_{safe}^H \subseteq \Sigma^\omega$ the corresponding finite-horizon safety property for $H \in \mathbb{Z}_+$,
 148 where for all words $w \in \Sigma^\omega \setminus P_{safe}^H$ there exists $w_{pref} \preceq w$ such that $|w_{pref}| \leq H$ and $w_{pref} \in$
 149 $BadPref(P_{safe})$. We model check regular safety properties by synchronizing the DFA and Markov
 150 chain in a standard way – by computing the product Markov chain.

151 **Definition 3.2** (Product Markov Chain). *Let $\mathcal{M} = (\mathcal{S}, \mathcal{P}, \mathcal{P}_0, AP, L)$ be a Markov chain and $\mathcal{D} =$
 152 $(\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$ be a DFA. The product Markov chain is $\mathcal{M} \otimes \mathcal{D} = (\mathcal{S} \times \mathcal{Q}, \mathcal{P}', \mathcal{P}'_0, \{\text{accept}\}, L')$,
 153 where $L'(\langle s, q \rangle) = \{\text{accept}\}$ if $q \in \mathcal{F}$ and $L'(\langle s, q \rangle) = \emptyset$ o/w, $\mathcal{P}'_0(\langle s, q \rangle) = \mathcal{P}_0(s)$ if $q =$
 154 $\Delta(\mathcal{Q}_0, L(s))$ and 0 o/w, and $\mathcal{P}'(\langle s', q' \rangle | \langle s, q \rangle) = \mathcal{P}(s' | s)$ if $q' = \Delta(q, L(s'))$ and 0 o/w.*

155 To compute the satisfaction probability of P_{safe} for a given state $s \in \mathcal{S}$ we consider the set of paths
 156 $\rho \in \mathcal{S}^\omega$ from s and the corresponding trace in the DFA. We provide the following definition.

157 **Definition 3.3** (Satisfaction probability for P_{safe}). *Let $\mathcal{M} = (\mathcal{S}, \mathcal{P}, \mathcal{P}_0, AP, L)$ be a Markov chain
 158 and let $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$ be the DFA such that $\mathcal{L}(\mathcal{D}) = BadPref(P_{safe})$. For a path $\rho \in \mathcal{S}^\omega$
 159 in the Markov chain, let $trace(\rho) = L(\rho[0]), L(\rho[1]), \dots \in \Sigma^\omega$ be the corresponding word over
 160 $\Sigma = Pow(AP)$. From a given state $s \in \mathcal{S}$ the satisfaction probability for P_{safe} is defined as follows,*

$$\Pr^{\mathcal{M}}(s \models P_{safe}) := \Pr^{\mathcal{M}}(\rho \in \mathcal{S}^\omega \mid \rho[0] = s, trace(\rho) \notin \mathcal{L}(\mathcal{D}))$$

161 *Perhaps more importantly, we note that this satisfaction probability can be written as the following*
 162 *reachability probability in the product Markov chain,*

$$\Pr^{\mathcal{M}}(s \models P_{safe}) = \Pr^{\mathcal{M} \otimes \mathcal{D}}(\langle s, q_s \rangle \not\models \diamond \text{accept})$$

163 *where $q_s = \Delta(\mathcal{Q}_0, L(s))$ and $\diamond \text{accept}$ is a PCTL path formula that reads, ‘eventually accept’ [9].*

164 For the corresponding finite-horizon safety property P_{safe}^H we state the following result.

165 **Proposition 3.4** (Satisfaction probability for P_{safe}^H). *Let \mathcal{M} and \mathcal{D} be the MDP and DFA in Defn. 3.3.
 166 For a path $\rho \in \mathcal{S}^\omega$ in the Markov chain, let $trace_H(\rho) = L(\rho[0]), L(\rho[1]) \dots, L(\rho[H])$ be the
 167 corresponding finite word over $\Sigma = Pow(AP)$. For a given state $s \in \mathcal{S}$ the finite horizon satisfaction
 168 probability for P_{safe}^H is defined as follows,*

$$\Pr^{\mathcal{M}}(s \models P_{safe}^H) := \Pr^{\mathcal{M}}(\rho \in \mathcal{S}^\omega \mid \rho[0] = s, trace_H(\rho) \notin \mathcal{L}(\mathcal{D}))$$

169 *where $H \in \mathbb{Z}_+$ is some fixed model checking horizon. Similar to before, we show that the finite*
 170 *horizon satisfaction probability can be written as the following bounded reachability probability,*

$$\Pr^{\mathcal{M}}(s \models P_{safe}^H) = \Pr^{\mathcal{M} \otimes \mathcal{D}}(\langle s, q_s \rangle \not\models \diamond^{\leq H} \text{accept})$$

171 *where $q_s = \Delta(\mathcal{Q}_0, L(s))$ is as before and $\diamond^{\leq H} \text{accept}$ is the corresponding step-bounded PCTL path*
 172 *formula that reads, ‘eventually accept in H timesteps’.*

173 The unbounded reachability probability can be computed by solving a system of linear equations, the
 174 bounded reachability probability can be computed with $\mathcal{O}(H)$ matrix multiplications, in both cases
 175 the time complexity of the procedure is a polynomial in the size of the product Markov chain [9].

176 **4 Problem Setup**

177 In this paper, we are interested in the quantitative model checking of regular safety properties for
 178 a fixed finite horizon H and in the context of episodic RL, i.e., where the length of the episode T
 179 is fixed. In particular, at every timestep we constrain the (step-bounded) reachability probability
 180 $\Pr(\langle s, q \rangle \not\models \diamond^{\leq H} \text{accept})$ in the product Markov chain $\mathcal{M}_\pi \otimes \mathcal{D}$. We assume that H is chosen so as
 181 to avoid any irrecoverable states [35, 64], i.e., those that lead to a violation of the safety property no
 182 matter the sequence of actions taken, the precise details of this notion are presented in Section 6. We
 183 specify the following constrained problem,

184 **Problem 4.1** (Step-wise bounded regular safety property constraint). *Let P_{safe} be a regular safety*
 185 *property, \mathcal{D} be the DFA such that $\mathcal{L}(\mathcal{D}) = \text{BadPref}(P_{\text{safe}})$ and \mathcal{M} be the MDP;*

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T]$$

186 *where all probability is taken under the product Markov Chain $\mathcal{M}_\pi \otimes \mathcal{D}$, $p_1 \in [0, 1]$ is a probability*
 187 *threshold, H is the model checking horizon and T is the fixed episode length.*

188 The hyperparameter p_1 is be directly used to trade-off safety and exploration in a semantically
 189 meaningful way; p_1 prescribes the probability of satisfying the finite-horizon safety property P_{safe}^H
 190 at each timestep. In particular, if p_1 is sufficiently small then we can guarantee (with high-probability)
 191 that the regular safety property P_{safe} is satisfied for the entire episode length T .

192 **Proposition 4.2.** *Let P_{safe}^T denote the (episodic) regular safety property for a fixed episode length*
 193 *T . Then satisfying $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T]$ guarantees that $\Pr(s_0 \models$
 194 $P_{\text{safe}}^T) \geq 1 - p_1 \cdot \lceil T/H \rceil$, where $s_0 \sim \mathcal{P}_0$ is the initial state.*

195 **Comparison to CMDP.** In the remainder of this section, we compare our problem setup to various
 196 CMDP settings [4], with the aim of unifying different perspectives from safe RL. The purpose of this
 197 is to show that our proposed method for solving Problem 4.1 can also be used to satisfy other more
 198 common CMDP constraints. First, we define the following cost function that prescribes a scalar cost
 199 $C > 0$ when the regular safety property P_{safe} is violated and 0 otherwise.

200 **Definition 4.3** (Cost function). *Let P_{safe} be a regular safety property and let \mathcal{D} be the DFA such*
 201 *that $\mathcal{L}(\mathcal{D}) = \text{BadPref}(P_{\text{safe}})$, modified such that for all $q \in \mathcal{F}$, $q \rightarrow \mathcal{Q}_0$. The cost function is then*
 202 *defined as,*

$$\mathcal{C}(\langle s, q \rangle) = \begin{cases} C & \text{if } \text{accept} \in L'(\langle s, q \rangle) \\ 0 & \text{otherwise} \end{cases}$$

203 *where $C > 0$ is some generic scalar cost and L' is the labelling function defined in Def. 3.2.*

204 **Resetting the DFA.** Rather than reset the environment, the DFA is reset once it reaches an accepting
 205 state, so as to measure the rate of constraint satisfaction over a fixed episode length T . This can easily
 206 be realized by replacing any outgoing transitions from the accepting states with transitions back to
 207 the initial state, i.e., for all $q \in \mathcal{F}$, $q \rightarrow \mathcal{Q}_0$.

208 **Non-Markovian costs.** The cost function is Markov on the product states $\langle s, q \rangle \in \mathcal{S} \times \mathcal{Q}$. However,
 209 in most cases the cost function is non-Markovian in the original state space \mathcal{S} , since the automaton
 210 state $q \in \mathcal{Q}$ could depend on some arbitrary history of states. Thus our problem setup generalizes the
 211 standard CMDP framework with non-Markovian safety constraints.

212 **Invariant properties.** Invariant properties $P_{\text{inv}}(\Phi)$, also written $\square\Phi$ ('always Φ '), where Φ is a
 213 propositional state formula, are the simplest type of safety properties where the cost function is still
 214 Markov in the original state space. In this case we are operating in the standard CMDP framework,
 215 we also note that checking invariant properties with a fixed model checking horizon has been studied
 216 in previous works, as *bounded safety* [29, 30] and *safety for a finite horizon* [45].

217 The most common type of CMDP constraints are *expected cumulative (cost) constraints*, which
 218 constrain the expected cost below a given threshold.

Problem 4.4 (Expected cumulative constraint [4, 58]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{E}_{\langle s_t, q_t \rangle \sim \mathcal{M}_\pi \otimes \mathcal{D}} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq d_1$$

219 *where $d_1 \in \mathbb{R}_+$ is the cost threshold and T is the fixed episode length.*

220 *Probabilistic cumulative (cost) constraints*, are a stricter class of constraints that constrain the
 221 cumulative cost with high probability, rather than in expectation.

Problem 4.5 (Probabilistic cumulative constraint [18, 56]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{P}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \leq d_2 \right] \geq 1 - \delta_2$$

222 where $d_2 \in \mathbb{R}_+$ is the cost threshold, δ_2 is a tolerance parameter, and T is the fixed episode length.

223 We also consider *instantaneous constraints*, which bound the cost ‘almost surely’ at each timestep
 224 $t \in [0, T]$. These are an even stricter type of constraint for highly safety-critical applications.

Problem 4.6 (Instantaneous constraint [23, 60, 69]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{P}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} [\mathcal{C}(\langle s_t, q_t \rangle) \leq d_3] = 1 \quad \forall t \in [0, T]$$

225 where $d_3 \in \mathbb{R}_+$ is the cost threshold and T is the fixed episode length.

226 In particular, these problems define a constrained set of feasible policies Π . We make the distinction
 227 here between a feasible policy and a solution to the problem, the former being any policy satisfying
 228 the constraints of the problem and the later being the optimal policy within the feasible set Π .

229 **Theorem 4.7.** *A feasible policy for Problem 4.1 is also a feasible policy for Problems 4.4, 4.5 and*
 230 *4.6 under specific parameter settings for p_1, d_1, d_2 and δ_2 , and d_3 .*

231 In Appendix G we provide a full set of statements that outline the relationships between the con-
 232 strained problems presented in this section. The significance of these results is that they demonstrate
 233 by solving Problem 4.1 with our proposed method we can obtain feasible policies for Problems 4.4,
 234 4.5 and 4.6, although for most of these problems there is no direct relationship between our problem
 235 setup, in particular we can say little about whether the optimal policy for one problem is necessarily
 236 optimal for another. Nevertheless, we find it interesting to explore the relationships between our setup
 237 and other perhaps more common constrained RL problems.

238 5 Model checking

239 In this section we outline several procedures for checking the finite-horizon satisfaction probability
 240 of regular safety properties and we summarise the settings in which they can be used.

241 **Assumption 5.1.** *We are given access to the ‘true’ transition probabilities \mathcal{P} .*

242 **Assumption 5.2.** *We are given access to a ‘black box’ model that perfectly simulates the ‘true’*
 243 *transition probabilities \mathcal{P} .*

244 **Assumption 5.3.** *We are given access to an approximate dynamic model $\hat{\mathcal{P}} \approx \mathcal{P}$, where the total*
 245 *variation (TV) distance $D_{TV}(\mathcal{P}_{\pi}(\cdot | s), \hat{\mathcal{P}}_{\pi}(\cdot | s)) \leq \epsilon/H$, for all $s \in \mathcal{S}$.¹*

246 **Exact model checking.** Under Assumption 5.1 we can precisely compute the (finite horizon)
 247 satisfaction probability of P_{safe} , in the Markov chain \mathcal{M}_{π} induced by the fixed policy π in time
 248 $\mathcal{O}(\text{poly}(\text{size}(\mathcal{M}_{\pi} \otimes \mathcal{D})) \cdot H)$ [9], where \mathcal{D} is the DFA such that $\mathcal{L}(\mathcal{D}) = \text{BadPref}(P_{safe})$ and H
 249 is the model checking horizon. H should not be too large and so the complexity of exact model
 250 checking ultimately depends on the size of the product $\mathcal{M}_{\pi} \otimes \mathcal{D}$, and so if the size of either the MDP
 251 or DFA is too large then exact model checking may be infeasible.

252 **Monte-Carlo model checking.** To address the limitations of exact model checking, we can drop
 253 Assumption 5.1. Rather, under Assumption 5.2, we can sample sufficiently many paths from a
 254 ‘black box’ model of the environment dynamics and estimate the reachability probability $\Pr(\langle s, q \rangle \models$
 255 $\diamond^{\leq H} \text{accept})$ in the product Markov chain $\mathcal{M}_{\pi} \otimes \mathcal{D}$, by computing the proportion of accepting paths.
 256 Using statistical bounds, such as Hoeffding’s inequality [40] or Bernstein-type bounds [52], we can
 257 bound the error of this estimate, with high probability.

258 **Proposition 5.4.** *Let $\epsilon > 0, \delta > 0, s \in \mathcal{S}$ and $H \geq 1$ be given. Under Assumption 5.2, we can*
 259 *obtain an ϵ -approximate estimate for the probability $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability at*
 260 *least $1 - \delta$, by sampling $m \geq \frac{1}{2\epsilon^2} \log\left(\frac{2}{\delta}\right)$ paths from the ‘black box’ model.*

¹For two discrete probability distributions μ_1 and μ_2 over the same space \mathcal{X} the TV distance is defined as:
 $D_{TV}(\mu_1(\cdot), \mu_2(\cdot)) = \frac{1}{2} \sum_{x \in \mathcal{X}} |\mu_1(x) - \mu_2(x)|$

261 We note that the time complexity of these statistical methods does not depend in the size of the
 262 product MDP or DFA, since the product states $\langle s, q \rangle \in \mathcal{S} \times \mathcal{Q}$ can be computed *on-the-fly*, rather the
 263 time complexity depends on the horizon H , the desired level of accuracy ϵ , failure probability δ .

264 **Model checking with approximate models.** In most realistic cases neither the ‘true’ transition
 265 probabilities nor a perfect ‘black box’ model is available to us before-hand. Under Assumption
 266 5.3 we can model check with an ‘approximate’ model of the MDP dynamics, which can either be
 267 constructed ahead of time (offline) or learned from experience, with maximum likelihood (or similar).
 268 We can then either exact model check in with the ‘approximate’ probabilities, or if the MDP is too
 269 large, we can leverage statistical model checking by sampling paths from the ‘approximate’ model.

270 **Proposition 5.5.** *Let $\epsilon > 0$, $\delta > 0$, $s \in \mathcal{S}$ and $H \geq 1$ be given. Under Assumption 5.3 we can make*
 271 *the following two statements:*

272 (1) *We can obtain an ϵ -approximate estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability 1 by*
 273 *exact model checking with the transition probabilities of $\widehat{\mathcal{P}}_\pi$ in time $\mathcal{O}(\text{poly}(\text{size}(\mathcal{M}_\pi \otimes \mathcal{D})) \cdot H)$.*

274 (2) *We can obtain an ϵ -approximate estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability at least*
 275 *$1 - \delta$, by sampling $m \geq \frac{2}{\epsilon^2} \log\left(\frac{2}{\delta}\right)$ paths from the ‘approximate’ dynamics model $\widehat{\mathcal{P}}_\pi$.*

276 6 Shielding the policy

277 At a high-level, the shielding meta-algorithm
 278 works by switching between the ‘task policy’
 279 trained with RL to maximize rewards and a
 280 ‘backup policy’, which typically constitutes a
 281 low-reward, possibly rule-based policy that is
 282 guaranteed to be safe. In some cases this
 283 ‘backup policy’ may be available to us before
 284 training, although in most realistic cases it will
 285 need to be learned. In our case we switch
 286 from the ‘task policy’ to the ‘backup policy’
 287 when the reachability probability $\Pr(\langle s, q \rangle \models$
 288 $\diamond^{\leq H} \text{accept})$ exceeds the probability threshold
 289 p_1 . To check this we can use any of the model
 290 checking procedures presented earlier. The
 291 ‘backup policy’ is used when the reachability
 292 probability exceeds p_1 . Intuitively if the ‘backup
 293 policy’ is guaranteed to be safe, then our system
 294 should satisfy the constraints of Problem 4.1,
 295 independent of the ‘task policy’.

296 **Backup policy.** In general we assume no knowl-
 297 edge of the safety dynamics before training and
 298 so the ‘backup policy’ needs to be learned. In
 299 particular, we can use the cost function defined
 300 in Defn. 4.3 and train the ‘backup policy’ with
 301 RL to minimize the *expected discounted cost*
 302 $(\mathbb{E}_\pi[\sum_{t=0}^T \gamma^t \mathcal{C}(s_t, q_t)])$. Importantly, we note that the cost function is defined on the product state
 303 space $\mathcal{S} \times \mathcal{Q}$ and so the ‘backup policy’ must also operate on this state space, possibly leading
 304 to slower convergence. However, we can eliminate this issue entirely by training the ‘backup pol-
 305 icy’ with *counterfactual experiences* [42, 43] – a method originally used for reward machines that
 306 generates additional synthetic data for the policy, by simulating experience from each automaton
 307 state.

308 **Meta Algorithm.** We now present the structure of the shielding meta-algorithm (see Algorithm
 309 1). The precise realization of this algorithm can vary depending on problem setting, tabular, deep
 310 RL, etc., however the main structure of the algorithm remains the same. In particular, during
 311 interaction with the environment we shield the agent by checking that the reachability probability
 312 $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ does not exceed threshold p_1 . Then, with the new accumulated experience
 313 we update the ‘task policy’ denoted π_{task} and the ‘backup policy’ denoted π_{safe} with RL, and if need be

Algorithm 1 Shielding (with runtime verification of regular safety properties)

Input: model checking parameters (ϵ, δ, p, H) ,
 labelling function L , DFA $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$.
Optional: probabilities \mathcal{P} , ‘backup policy’ π_{safe} .

Initialize: ‘task policy’ π_{task} , ‘backup policy’ π_{safe}
 and (approximate) probabilities $\widehat{\mathcal{P}}$.

for each episode do

Observe $s_0, L(s_0)$ and $q_0 \leftarrow \Delta(\mathcal{Q}_0, L(s_0))$

for $t = 0, \dots, T$ **do** \triangleright Fixed episode length

Sample action $a \sim \pi_{\text{task}}(\cdot \mid s_t)$

if $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ **then**

// Use the proposed action

$a_t \leftarrow a$

else

// Override the action

$a_t \sim \pi_{\text{safe}}(\cdot \mid s_t, q_t)$

Play a_t and observe $s_{t+1}, L(s_{t+1}), r_t$

$q_{t+1} \leftarrow \Delta(q_t, L(s_{t+1}))$,

$c_t \leftarrow 1[q_{t+1} \in \mathcal{F}]$

Update π_{task} with (s_t, a_t, s_{t+1}, r_t)

Update π_{safe} with $(s_t, q_t, a_t, s_{t+1}, q_{t+1}, c_t)$

Update $\widehat{\mathcal{P}}$ with (s_t, a_t, s_{t+1})

314 we update our (approximate) dynamics model accordingly. In principle, the underlying RL algorithm
 315 used to train either ‘task policy’ or ‘backup policy’ can differ, and the dynamics model can be a
 316 simple maximum likelihood estimate or something more complex, e.g., Gaussian Process model
 317 [25, 70], ensemble of parametric neural networks [21, 44] or a world model [32, 33].

318 **Global Safety Guarantees.** In the tabular setting we can guarantee the safety of the system described
 319 in Algorithm 1 under various assumptions, even when doing Monte-Carlo model checking on an
 320 ‘approximate’ model of the environment dynamics. First, we provide the following definitions.

321 **Definition 6.1** (Non-critical state). *A product state $\langle s, q \rangle \in \mathcal{S} \times \mathcal{Q}$ is said to be non-critical for a*
 322 *given model checking horizon H if for all policies π we have $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) = 0$.*

323 **Definition 6.2** (Irrecoverable). *A critical state $\langle s, q \rangle \in \mathcal{S} \times \mathcal{Q}$ is said to be irrecoverable with*
 324 *probability p_1 if for all policies π we have $\Pr(\langle s, q \rangle \models \diamond \text{accept}) \geq p_1$. In other words, for any*
 325 *sequence of actions a_0, a_1, \dots the minimum probability $\Pr^{\min}(\langle s, q \rangle \models \diamond \text{accept})$ of reaching an*
 326 *accepting state is p_1 , where $\Pr^{\min}(\langle s, q \rangle \models \diamond \text{accept}) = \inf_{\pi} \Pr^{\mathcal{M}_{\pi} \otimes \mathcal{D}}(\langle s, q \rangle \models \diamond \text{accept})$*

327 The safety-guarantees for Algorithm 1 rely on the following assumptions.

328 **Assumption 6.3.** *We assume H is sufficiently large so that it is not possible to transition from any*
 329 *non-critical state to an irrecoverable state. Furthermore we assume that there exists some $H^* < H$*
 330 *such that if $\Pr^{\min}(\langle s, q \rangle \models \diamond \text{accept}) = p_1$ then $\Pr^{\min}(\langle s, q \rangle \models \diamond^{\leq H^*} \text{accept}) = p_1$.*

331 **Assumption 6.4.** *The initial state $\langle s_0, L(s_0) \rangle$ is non-critical and for any state $\langle s, q \rangle \in \mathcal{S} \times \mathcal{Q}$ that is*
 332 *not irrecoverable, the ‘backup policy’ π_{safe} satisfies $\Pr^{\mathcal{M}_{\pi_{\text{safe}}} \otimes \mathcal{D}}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$*

333 **Theorem 6.5.** *Under Assumption 6.3 and 6.4, and provided that every state action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$*
 334 *has been visited at least $\mathcal{O}\left(\frac{H^2 |S|^2}{\epsilon^2} \log\left(\frac{|A||S|^2}{\delta}\right)\right)$ times. Then with probability $1 - \delta$ the system*
 335 *satisfies the constraints of Problem 4.1, independent of the ‘task policy’.*

336 The theory is quite conservative here due to the strong dependence on $|S|$, in practice we can replace
 337 the outer $|S|^2$ by the maximum number of successor states from any given state. With regards to our
 338 assumptions, both are not overly restrictive. Assumption 6.3 essentially states that any irrecoverable
 339 states, will reach the accepting state with some probability $> g$ within a fixed horizon H^* . Similar
 340 statements have been considered in prior work [35, 64]. Assumption 6.4 states that the ‘backup
 341 policy’ satisfies $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ if possible, we would expect this to be the case when
 342 training the ‘backup policy’ with RL to minimize cost. The analysis for Theorem 6.5 then follows
 343 by showing that the system can be recovered to a non-critical state after entering a critical but not
 344 irrecoverable state.

345 7 Empirical Evaluation

346 We implemented two separate realizations of Algorithm 1, the first adapted to tabular environments
 347 which implements both exact or statistical model checking over the learned transition probabilities, the
 348 second is adapted to (visual) deep RL, making use of *world models* [32, 33], specifically DreamerV3
 349 [34], to learn a latent dynamics model for model checking and policy optimization.

350 **Tabular RL.** We conduct experiments on a simple ‘colour’ grid-
 351 world environment, with regular safety properties of increasing
 352 difficulty. In short, the goal is to navigate from a starting state
 353 to a goal position as frequently as possible, while respecting a
 354 given regular safety property during training. The environment is
 355 stochastic – with some probability p the agent’s action is ignored
 356 and another action is chosen uniformly instead. For smaller p val-
 357 ues the environment becomes more deterministic and the safety property typically becomes easier to
 358 satisfy with higher probability, we refer the reader to Appendix D.1 for more details. Table 1 outlines
 359 the three safety properties used for our environments. We use PCTL-like notation to describe the
 360 safety properties, although strictly speaking (2) and (3) are actually PCTL* path formula. Regardless
 361 of this slight technical detail, properties (1)-(3) are valid regular safety properties, as we can come up
 362 with a DFA that accepts the bad prefixes for them.

363 We compare our approach to Q-learning (without any penalties), and Q-learning on the product
 364 state space, with penalties provided by the cost function (Defn. 4.3) and trained with counterfactual

Table 1: Safety properties

property P_{safe}
(1) $\Box \neg \text{green}$
(2) $\Box \text{goal} \rightarrow \diamond^{\leq 10} \text{blue}$
(3) $\Box \text{goal} \rightarrow \diamond^{\leq 10} \Box^{\leq 5} \text{purple}$

365 experiences [43]. In all cases, by separating reward and safety into two distinct policies, we are able
 366 to effectively trade-off the two objectives. Q-learning simply finds the best policy ignoring the costs,
 367 and Q-learning with penalties is able to find a safe policy, but struggles to meaningfully balance both
 368 objectives (see Fig. 2). Hyperparameter settings for all experiments are detailed in Appendix E. In
 369 addition, we provide an extensive series of ablation studies in Appendix F for these experiments. For
 370 example, we show that we don't lose much by using Monte Carlo model checking as opposed to
 371 exact model checking with the 'true' probabilities. We also show that tuning the cost coefficient C
 372 offers no meaningful way to trade-off reward and the probability of constraint satisfaction.

373 **Deep RL.** We deploy our version of Algo-
 374 rithm 1 built on DreamerV3 [34] on Atari
 375 Seaquest, provided as part of the Arcade Learning
 376 Environment (ALE)[10, 50]. We experiment
 377 with two different regular safety properties:
 378 (1) $(\Box \neg \text{surface} \rightarrow \Box (\text{surface} \rightarrow \text{diver})) \wedge$
 379 $(\Box \neg \text{out-of-oxygen}) \wedge (\Box \neg \text{hit})$ and (2) $\Box \text{diver} \wedge$
 380 $\neg \text{surface} \rightarrow \Diamond_{\leq 30} \text{surface}$. We compare our approach
 381 to the base DreamerV3 algorithm and a version of
 382 DreamerV3 that implements the augmented Lagrangian
 383 penalty framework, similarly to [7, 41], for additional details see Ap-
 384 pendix B.1.

386 Again our approach is able to effectively trade-off
 387 both objectives, while (base) DreamerV3 ignores
 388 the cost, the Lagrangian approach appears to
 389 learn a safe policy that is not always efficient
 390 in terms of reward (see Fig. 3). We refer the
 391 reader to Appendix D.2 for more details of the
 392 environment and an extended discussion.

393 **Separating Reward and Safety.** The separation
 394 of reward and safety objectives into two distinct
 395 policies has been demonstrated as an effective
 396 strategy towards safety-aware decision making
 397 [3, 30, 46, 63], in many cases the safety ob-
 398 jective is simpler and can be more quickly learnt
 399 [46]. In our experiments it is clear that when
 400 the system enters a critical state, the 'backup
 401 policy' is able to efficiently guide the system
 402 back to a non-critical state where the task policy
 403 can continue collecting reward. However, there
 404 is evidence that the complete separation of poli-
 405 cies is not always appropriate [31] and penalties
 406 or a slight coupling of the policies is required
 407 to stop the 'task' and 'backup policy' fighting
 408 for control of the system. Furthermore, by separating
 409 reward and safety, we typically lose any
 410 asymptotic convergence guarantees, similar to the
 situation faced for hierarchical RL [61], although
 there has been recent work to develop convergence
 guarantees for shielding [75].

411 8 Conclusion

412 In this paper we propose a shielding meta-algorithm for the runtime verification of regular safety
 413 properties, given as a probabilistic constraint on the system. We provide a thorough theoretical
 414 examination of the problem and develop probabilistic safety guarantees for the meta-algorithm,
 415 which hold under reasonable assumptions. Empirically, we demonstrate that shielding is able to
 416 effectively balance both reward and safety, in both the tabular and deep RL setting. A more thorough
 417 theoretical and empirical examinations of the conditions for when shielding is appropriate would be
 418 an interesting direction for future work.

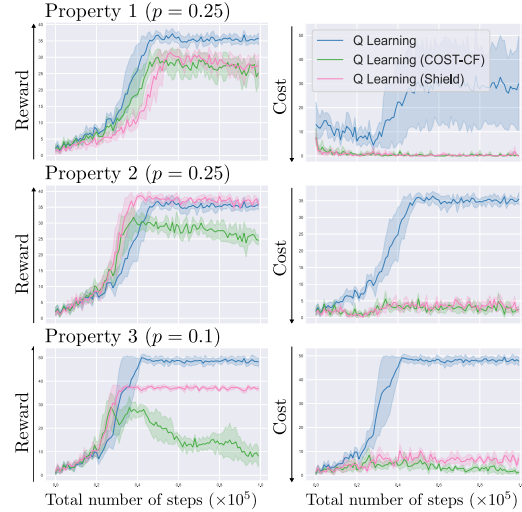


Figure 2: Episode reward and cost for tabular RL 'colour' gridworld environment.

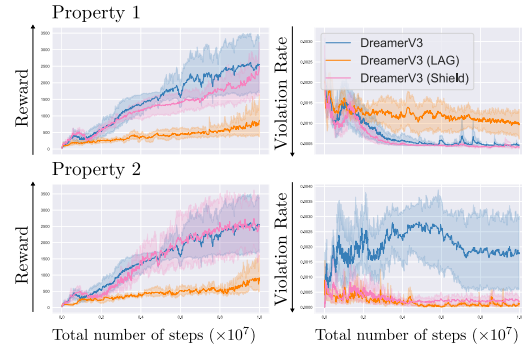


Figure 3: Episode reward and violation rate for deep RL Atari Seaquest.

References

- 419
- 420 [1] Pieter Abbeel and Andrew Y Ng. 2005. Exploration and apprenticeship learning in reinforcement
421 learning. In *Proceedings of the 22nd international conference on Machine learning*. 1–8.
- 422 [2] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. 2017. Constrained policy opti-
423 mization. In *International conference on machine learning*. PMLR, 22–31.
- 424 [3] Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum,
425 and Ufuk Topcu. 2018. Safe reinforcement learning via shielding. In *Proceedings of the AAAI*
426 *Conference on Artificial Intelligence*, Vol. 32.
- 427 [4] Eitan Altman. 1999. *Constrained Markov decision processes: stochastic modeling*. Routledge.
- 428 [5] Aaron D Ames, Samuel Coogan, Magnus Egerstedt, Gennaro Notomista, Koushil Sreenath,
429 and Paulo Tabuada. 2019. Control barrier functions: Theory and applications. In *2019 18th*
430 *European control conference (ECC)*. IEEE, 3420–3431.
- 431 [6] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
432 2016. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565* (2016).
- 433 [7] Yarden As, Ilnura Usmanova, Sebastian Curi, and Andreas Krause. 2022. Constrained policy
434 optimization via bayesian world models. *arXiv preprint arXiv:2201.09802* (2022).
- 435 [8] Fahiem Bacchus, Craig Boutilier, and Adam Grove. 1996. Rewarding behaviors. In *Proceedings*
436 *of the National Conference on Artificial Intelligence*. 1160–1167.
- 437 [9] Christel Baier and Joost-Pieter Katoen. 2008. *Principles of model checking*. MIT press.
- 438 [10] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. 2013. The Arcade Learning Environ-
439 ment: An Evaluation Platform for General Agents. *Journal of Artificial Intelligence Research*
440 47 (jun 2013), 253–279.
- 441 [11] Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. 2017. Safe
442 model-based reinforcement learning with stability guarantees. *Advances in neural information*
443 *processing systems* 30 (2017).
- 444 [12] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
445 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
446 Qiao Zhang. 2018. *JAX: composable transformations of Python+NumPy programs*. [http:
447 //github.com/google/jax](http://github.com/google/jax)
- 448 [13] Ronen I Brafman, Giuseppe De Giacomo, et al. 2019. Regular Decision Processes: A Model
449 for Non-Markovian Domains.. In *IJCAI*. 5516–5522.
- 450 [14] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
451 and Wojciech Zaremba. 2016. Openai gym. *arXiv preprint arXiv:1606.01540* (2016).
- 452 [15] Lukas Brunke, Melissa Greeff, Adam W Hall, Zhacong Yuan, Siqi Zhou, Jacopo Panerati,
453 and Angela P Schoellig. 2022. Safe learning in robotics: From learning-based control to safe
454 reinforcement learning. *Annual Review of Control, Robotics, and Autonomous Systems* 5 (2022),
455 411–444.
- 456 [16] Emma Brunskill, Bethany R Leffler, Lihong Li, Michael L Littman, and Nicholas Roy. 2009.
457 Provably efficient learning with typed parametric models. (2009).
- 458 [17] Mingyu Cai, Shaoping Xiao, Zhijun Li, and Zhen Kan. 2021. Optimal probabilistic motion
459 planning with potential infeasible LTL constraints. *IEEE transactions on automatic control* 68,
460 1 (2021), 301–316.
- 461 [18] Weiqin Chen, Dharmashankar Subramanian, and Santiago Paternain. 2024. Probabilistic
462 constraint for safety-critical reinforcement learning. *IEEE Trans. Automat. Control* (2024).
- 463 [19] Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. 2019. End-to-end safe
464 reinforcement learning through barrier functions for safety-critical continuous control tasks. In
465 *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 3387–3395.

- 466 [20] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. 2018. Risk-
467 constrained reinforcement learning with percentile risk criteria. *Journal of Machine Learning*
468 *Research* 18, 167 (2018), 1–51.
- 469 [21] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. 2018. Deep rein-
470 forcement learning in a handful of trials using probabilistic dynamics models. *Advances in*
471 *neural information processing systems* 31 (2018).
- 472 [22] Roberto Cipollone, Anders Jonsson, Alessandro Ronca, and Mohammad Sadegh Talebi. 2024.
473 Provably Efficient Offline Reinforcement Learning in Regular Decision Processes. *Advances in*
474 *Neural Information Processing Systems* 36 (2024).
- 475 [23] Gal Dalal, Krishnamurthy Dvijotham, Matej Vecerik, Todd Hester, Cosmin Paduraru, and Yuval
476 Tassa. 2018. Safe exploration in continuous action spaces. *arXiv preprint arXiv:1801.08757*
477 (2018).
- 478 [24] Giuseppe De Giacomo, Luca Iocchi, Marco Favorito, and Fabio Patrizi. 2020. Restraining
479 bolts for reinforcement learning agents. In *Proceedings of the AAAI Conference on Artificial*
480 *Intelligence*, Vol. 34. 13659–13662.
- 481 [25] Marc Deisenroth and Carl E Rasmussen. 2011. PILCO: A model-based and data-efficient
482 approach to policy search. In *Proceedings of the 28th International Conference on machine*
483 *learning (ICML-11)*. 465–472.
- 484 [26] Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. 2019. Challenges of real-world
485 reinforcement learning. *arXiv preprint arXiv:1904.12901* (2019).
- 486 [27] Jie Fu and Ufuk Topcu. 2014. Probably approximately correct MDP learning and control with
487 temporal logic constraints. *arXiv preprint arXiv:1404.7073* (2014).
- 488 [28] Javier Garcia and Fernando Fernández. 2015. A comprehensive survey on safe reinforcement
489 learning. *Journal of Machine Learning Research* 16, 1 (2015), 1437–1480.
- 490 [29] M Giacobbe, Mohammadhosein Hasanbeig, Daniel Kroening, and Hjalmar Wijk. 2021. Shield-
491 ing atari games with bounded prescience. In *Proceedings of the International Joint Conference*
492 *on Autonomous Agents and Multiagent Systems, AAMAS*.
- 493 [30] Alexander W Goodall and Francesco Belardinelli. 2023. Approximate Model-Based Shielding
494 for Safe Reinforcement Learning. *arXiv preprint arXiv:2308.00707* (2023).
- 495 [31] Alexander W Goodall and Francesco Belardinelli. 2024. Leveraging Approximate Model-based
496 Shielding for Probabilistic Safety Guarantees in Continuous Environments. *arXiv preprint*
497 *arXiv:2402.00816* (2024).
- 498 [32] David Ha and Jürgen Schmidhuber. 2018. Recurrent World Models Facilitate Policy Evo-
499 lution. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach,
500 H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran As-
501 sociates, Inc. [https://proceedings.neurips.cc/paper_files/paper/2018/file/](https://proceedings.neurips.cc/paper_files/paper/2018/file/2de5d16682c3c35007e4e92982f1a2ba-Paper.pdf)
502 [2de5d16682c3c35007e4e92982f1a2ba-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/2de5d16682c3c35007e4e92982f1a2ba-Paper.pdf)
- 503 [33] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and
504 James Davidson. 2019. Learning latent dynamics for planning from pixels. In *International*
505 *conference on machine learning*. PMLR, 2555–2565.
- 506 [34] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. 2023. Mastering diverse
507 domains through world models. *arXiv preprint arXiv:2301.04104* (2023).
- 508 [35] Alexander Hans, Daniel Schneegass, Anton Schäfer, and Steffen Udluft. 2008. Safe Exploration
509 for Reinforcement Learning. 143–148.
- 510 [36] Mohammadhosein Hasanbeig, Alessandro Abate, and Daniel Kroening. 2018. Logically-
511 constrained reinforcement learning. *arXiv preprint arXiv:1801.08099* (2018).
- 512 [37] Mohammadhosein Hasanbeig, Alessandro Abate, and Daniel Kroening. 2020. Cautious rein-
513 forcement learning with logical constraints. *arXiv preprint arXiv:2002.12156* (2020).

- 514 [38] Mohammadhosein Hasanbeig, Daniel Kroening, and Alessandro Abate. 2020. Deep reinforcement
515 learning with temporal logics. In *Formal Modeling and Analysis of Timed Systems: 18th*
516 *International Conference, FORMATS 2020, Vienna, Austria, September 1–3, 2020, Proceedings*
517 *18*. Springer, 1–22.
- 518 [39] Chloe He, Borja G León, and Francesco Belardinelli. [n.d.]. Do androids dream of electric
519 fences? Safety-aware reinforcement learning with latent shielding. *CEUR Workshop Proceed-*
520 *ings*. https://ceur-ws.org/Vol-3087/paper_50.pdf
- 521 [40] Wassily Hoeffding. 1963. Probability Inequalities for Sums of Bounded Random Variables. *J.*
522 *Amer. Statist. Assoc.* 58, 301 (1963), 13–30.
- 523 [41] Weidong Huang, Jiaming Ji, Borong Zhang, Chunhe Xia, and Yaodong Yang. 2023. Safe Dream-
524 erV3: Safe Reinforcement Learning with World Models. *arXiv preprint arXiv:2307.07176*
525 (2023).
- 526 [42] Rodrigo Toro Icarte, Toryn Klassen, Richard Valenzano, and Sheila McIlraith. 2018. Using
527 reward machines for high-level task specification and decomposition in reinforcement learning.
528 In *International Conference on Machine Learning*. PMLR, 2107–2116.
- 529 [43] Rodrigo Toro Icarte, Toryn Q Klassen, Richard Valenzano, and Sheila A McIlraith. 2022.
530 Reward machines: Exploiting reward function structure in reinforcement learning. *Journal of*
531 *Artificial Intelligence Research* 73 (2022), 173–208.
- 532 [44] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. 2019. When to trust your model:
533 Model-based policy optimization. *Advances in neural information processing systems* 32
534 (2019).
- 535 [45] Nils Jansen, Bettina Könighofer, Sebastian Junges, Alex Serban, and Roderick Bloem. 2020.
536 Safe reinforcement learning using probabilistic shields. In *31st International Conference on*
537 *Concurrency Theory (CONCUR 2020)*. Schloss-Dagstuhl-Leibniz Zentrum für Informatik.
- 538 [46] Nils Jansen, Bettina Könighofer, Sebastian Junges, Alexandru C Serban, and Roderick Bloem.
539 2018. Safe reinforcement learning via probabilistic shields. *arXiv preprint arXiv:1807.06096*
540 (2018).
- 541 [47] Sham Kakade, Michael J Kearns, and John Langford. 2003. Exploration in metric state spaces. In
542 *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*. 306–312.
- 543 [48] Michael Kearns and Satinder Singh. 2002. Near-optimal reinforcement learning in polynomial
544 time. *Machine learning* 49 (2002), 209–232.
- 545 [49] Qingkai Liang, Fanyu Que, and Eytan Modiano. 2018. Accelerated primal-dual policy opti-
546 mization for safe reinforcement learning. *arXiv preprint arXiv:1802.06480* (2018).
- 547 [50] Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, and
548 Michael Bowling. 2018. Revisiting the Arcade Learning Environment: Evaluation Protocols
549 and Open Problems for General Agents. *Journal of Artificial Intelligence Research* 61 (2018),
550 523–562.
- 551 [51] Sultan Javed Majeed, Marcus Hutter, et al. 2018. On Q-learning Convergence for Non-Markov
552 Decision Processes.. In *IJCAI*, Vol. 18. 2546–2552.
- 553 [52] Andreas Maurer and Massimiliano Pontil. 2009. Empirical bernstein bounds and sample
554 variance penalization. *arXiv preprint arXiv:0907.3740* (2009).
- 555 [53] Stephen McIlvanna, Nhat Nguyen Minh, Yuzhu Sun, Mien Van, and Wasif Naem. 2022.
556 Reinforcement learning-enhanced control barrier functions for robot manipulators. *arXiv*
557 *preprint arXiv:2211.11391* (2022).
- 558 [54] Rohan Mitta, Hosein Hasanbeig, Jun Wang, Daniel Kroening, Yiannis Kantaros, and Alessandro
559 Abate. 2024. Safeguarded Progress in Reinforcement Learning: Safe Bayesian Exploration
560 for Control Policy Synthesis. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
561 Vol. 38. 21412–21419.

- 562 [55] Andrew Y Ng, Daishi Harada, and Stuart Russell. 1999. Policy invariance under reward
563 transformations: Theory and application to reward shaping. In *Icml*, Vol. 99. 278–287.
- 564 [56] Santiago Paternain, Miguel Calvo-Fullana, Luiz FO Chamon, and Alejandro Ribeiro. 2022. Safe
565 policies for reinforcement learning via primal-dual methods. *IEEE Trans. Automat. Control* 68,
566 3 (2022), 1321–1336.
- 567 [57] Aravind Rajeswaran, Igor Mordatch, and Vikash Kumar. 2020. A game theoretic framework for
568 model based reinforcement learning. In *International conference on machine learning*. PMLR,
569 7953–7963.
- 570 [58] Alex Ray, Joshua Achiam, and Dario Amodei. 2019. Benchmarking safe exploration in deep
571 reinforcement learning. *arXiv preprint arXiv:1910.01708* 7, 1 (2019), 2.
- 572 [59] Alessandro Ronca and Giuseppe De Giacomo. 2021. Efficient PAC reinforcement learning in
573 regular decision processes. *arXiv preprint arXiv:2105.06784* (2021).
- 574 [60] Yanan Sui, Alkis Gotovos, Joel Burdick, and Andreas Krause. 2015. Safe exploration for
575 optimization with Gaussian processes. In *International conference on machine learning*. PMLR,
576 997–1005.
- 577 [61] Richard S Sutton, Doina Precup, and Satinder Singh. 1999. Between MDPs and semi-MDPs: A
578 framework for temporal abstraction in reinforcement learning. *Artificial intelligence* 112, 1-2
579 (1999), 181–211.
- 580 [62] Chen Tessler, Daniel J Mankowitz, and Shie Mannor. 2018. Reward constrained policy opti-
581 mization. *arXiv preprint arXiv:1805.11074* (2018).
- 582 [63] Brijen Thananjeyan, Ashwin Balakrishna, Suraj Nair, Michael Luo, Krishnan Srinivasan, Minh
583 Hwang, Joseph E Gonzalez, Julian Ibarz, Chelsea Finn, and Ken Goldberg. 2021. Recovery
584 rl: Safe reinforcement learning with learned recovery zones. *IEEE Robotics and Automation*
585 *Letters* 6, 3 (2021), 4915–4922.
- 586 [64] Garrett Thomas, Yuping Luo, and Tengyu Ma. 2021. Safe reinforcement learning by imagining
587 the near future. *Advances in Neural Information Processing Systems* 34 (2021), 13859–13869.
- 588 [65] Rodrigo Toro Icarte, Ethan Waldie, Toryn Klassen, Rick Valenzano, Margarita Castro, and
589 Sheila McIlraith. 2019. Learning reward machines for partially observable reinforcement
590 learning. *Advances in neural information processing systems* 32 (2019).
- 591 [66] Cameron Voloshin, Hoang Le, Swarat Chaudhuri, and Yisong Yue. 2022. Policy optimization
592 with linear temporal logic constraints. *Advances in Neural Information Processing Systems* 35
593 (2022), 17690–17702.
- 594 [67] Kim P Wabersich and Melanie N Zeilinger. 2018. Linear model predictive safety certification
595 for learning-based control. In *2018 IEEE Conference on Decision and Control (CDC)*. IEEE,
596 7130–7135.
- 597 [68] Kim Peter Wabersich and Melanie N Zeilinger. 2021. A predictive safety filter for learning-based
598 control of constrained nonlinear dynamical systems. *Automatica* 129 (2021), 109597.
- 599 [69] Akifumi Wachi, Yanan Sui, Yisong Yue, and Masahiro Ono. 2018. Safe exploration and opti-
600 mization of constrained mdps using gaussian processes. In *Proceedings of the AAAI Conference*
601 *on Artificial Intelligence*, Vol. 32.
- 602 [70] Christopher KI Williams and Carl Edward Rasmussen. 2006. *Gaussian processes for machine*
603 *learning*. Vol. 2. MIT press Cambridge, MA.
- 604 [71] Eric M Wolff, Ufuk Topcu, and Richard M Murray. 2012. Robust control of uncertain Markov
605 decision processes with temporal logic specifications. In *2012 IEEE 51st IEEE Conference on*
606 *decision and control (CDC)*. IEEE, 3372–3379.
- 607 [72] Jorge Nocedal Stephen J Wright. 2006. Numerical optimization.

- 608 [73] Wenli Xiao, Yiwei Lyu, and John Dolan. 2023. Model-based Dynamic Shielding for Safe
609 and Efficient Multi-agent Reinforcement Learning. In *Proceedings of the 2023 International*
610 *Conference on Autonomous Agents and Multiagent Systems*. 1587–1596.
- 611 [74] Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J Ramadge. 2020. Projection-
612 based constrained policy optimization. *arXiv preprint arXiv:2010.03152* (2020).
- 613 [75] Wen-Chi Yang, Giuseppe Marra, Gavin Rens, and Luc De Raedt. 2023. Safe reinforcement
614 learning via probabilistic logic shields. *arXiv preprint arXiv:2303.03226* (2023).

Algorithm 2 Exact Model Checking [9]**Input:** model checking parameters (p, H) , current state $\langle s, q \rangle$, current action a , product MC $\mathcal{M}_\pi \otimes \mathcal{D} = (\mathcal{S} \times \mathcal{Q}, \mathcal{P}', \mathcal{P}'_0, \{\text{accept}\}, L')$ **Output:** *true* if $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ Initialize zero vector $\mathbf{x}^{(0)} \leftarrow \mathbf{0}$ with size $|\mathcal{S}| \times |\mathcal{Q}|$ Initialize probability matrix $\mathbf{A} \leftarrow (\mathcal{P}'(s, t))_{s, t \notin \text{accept}}$ (ignoring accepting states)Initialize probability vector $\mathbf{b} \leftarrow (\mathcal{P}'(s, \text{accept}))_{s \notin \text{accept}}$ (going to accepting states)

// Iterate over the model checking horizon

for $i = 1, \dots, H$ **do** Compute $\mathbf{x}^{(i)} = \mathbf{A}\mathbf{x}^{(i-1)} + \mathbf{b}$

// Get the corresponding probability

Let $X \leftarrow \mathbf{x}_{\langle s, q \rangle}$ **If** $X < p$ **return true** **else return false**

Algorithm 3 Monte-Carlo Model Checking**Input:** model checking parameters (ϵ, δ, p, H) , current state $\langle s, q \rangle$, current action a , policy π , labelling function L , DFA $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$ and (approximate) transition probabilities \mathcal{P} **Output:** *true* if $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ Choose $m \geq 2/(\epsilon^2) \log(2/\delta)$ **for** $i = 1, \dots, m$ **do** Set $s_0 \leftarrow s, q_0 \leftarrow q$ and $a_0 \leftarrow a$

// Sample a path through the model

for $j = 1, \dots, H$ **do** Sample next state $s_j \sim \mathcal{P}(\cdot \mid s_{j-1}, a_{j-1})$, Compute $q_j \leftarrow \Delta(q_{j-1}, L(s_j))$, Sample action $a_j \sim \pi(\cdot \mid s_j)$

// Check if the path is accepting

 Let $X_i \leftarrow 1[q_H \in \mathcal{F}]$

// Construct probability estimate

Let $\tilde{X} \leftarrow \frac{1}{m} \sum_{i=1}^m X_i$ **If** $\tilde{X} < p - \epsilon$ **return true** **else return false**

Algorithm 4 Tabular Q-learning (Regular Safety Property) with Counter Factual Experiences [65]**Input:** MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{P}_0, \mathcal{R}, AP, L)$, DFA $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$, discount factor $\gamma \in (0, 1]$, learning rate $\alpha \in (0, 1]$, temperature $\tau > 0$, cost coefficient C and fixed episode length T **Initialize:** (Q-table) $\hat{Q}(s, q, a) \leftarrow 0 \forall s \in \mathcal{S}, q \in \mathcal{Q}, a \in \mathcal{A}$ **for each episode do** Observe $s_0, L(s_0)$ and $q_0 \leftarrow \Delta(\mathcal{Q}_0, L(s_0))$ **for** $t = 0, \dots, T$ **do** Sample action a_t from $\langle s_t, q_t \rangle$ using the Boltzmann policy derived from \hat{Q} with temp. τ Play action a_t and observe $s_{t+1}, L(s_{t+1})$ and r_t (reward is optional).

// Generate synthetic data by simulating all automaton transitions

for $\bar{q} \in \mathcal{Q}$ **do** Compute $\bar{q}' \leftarrow \Delta(q', L(s_{t+1}))$ Compute cost $\bar{c}' \leftarrow C \cdot 1[\bar{q}' \in \mathcal{F}]$ Compute *done* $\leftarrow 1[\bar{q}' \in \mathcal{F}]$

// Q-learning step

 $\hat{Q}(s_t, \bar{q}, a_t) \leftarrow (1 - \alpha) \cdot \hat{Q}(s_t, \bar{q}, a_t) + \alpha \cdot (r_t + \bar{c}' + \gamma \cdot \text{done} \cdot \max_{a' \in \mathcal{A}} \hat{Q}(s_{t+1}, \bar{q}', a'))$ Compute $q_{t+1} \leftarrow \Delta(q_t, L(s_{t+1}))$ and continue

Algorithm 5 DreamerV3 [34] with Shielding (Regular Safety Property)

Initialize: replay buffer D with S random episodes, world model parameters θ , ‘task policy’ π_{task} and ‘backup policy’ π_{safe} randomly.

for each episode do

Observe $o_0, L(s_0)$ and $q_0 \leftarrow \Delta(Q_0, L(s_0))$

for $t = 1, \dots, T$ **do**

Sample action $a \sim \pi_{task}$ from the task policy

// Estimate the reachability probability using the world model p_θ

if $\Pr(\langle s, q \mid \diamond^{\leq H} \text{accept} \rangle \leq p_1)$ **then**

Use proposed action

$a_t \leftarrow a$

else

// Override action

$a_t \sim \pi_{safe}$

Play action a_t and observe $o_{t+1}, L(s_{t+1})$ and r_t

Compute $q_{t+1} \leftarrow \Delta(q_t, L(s_{t+1}))$,

Compute cost $c_t \leftarrow 1[q_{t+1} \in \mathcal{F}]$

Append $(o_t, a_t, r_t, c_t, o_{t+1})$ to the replay buffer D

if update then

// World model learning

Sample a batch B of transition sequences $\{(o_{t'}, a_{t'}, r_{t'}, c_{t'}, o_{t'+1})\} \sim \mathcal{D}$.

Update the world model parameters θ with maximum likelihood.

// Task policy optimization

‘Imagine’ sequences $\{\hat{o}_{t':t'+H}, \hat{r}_{t':t'+H}, \hat{c}_{t':t'+H}\}$ with the ‘task policy’ π_{task}

Update the ‘task policy’ π_{task} with RL (to maximize reward).

Update the corresponding value critics with maximum likelihood

// Backup policy optimization

‘Imagine’ sequences $\{\hat{o}_{t':t'+H}, \hat{r}_{t':t'+H}, \hat{c}_{t':t'+H}\}$ with the ‘backup policy’ π_{safe}

Update the ‘backup policy’ π_{safe} with RL (to minimize cost)

Update the corresponding value critics with maximum likelihood

616 B Technical Details

617 B.1 Augmented Lagrangian

618 We first define the following objective functions,

$$J_{\mathcal{R}}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \mathcal{R}(s_t, a_t) \right] \quad (1)$$

$$J_{\mathcal{C}}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \mathcal{C}(s_t, a_t) \right] \quad (2)$$

619 The augmented Lagrangian [72] is an adaptive penalty-based technique for the following constrained
620 optimization problem,

$$\max_{\pi} J_{\mathcal{R}}(\pi) \quad \text{subject to} \quad J_{\mathcal{C}}(\pi) \leq d \quad (4)$$

621 where d is some cost threshold. The corresponding Lagrangian is given by,

$$\max_{\pi} \min_{\lambda \geq 0} [J_{\mathcal{R}}(\pi) - \lambda (J_{\mathcal{C}}(\pi) - d)] = \max_{\pi} \begin{cases} J_{\mathcal{R}}(\pi) & \text{if } J_{\mathcal{C}}(\pi) < d \\ -\infty & \text{otherwise} \end{cases} \quad (5)$$

622 The LHS is an equivalent form for the constrained optimization problem (RHS), since if π is feasible,
623 i.e. $J_{\mathcal{C}}(\pi) < d$ then the maximum value for λ is $\lambda = 0$. If π is not feasible then λ can be arbitrarily
624 large to solve this equation. Unfortunately this form of the objective function is non-smooth when
625 moving from feasible to infeasible policies, thus we introduce a proximal relaxation of the augmented

626 Lagrangian [72],

$$\max_{\pi} \min_{\lambda \geq 0} \left[J_{\mathcal{R}}(\pi) - \lambda (J_{\mathcal{C}}(\pi) - d) + \frac{1}{\mu_k} (\lambda - \lambda_k)^2 \right] \quad (6)$$

627 where μ_k is a non-decreasing penalty multiplier dependent on the gradient step k . The new term
628 that has been introduced here encourages the λ to stay close to the previous value λ_k , resulting in a
629 smooth and differentiable function. The derivative w.r.t λ gives us the following gradient update step,

$$\lambda_{k+1} = \begin{cases} \lambda_k + \mu_k (J_{\mathcal{C}}(\pi) - d) & \text{if } \lambda_k + \mu_k (J_{\mathcal{C}}(\pi) - d) \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

630 At each gradient step, the penalty multiplier μ_k is updated in a non-decreasing way by using some
631 small fixed (power) parameter σ ,

$$\mu_{k+1} = \max\{(\mu_k)^{1+\sigma}, 1\} \quad (8)$$

632 The policy π is then updated by taking gradient steps of the following unconstrained objective,

$$\tilde{J}(\pi, \lambda_k, \mu_k) = J_{\mathcal{R}}(\pi) - \Psi_{\mathcal{C}}(\pi, \lambda_k, \mu_k)$$

633 where,

$$\Psi_{\mathcal{C}}(\pi, \lambda_k, \mu_k) = \begin{cases} \lambda_k (J_{\mathcal{C}}(\pi) - d) + \frac{\mu_k}{2} (J_{\mathcal{C}}(\pi) - d)^2 & \text{if } \lambda_k + \mu_k (J_{\mathcal{C}}(\pi) - d) \geq 0 \\ -\frac{(\lambda_k)^2}{2\mu_k} & \text{otherwise} \end{cases}$$

634 C Technical Proofs

635 C.1 Proof of Proposition 3.4

636 **Proposition 3.4 (restated)** (Satisfaction probability for P_{safe}^H). *Let \mathcal{M} and \mathcal{D} be the MDP and*
637 *DFA from before (Defn. 3.3). For a path $\rho \in \mathcal{S}^\omega$ in the Markov chain, let $trace_H(\rho) =$*
638 *$L(\rho[0]), L(\rho[1]) \dots, L(\rho[H])$ be the corresponding finite word over $\Sigma = Pow(AP)$. For a given*
639 *state $s \in \mathcal{S}$ the finite horizon satisfaction probability for P_{safe} is defined as follows,*

$$\Pr^{\mathcal{M}}(s \models P_{safe}^H) := \Pr^{\mathcal{M}}(\rho \in \mathcal{S}^\omega \mid \rho[0] = s, trace_H(\rho) \notin \mathcal{L}(\mathcal{D}))$$

640 where $H \in \mathbb{Z}_+$ is some fixed model checking horizon. Similar to before, we show that the finite
641 horizon satisfaction probability can be written as the following bounded reachability probability,

$$\Pr^{\mathcal{M}}(s \models P_{safe}^H) = \Pr^{\mathcal{M} \otimes \mathcal{D}}(\langle s, q_s \rangle \not\models \diamond^{\leq H} \text{accept})$$

642 where $q_s = \Delta(\mathcal{Q}_0, L(s))$ is as before and $\diamond^{\leq H} \text{accept}$ is the corresponding step-bounded PCTL path
643 formula that reads, ‘eventually accept in H timesteps’.

644 *Proof.* Let P_{safe} be a regular safety property and let $\mathcal{D} = (\mathcal{Q}, \Sigma, \Delta, \mathcal{Q}_0, \mathcal{F})$ be the DFA such that
645 $\mathcal{L}(\mathcal{D}) = BadPref(P_{safe})$. We provide a formal definition for P_{safe} and the corresponding finite
646 horizon property P_{safe}^H , respectively:

$$P_{safe} = \{w \in \Sigma^\omega \mid \forall w_{pref} \in \Sigma^\omega \text{ s.t. } w_{pref} \preceq w, w_{pref} \notin \mathcal{L}(\mathcal{D})\} \quad (9)$$

$$P_{safe}^H = \{w \in \Sigma^\omega \mid \forall w_{pref} \in \Sigma^\omega \text{ s.t. } w_{pref} \preceq w \wedge |w_{pref}| \leq H + 1, w_{pref} \notin \mathcal{L}(\mathcal{D})\} \quad (10)$$

647 Let $\mathcal{M} = (\mathcal{S}, \mathcal{P}, \mathcal{P}_0, AP, L)$ be a Markov chain and consider the product Markov chain $\mathcal{M} \otimes \mathcal{D}$
648 from Defn. 3.2. For any path $\rho = s_0, s_1, s_2, \dots$, there exists a unique run q_0, q_1, q_2, \dots for the trace
649 $trace(\rho) = L(s_0), L(s_1), L(s_2) \dots$, and denote,

$$\rho^+ = \langle s_0, q_0 \rangle, \langle s_1, q_1 \rangle, \langle s_2, q_2 \rangle \dots \quad (11)$$

650 where start state is $\langle s_0, \Delta(\mathcal{Q}_0, L(s_0)) \rangle$. Before we deal with probabilities let’s just consider a
651 fixed path $\rho \in \mathcal{S}^\omega$, the finite trace $trace_H(\rho) = L(\rho[0]), L(\rho[1]) \dots, L(\rho[H])$, the unique run
652 $q_0, q_1, q_2, \dots, q_H$ and the path $\rho^+ \in \Sigma^\omega \times \mathcal{Q}^\omega$ in the product Markov chain. We prove the following
653 statement,

$$\rho \not\models P_{safe}^H \text{ if and only if } \rho^+ \models \diamond \text{accept}^{\leq H} \quad (12)$$

654 We start with the (\rightarrow) direction, in particular, $\rho \not\models P_{safe}^H$ if and only if $trace_H(\rho) \in \mathcal{L}(\mathcal{D})$. Recall
655 that by definition $\mathcal{L}(\mathcal{D}) = \{w \in \Sigma^* \mid \Delta^*(\mathcal{Q}_0, w) \in \mathcal{F}\}$, and so $trace_H(\rho) \in \mathcal{L}(\mathcal{D})$ implies that
656 $q_H = \Delta^*(\mathcal{Q}_0, trace_H(\rho)) \in \mathcal{F}$, which by construction implies that $\rho^+ \models \diamond_{accept}^{\leq H}$.

657 The opposite direction (\leftarrow) is a little more involved, in particular, $\rho^+ \models \diamond_{accept}^{\leq H}$ implies that
658 for the unique run $q_0, q_1, q_2, \dots, q_H$ there exists $t \leq H$ such that $q_t \in \mathcal{F}$. We notice that since
659 $\mathcal{L}(\mathcal{D}) = BadPref(P_{safe})$ then once the DFA reaches an accepting state it will remain in an accepting
660 state for the rest of the run. Therefore, $q_t \in \mathcal{F}$ for $t \leq H$ implies that $q_H \in \mathcal{F}$. Then by definition
661 the trace $trace_H(\rho)$ that determined the unique run $q_0, q_1, q_2, \dots, q_H$ must be in the language $\mathcal{L}(\mathcal{D})$,
662 which again by definition implies that $\rho \not\models P_{safe}^H$.

663 We now deal with the probabilities. First we note that the DFA \mathcal{D} does not affect the probabilities of
664 the product Markov chain – it can be shown that for every measurable set P of paths in \mathcal{M} ,

$$\Pr^{\mathcal{M}}(P) = \Pr^{\mathcal{M} \otimes \mathcal{A}}(\rho^+ \mid \rho \in P) \quad (13)$$

665 see [9]. It now remains to construct this set P in the proper way. In particular, if P is the set of paths
666 starting in some state $s \in \mathcal{S}$ and that refute P_{safe} in the next H timesteps, i.e.,

$$P = \{\rho \in \mathcal{S}^\omega \mid \rho[0] = s, \{w' \in \Sigma^* \mid w_{pref} \preceq trace(\rho) \wedge |w_{pref}| \leq H + 1\} \cap \mathcal{L}(\mathcal{D}) \neq \emptyset\} \quad (14)$$

667 and P^+ is defined as the set of paths starting from the corresponding state $\langle s, q_s \rangle$ (where $q_s =$
668 $\Delta(\mathcal{Q}_0, L(s))$) in $\mathcal{M} \otimes \mathcal{D}$ that eventually reach an accepting state of \mathcal{D} in the next H steps, i.e.

$$P^+ = \{\rho^+ \in (\mathcal{S} \times \mathcal{Q})^\omega \mid \rho^+[0] = \langle s, q_s \rangle \wedge \rho^+ \models \diamond^{\leq H} accept\} \quad (15)$$

669 Then by construction we have,

$$\Pr^{\mathcal{M}}(P) = \Pr^{\mathcal{M} \otimes \mathcal{D}}(\rho^+ \mid \rho[0] = s, \rho \in P) = \Pr^{\mathcal{M} \otimes \mathcal{D}}(P^+) \quad (16)$$

670 Finally the probability $\Pr^{\mathcal{M}}(P)$ and $\Pr^{\mathcal{M}}(s \models P_{safe}^H)$ are related as follows,

$$\Pr^{\mathcal{M}}(s \models P_{safe}^H) = 1 - \Pr^{\mathcal{M}}(P) \quad (17)$$

$$= 1 - \Pr^{\mathcal{M} \otimes \mathcal{D}}(P^+) \quad (18)$$

$$= 1 - \Pr^{\mathcal{M} \otimes \mathcal{D}}(\langle s, q_s \rangle \models \diamond^{\leq H} accept) \quad (19)$$

$$= \Pr^{\mathcal{M} \otimes \mathcal{D}}(\langle s, q_s \rangle \not\models \diamond^{\leq H} accept) \quad (20)$$

671 □

672 C.2 Proof of Proposition 4.2

673 **Proposition 4.2 (restated).** *Let P_{safe}^T denote the (episodic) regular safety property for a fixed episode*
674 *length T . Then satisfying $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} accept) \leq p_1$ for all $t \in [0, T]$ guarantees that*
675 *$\Pr(s_0 \models P_{safe}^T) \geq 1 - p_1 \cdot \lceil T/H \rceil$, where $s_0 \sim \mathcal{P}_0$ is the initial state.*

676 *Proof.* Consider splitting up the episode in to $\lceil T/H \rceil$ chunks with length at most H . Let
677 $X_0, X_1, \dots, X_{\lceil T/H \rceil - 1}$ be the indicator random variables defined as follows,

$$X_i = \begin{cases} 1 & \text{if } \langle s_{i \cdot H}, q_{i \cdot H} \rangle \models \diamond^{\leq H} accept \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

678 Since $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} accept) \leq p_1$ for all $t \in [0, T]$ then the probability $\Pr(X_i = 1) \leq p_1$. By
679 construction we have,

$$\text{if } \bigcap_{i=0}^{\lceil T/H \rceil - 1} X_i = 0 \text{ then } s_0 \models P_{safe}^T \quad (22)$$

680 Intuitively we satisfy P_{safe} for the entire episode length if we never enter an accepting state in each of
 681 the $\lceil T/H \rceil$ chunks. The final result is then obtained by taking a union bound as follows,

$$\Pr(s_0 \models P_{safe}^T) \geq \Pr\left(\bigcap_{i=0}^{\lceil T/H \rceil - 1} X_i = 0\right) \quad (23)$$

$$= 1 - \Pr\left(\bigcup_{i=0}^{\lceil T/H \rceil - 1} X_i = 1\right) \quad (24)$$

$$\geq 1 - \sum_{i=0}^{\lceil T/H \rceil - 1} \Pr(X_i = 1) \quad (25)$$

$$\geq 1 - p_1 \cdot \lceil T/H \rceil \quad (26)$$

$$(27)$$

682

□

683 C.3 Proof of Proposition 5.4

684 **Proposition 5.4 (restated).** *Let $\epsilon > 0$, $\delta > 0$, $s \in \mathcal{S}$ be given. Under Assumption 5.2, we can obtain*
 685 *an ϵ -approximate estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability at least $1 - \delta$, by sampling*
 686 *$m \geq \frac{1}{2\epsilon^2} \log\left(\frac{2}{\delta}\right)$ paths from the ‘black box’ model.*

687 *Proof.* In words, we estimate $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ by sampling m paths from a ‘black box’
 688 model of the environment dynamics. We label each path as satisfying or not and return the proportion
 689 of satisfying traces as an estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$. We proceed as follows, let ρ_1, \dots, ρ_m
 690 be a sequence of paths sampled from the ‘black box’ model and let $\text{trace}(\rho_1), \dots, \text{trace}(\rho_m)$ be the
 691 corresponding traces. Furthermore, let X_1, \dots, X_m be indicator r.v.s such that,

$$X_i = \begin{cases} 1 & \text{if } \text{trace}(\rho_i) \models \diamond^{\leq H} \text{accept}, \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

692 Recall that $\text{trace}(\rho_i) \models \diamond^{\leq H} \text{accept}$ can be checked in time $O(\text{poly}(H))$. Now let,

$$\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i \text{ where } \mathbb{E}[\bar{X}] = \Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \quad (29)$$

693 then by Hoeffding’s inequality [40],

$$\mathbb{P}\left[|\bar{X} - \mathbb{E}[\bar{X}]| \geq \epsilon\right] \leq 2 \exp(-2m\epsilon^2) \quad (30)$$

694 Bounding the RHS from above by δ and rearranging gives the desired result. □

695 C.4 Proof of Proposition 5.5

696 We start by introducing the following lemma.

697 **Lemma C.1** (Error amplification for trace distributions). *Let $\hat{\mathcal{P}} \approx \mathcal{P}$ be such that,*

$$D_{TV}\left(\mathcal{P}(\cdot | s), \hat{\mathcal{P}}(\cdot | s)\right) \leq \alpha \forall s \in \mathcal{S} \quad (31)$$

698 *Let the start state $s_0 \in \mathcal{S}$ be given, and let $\mathcal{P}_t(\cdot)$ and $\hat{\mathcal{P}}_t(\cdot)$ denote the path distribution (at time t) for*
 699 *the two transition probabilities \mathcal{P} and $\hat{\mathcal{P}}$ respectively. Then the total variation distance between the*
 700 *two path distributions (at time t) are bounded as follows,*

$$D_{TV}\left(\mathcal{P}_t(\cdot), \hat{\mathcal{P}}_t(\cdot)\right) \leq \alpha t \forall t \quad (32)$$

701 *Proof.* We will prove this fact by doing an induction on t . We recall that $\mathcal{P}_t(\cdot)$ and $\widehat{\mathcal{P}}_t(\cdot)$ denote the
702 path distribution (at time t) for the two transition probabilities \mathcal{P} and $\widehat{\mathcal{P}}$ respectively. Formally we
703 define them as follows,

$$\mathcal{P}_t(\rho) = \Pr(s_0, \dots, s_t \preceq \rho \mid s_0 = s, \mathcal{P}) \quad (33)$$

$$\widehat{\mathcal{P}}_t(\rho) = \Pr(s_0, \dots, s_t \preceq \rho \mid s_0 = s, \widehat{\mathcal{P}}) \quad (34)$$

704 These probabilities read as follows, ‘the probability of the sequence $s_0, \dots, s_t \preceq \rho$ at time t ’, or
705 similarly ‘the probability that the sequence s_0, \dots, s_t is a prefix of ρ at time t ’. Since the start state
706 $s_0 \in \mathcal{S}$ is given we note that,

$$\mathcal{P}_0(\cdot) = \widehat{\mathcal{P}}_0(\cdot) \quad (35)$$

707 Before we continue with the induction on t we make the following observation, for any path $\rho \in \mathcal{S}^\omega$
708 we have by the triangle inequality,

$$\left| \mathcal{P}_t(\rho) - \widehat{\mathcal{P}}_t(\rho) \right| = \left| \mathcal{P}(s_t \mid s_{t-1})\mathcal{P}_{t-1}(\rho) - \widehat{\mathcal{P}}(s_t \mid s_{t-1})\widehat{\mathcal{P}}_{t-1}(\rho) \right| \quad (36)$$

$$\leq \mathcal{P}_{t-1}(\rho) \left| \mathcal{P}(s_t \mid s_{t-1}) - \widehat{\mathcal{P}}(s_t \mid s_{t-1}) \right| + \widehat{\mathcal{P}}(s_t \mid s_{t-1}) \left| \mathcal{P}_{t-1}(\rho) - \widehat{\mathcal{P}}_{t-1}(\rho) \right| \quad (37)$$

709 Now we continue with the induction on t ,

$$2D_{TV}(\mathcal{P}_t(\cdot), \widehat{\mathcal{P}}_t(\cdot)) = \sum_{\rho \in \mathcal{S}^\omega} \left| \mathcal{P}_t(\rho) - \widehat{\mathcal{P}}_t(\rho) \right| \quad (38)$$

$$\leq \sum_{\rho \in \mathcal{S}^\omega} \mathcal{P}_{t-1}(\rho) \left| \mathcal{P}(s_t \mid s_{t-1}) - \widehat{\mathcal{P}}(s_t \mid s_{t-1}) \right| + \sum_{\rho \in \mathcal{S}^\omega} \widehat{\mathcal{P}}(s_t \mid s_{t-1}) \left| \mathcal{P}_{t-1}(\rho) - \widehat{\mathcal{P}}_{t-1}(\rho) \right| \quad (39)$$

$$\leq \sum_{\rho \in \mathcal{S}^\omega} \mathcal{P}_{t-1}(\rho) \cdot (2\alpha) + \sum_{\rho \in \mathcal{S}^\omega} \left| \mathcal{P}_{t-1}(\rho) - \widehat{\mathcal{P}}_{t-1}(\rho) \right| \quad (40)$$

$$= 2\alpha + 2D_{TV}(\mathcal{P}_{t-1}(\cdot), \widehat{\mathcal{P}}_{t-1}(\cdot)) \quad (41)$$

$$\leq 2\alpha t \quad (42)$$

710 The final result is obtained by an induction on t where the base case comes from $\mathcal{P}_0(\cdot) = \widehat{\mathcal{P}}_0(\cdot)$. \square

711 **Proposition 5.5 (restated).** Let $\epsilon > 0$, $\delta > 0$, $s \in \mathcal{S}$ and horizon $H \geq 1$ be given. Under Assumption
712 5.3 we can make the following two statements:

713 (1) We can obtain an ϵ -approximate estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability 1 by
714 exact model checking with the transition probabilities of $\widehat{\mathcal{P}}_\pi$ in time $\mathcal{O}(\text{poly}(\text{size}(\mathcal{M}_\pi \otimes \mathcal{D})) \cdot H)$.

715 (2) We can obtain an ϵ -approximate estimate for $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with probability at least
716 $1 - \delta$, by sampling $m \geq \frac{2}{\epsilon^2} \log\left(\frac{2}{\delta}\right)$ paths from the ‘approximate’ dynamics model $\widehat{\mathcal{P}}_\pi$.

717 *Proof.* We start by proving statement (1) and then statement (2) will follow quickly. First let
718 $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ and $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ denote the acceptance probabilities for the
719 two transition probabilities \mathcal{P} and $\widehat{\mathcal{P}}$ respectively. We also let $g(\cdot)$ and $\widehat{g}(\cdot)$ denote the average trace
720 distribution (over the next H timesteps) for the two transition probabilities \mathcal{P} and $\widehat{\mathcal{P}}$ respectively,
721 where,

$$g(\rho) = \frac{1}{H} \sum_{t=1}^H \mathcal{P}_t(\rho) \quad (43)$$

$$\widehat{g}(\rho) = \frac{1}{H} \sum_{t=1}^H \widehat{\mathcal{P}}_t(\rho) \quad (44)$$

722 Before we continue with the proof of (1) we make the following observations,

723 • $\max_{\langle s, q \rangle} \left| \Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) - \widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \right| \leq 1$

724 • Let $f(x) : x \in \mathcal{X} \rightarrow [0, 1]$ be a real-valued function. Let $\mathcal{P}_1(\cdot)$ and $\mathcal{P}_2(\cdot)$ be probability
725 distributions over the space \mathcal{X} , then.

$$\left| \mathbb{E}_{x \sim \mathcal{P}_1(\cdot)}[f(x)] - \mathbb{E}_{x \sim \mathcal{P}_2(\cdot)}[f(x)] \right| \leq D_{TV}(\mathcal{P}_1(\cdot), \mathcal{P}_2(\cdot))$$

726 We continue by showing the following,

$$\left| \Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) - \widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \right| \quad (45)$$

$$= \left| \mathbb{E}_{\rho \sim g} [1[\langle s, q \rangle \models \diamond^{\leq H} \text{accept}]] - \mathbb{E}_{\rho \sim \widehat{g}} [1[\langle s, q \rangle \models \diamond^{\leq H} \text{accept}]] \right| \quad (46)$$

$$\leq D_{TV}(g(\cdot), \widehat{g}(\cdot)) \quad (47)$$

$$= \frac{1}{2} \sum_{\rho \in \mathcal{S}^\omega} |g(\rho) - \widehat{g}(\rho)| \quad (48)$$

$$= \frac{1}{2H} \sum_{\rho \in \mathcal{S}^\omega} \left| \sum_{t=1}^H \mathcal{P}_t(\rho) - \widehat{\mathcal{P}}_t(\rho) \right| \quad (49)$$

$$\leq \frac{1}{2H} \sum_{t=1}^H \left| \sum_{\rho \in \mathcal{S}^\omega} \mathcal{P}_t(\rho) - \widehat{\mathcal{P}}_t(\rho) \right| \quad (50)$$

$$\leq \frac{1}{2H} \sum_{t=1}^H H(\epsilon/H) \quad (51)$$

$$= \epsilon/2 \quad (52)$$

$$(53)$$

727 The first inequality (Eq. 47) comes from our earlier observations. The second inequality (Eq. 50) is
728 straightforward and the final inequality (Eq. 51) is obtained by applying Lemma C.1 and Assumption
729 5.3. We note that this result is similar to the *simulation lemma* [48], which has been proved many
730 times for several different settings [1, 16, 47, 57].

731 This concludes the proof of statement (1), since we have shown that $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ is an
732 $\epsilon/2$ -approximate estimate of $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$, under the Assumption 5.3.

733 The proof of statement (2) follows quickly. We have established that,

$$\left| \Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) - \widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \right| \leq \epsilon/2 \quad (54)$$

734 It remains to obtain an $\epsilon/2$ -approximate estimate of $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$. By using the
735 same reasoning as in the proof of Proposition 5.4. We can obtain an $\epsilon/2$ -approximate estimate
736 of $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ by sampling m paths, ρ_1, \dots, ρ_m , from the approximate dynamics model
737 $\widehat{\mathcal{P}}$. Then provided,

$$m \geq \frac{2}{\epsilon^2} \log \left(\frac{2}{\delta} \right) \quad (55)$$

738 with probability $1 - \delta$ we can obtain $\epsilon/2$ -approximate estimate of $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ and by
739 extension an ϵ -approximate estimate of $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$. This concludes the proof. \square

740 C.5 Proof of Theorem 6.5

741 **Theorem 6.5 (restated).** *Under Assumption 6.3 and 6.4, and provided that every state action pair*
742 *$(s, a) \in \mathcal{S} \times \mathcal{A}$ has been visited at least $\mathcal{O} \left(\frac{H^2 |\mathcal{S}|^2}{\epsilon^2} \log \left(\frac{|\mathcal{A}| |\mathcal{S}|^2}{\delta} \right) \right)$ times. Then with probability $1 - \delta$*
743 *the system satisfies the constraints of Problem 4.1, independent of the ‘task policy’.*

744 *Proof.* We split the proof up in to three parts, (1), (2) and (3). In part (1) we show that the given
745 sample complexity bound gives us an approximate model of the environment dynamics with high

746 probability. In part **(2)** we use our assumptions to reason about the probabilistic recoverability
 747 of the system when it enters a critical state. In part **(3)** we put everything together and deal with
 748 approximation error ϵ the remaining failure probability that are both unavoidable for the statistical
 749 model checking procedures used to shield the system.

750 **(1)** We show that the following holds with probability $1 - \delta/2$,

$$D_{TV} \left(\mathcal{P}_\pi(\cdot | s), \widehat{\mathcal{P}}_\pi(\cdot | s) \right) \leq \epsilon/H \quad \forall s \in \mathcal{S} \quad (56)$$

751 when every state action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ has been visited at least,

$$\mathcal{O} \left(\frac{H^2 |\mathcal{S}|^2}{\epsilon^2} \log \left(\frac{|\mathcal{A}| |\mathcal{S}|^2}{\delta} \right) \right)$$

752 times. First we let $\#(s, a)$ denote the total number of times that (s, a) has been observed, similarly
 753 we let $\#(s', s, a)$ denote the total number of times that (s', s, a) has been observed. The maximum
 754 likelihood estimate for the unknown probability $\mathcal{P}(s' | s, a)$ is $\widehat{\mathcal{P}}(s' | s, a) = \#(s', s, a) / \#(s, a)$.
 755 Let us fix some $(s, a) \in \mathcal{S} \times \mathcal{A}$, and $s' \in \mathcal{S}$, we let $p_{s'} = \mathcal{P}(s' | s, a)$ denote the true probability of
 756 transitioning to s' from (s, a) and we let $\hat{p}_{s'} = \#(s', s, a) / \#(s, a)$ denote our estimate. We note that
 757 $\mathbb{E}[\hat{p}_{s'}] = p_{s'}$, i.e. $\hat{p}_{s'}$ is an unbiased estimator for $p_{s'}$. Let $m = \#(s, a)$ also be the number of times
 758 that (s, a) has been observed, then by Hoeffding's inequality [40] we have,

$$\mathbb{P} \left[|p_{s'} - \hat{p}_{s'}| \geq \frac{\epsilon}{H|\mathcal{S}|} \right] \leq 2 \exp \left(-2m \frac{\epsilon^2}{H^2 |\mathcal{S}|^2} \right) \quad (57)$$

759 Bounding the LHS from above by $1 - \delta/2(|\mathcal{A}| |\mathcal{S}|^2)$ and rearranging gives the following lower bound
 760 for m ,

$$m \geq \frac{H^2 |\mathcal{S}|^2}{2\epsilon^2} \log \left(\frac{4|\mathcal{A}| |\mathcal{S}|^2}{\delta} \right) \quad (58)$$

761 Taking a union bound over all $(s', s, a) \in \mathcal{S} \times \mathcal{S} \times \mathcal{A}$, then for all state action pairs $(s, a) \in \mathcal{S} \times \mathcal{A}$
 762 we have the following with probability at least $1 - \delta$.

$$2D_{TV} \left(\mathcal{P}(\cdot | s, a), \widehat{\mathcal{P}}(\cdot | s, a) \right) = \sum_{s' \in \mathcal{S}} |p_{s'} - \hat{p}_{s'}| \leq \sum_{s' \in \mathcal{S}} \frac{\epsilon}{H|\mathcal{S}|} \leq \epsilon/H \quad (59)$$

763 Now fix some $s \in \mathcal{S}$ and we observe the following,

$$2D_{TV} \left(\mathcal{P}_\pi(\cdot | s), \widehat{\mathcal{P}}_\pi(\cdot | s) \right) = \sum_{s' \in \mathcal{S}} |\mathcal{P}_\pi(s' | s) - \widehat{\mathcal{P}}_\pi(s' | s)| \quad (60)$$

$$= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}} |\mathcal{P}(s' | s, a) \pi(a | s) - \widehat{\mathcal{P}}(s' | s, a) \pi(a | s)| \quad (61)$$

$$= \sum_{a \in \mathcal{A}} \pi(a | s) \sum_{s' \in \mathcal{S}} |\mathcal{P}(s' | s, a) - \widehat{\mathcal{P}}(s' | s, a)| \quad (62)$$

$$= \sum_{a \in \mathcal{A}} \pi(a | s) 2D_{TV} \left(\mathcal{P}(\cdot | s, a), \widehat{\mathcal{P}}(\cdot | s, a) \right) \quad (63)$$

$$\leq \epsilon/H \quad (64)$$

764 Thus with probability at least $1 - \delta/2$ we have for all $s \in \mathcal{S}$ that,

$$D_{TV} \left(\mathcal{P}_\pi(\cdot | s), \widehat{\mathcal{P}}_\pi(\cdot | s) \right) \leq \epsilon/H \quad (65)$$

765 **(2)** Using Assumption 6.3 and 6.4 we can argue about the safety of the system. Suppose firstly,
 766 that we can check the condition $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$, precisely and without any failure
 767 probability (we will deal with statistical model checking in part **(3)**). From any non-critical state we
 768 can transition arbitrarily to a critical state, although under Assumption 6.3 this critical state is not
 769 irrecoverable with probability $\geq p_1$. We now consider the following two cases:

770 (i) $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ under the ‘task’ policy.

771 (ii) $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) > p_1$ under the ‘task’ policy.

772 For case (i) we can safely use the ‘task’ policy and return to a non-critical state within H timesteps
 773 with probability at least $1 - p_1$. For case (ii) we deploy the ‘safe’ policy and under Assumption 6.4
 774 we can return to a non-critical state within H timesteps with probability at least $1 - p_1$. We have now
 775 established an invariant, since from every non-critical state we can return to a non-critical state with
 776 probability $1 - p_1$ and thus satisfy $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ at every timestep $t \in [0, T]$.

777 **(3)** We now make a similar argument but for the statistical model checking procedure where we
 778 can only obtain an ϵ -approximate estimate for the probability $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ with high
 779 probability. Let us denote our ϵ -approximate estimate $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$, rather than check
 780 the condition $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$, we can check condition $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq$
 781 $p_1 - \epsilon$, and if $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept})$ is indeed an ϵ -approximate estimate then this guarantees
 782 $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$. Consider the following two cases:

783 (i) Our estimate $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 - \epsilon$

784 (ii) Our estimate $\widehat{\Pr}(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) > p_1 - \epsilon$

785 For case (i) we can safely use the ‘task’ policy and return to a non-critical state within H timesteps
 786 with probability at least $1 - p_1$. For case (ii) we deploy the ‘safe’ policy and under Assumption 6.4
 787 we can return to a non-critical state within H timesteps with probability at least $1 - p_1$. Again we
 788 have established an invariant, since from every non-critical state we can return to a non-critical state
 789 with probability $1 - p_1$ and thus satisfy $\Pr(\langle s, q \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ at every timestep $t \in [0, T]$.

790 We still need to deal with the failure probability of the statistical model checking procedure at
 791 each timestep, by choosing failure probability $1 - \delta/2T$ we can guarantee (by a union bound) an
 792 ϵ -approximate estimate for each timestep with probability $1 - \delta/2$. Finally, taking a union bound
 793 over part **(1)** and **(2)** gives the desired total failure probability $1 - \delta$.

794

□

795 D Environment Details

796 D.1 Colour Gridworld

797 The colour gridworld environment is a simple 9×9 grid, with
 798 state space $|\mathcal{S}| = 81$ and action space $|\mathcal{A}| = 5$, where each action
 799 corresponds to the following movements: *Left, Right, Up, Down, Stay*.
 800 The objective is to navigate from the start state in one corner of the
 801 grid, to the goal state in the other corner, after reaching the goal state
 802 the agent is then sent back to the start state. The agent must navigate
 803 to the goal state as many times as possible in a fixed episode length
 804 of $T = 1000$. The reward function is a sparse reward that gives the
 805 agent $+1$ reward for reaching the goal and 0 otherwise. When the
 806 environment is fully deterministic the maximum achievable reward
 807 is 58 .

808 In addition to the goal state, there are three other distinct states,
 809 *green, blue* and *purple*, each labelled with their corresponding
 810 colours, see Fig. 4. The set of atomic propositions is thus $AP =$
 811 $\{\text{green, blue, purple, goal}\}$, the safety properties are specified over
 812 the set AP , in particular we conduct experiments with 3 different
 813 safety properties of increasing complexity:

- 814 • (1) $\square \neg \text{green}$
- 815 • (2) $\square \text{goal} \rightarrow \diamond^{\leq 10} \text{blue}$
- 816 • (3) $\square \text{goal} \rightarrow \diamond^{\leq 10} \square^{\leq 5} \text{purple}$

817 Property (1) is a simple invariant property $P_{inv}(\neg \text{green})$ that states the green state must always be
 818 avoided. Property (2) and (3) are more complex safety properties that interfere with the goal state. In

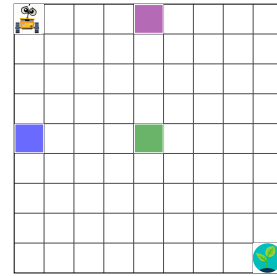


Figure 4: Colour gridworld environment. Top left hand corner (*agent*) is the start position. The agent must navigate to the *goal* position in the bottom right hand corner of gridworld. The coloured states labelled *blue, green* and *purple* correspondingly.

819 particular, property (2) states that once the *goal* state is reached then the *blue* state must be reached
 820 within 10 steps, this actually has no direct consequences on the maximum reward achievable but may
 821 interfere with convergence as the goal state seemingly leads to a high penalty if the *blue* state is not
 822 reached.

823 Property (3) states that once the *goal* state is reached then the *purple* state must be reached within 10
 824 steps and then *purple* must hold for the next 5 timesteps. In safety property both interferes with the
 825 goal and has direct consequences on the maximum achievable reward as staying in purple for 5 steps
 826 does not lead to progress towards the goal state. In terms of the size of the DFA $|\mathcal{Q}|$, property (1) is
 827 an invariant so the cost function is Markov and the size of the DFA is 2, for property (2) and (3) the
 828 size of the DFA is 12 and 62 respectively.

829 Each of the safety properties are tested with the cor-
 830 responding p value for the environment, detailed in
 831 Table 1, which is repeated here for reference. The p
 832 value corresponds to the level of stochasticity in the
 833 environment. In particular, if $p = 0.25$ then there is
 834 a 25% chance of the agents action being overridden
 835 with another random action chosen uniformly. Given
 836 the environment is stochastic then it is difficult to satisfy the safety properties with probability 1.
 837 Through preliminary statistical analysis we computed the maximum satisfaction probabilities for
 838 each property, to help inform an appropriate p value to test with. With $p = 0.25$, property (1) can be
 839 satisfied with very high probability close to 1, while still achieving maximum reward. With $p = 0.25$
 840 property (2) can be satisfied with probability ≈ 0.93 while still achieving maximum reward. With
 841 $p = 0.1$ property (3) can be satisfied with probability ≈ 0.75 while still achieving good reward.

Table 2: Safety properties and p value

property	rand. act. p
(1) $\Box \neg green$	0.25
(2) $\Box goal \rightarrow \Diamond^{\leq 10} blue$	0.25
(3) $\Box goal \rightarrow \Diamond^{\leq 10} \Box^{\leq 5} purple$	0.1

842 **Hyperparameter settings.** We discuss some of the hyperparameter settings for our shielding
 843 approach that are not detailed in Table 5.

844 Property (1): we use a model checking horizon of $H = 3$, and probability threshold $p_1 = 1.0$, with
 845 the number of samples $m = 4096$, we can obtain a roughly $\epsilon = 0.05$ approximate estimate of the
 846 finite horizon satisfaction probability with failure probability $\delta = 0.01$.

847 Property (2): we use a model checking horizon of $H = 10$, and probability threshold $p_1 = 0.9$, with
 848 the number of samples $m = 8192$, we can obtain a roughly $\epsilon = 0.05$ approximate estimate of the
 849 finite horizon satisfaction probability with a smaller failure probability $\delta = 0.001$.

850 Property (3): again we use a model checking horizon of $H = 10$, and probability threshold $p_1 = 0.6$,
 851 with the number of samples $m = 1024$, we can obtain roughly a $\epsilon = 0.1$ approximate estimate of the
 852 finite horizon satisfaction probability with failure probability $\delta = 0.01$.

853 **Extended discussion of results.** First we provide slightly larger figures that than provided in the
 854 main paper, see Figure 5.

855 In general we observe that our shielding method is able to effectively trade-off reward and safety, in
 856 all cases converging to a system that obtains superior or comparable performance with the baseline.
 857 For property (1) we might expect our method to be able to recover the optimal policy that avoids the
 858 green state, it is clear in this case that the shielding procedure has harmed convergence and perhaps
 859 further investigation and hyperparameter tuning will encourage improvements. For property (2) and
 860 (3) the results are what we expect – we can recover the best policy that satisfies the step-wise bounded
 861 safety property with the desired probability p_1 .

862 The intuitive reason for why simply penalising Q-learning doesn’t work, is that tuning the cost
 863 coefficient C is challenging for stochastic environments, where safety cannot be enforced ‘almost
 864 surely’ (with probability 1), and the precise value of C offers little to no semantic meaning. For
 865 different levels of stochasticity p values it is hard to know what desired level of safety we can achieve
 866 while still converging to a high reward policy, making tuning C even harder without knowing more
 867 about the structure of the environment. In Appendix F we study more closely the effect of C and
 868 p . Furthermore, we note the sensitivity of our method to the chosen model checking horizon H . In
 869 particular, if H is too large we might expect the system to be overly conservative, we also address
 870 this in more detail in Appendix F.

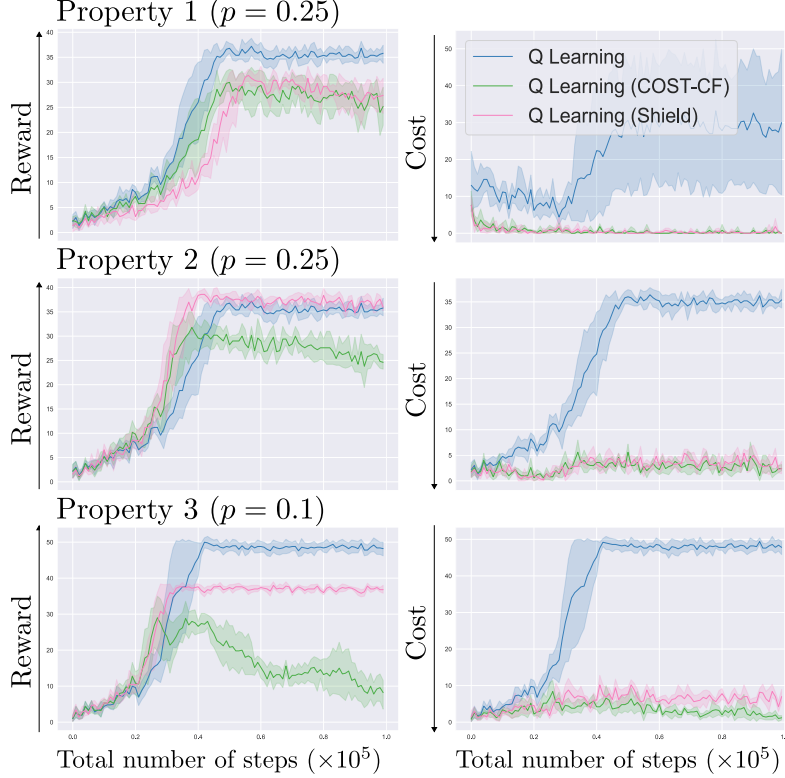


Figure 5: Episode reward and cost for tabular RL ‘colour’ gridworld environment.

871 D.2 Atari Seaquest

872 Our DreamerV3 [34] based shielding procedure is tested
 873 on Atari Seaquest, provided as part of the Arcade Learning
 874 Environment (ALE)[10, 50]. Seaquest is a partially
 875 observable environment meaning we do not have direct
 876 access to the underlying state space \mathcal{S} , we are however
 877 provided with observations $o \in \mathcal{O}$ as pixel images which
 878 correspond to $64 \times 64 \times 3$ tensors. Fortunately Dreamer
 879 V3 is specifically designed to operate in visual settings
 880 and is able to effectively learn a predictive world model
 881 that closely approximate the environment dynamics. The
 882 action space of Seaquest is finite, specifically $|\mathcal{A}| = 18$,
 883 where each action corresponds to a joystick movement and fire
 884 button interaction. Rewards are obtained by ‘shooting’ an
 885 enemy shark or submarine, or by rescuing divers and
 886 returning them to the surface. In addition, the agent must
 887 manage its oxygen resources and avoid being hit by sharks
 and the enemy submarines which fire back, see Fig. 6. The environment is also made stochastic by using ‘sticky actions’ [50], where the agents previous action is repeated with probability $p = 0.25$.



Figure 6: Atari Seaquest environment [10, 50]. The goal is to rescue divers (*small blue people*), while shooting enemy *sharks* and *submarines*.

888 In terms of safety properties we experiment with the following two properties,

- 889 • (1) $(\Box \neg \text{surface} \rightarrow \Box (\text{surface} \rightarrow \text{diver})) \wedge (\Box \neg \text{out-of-oxygen}) \wedge (\Box \neg \text{hit})$
- 890 • (2) $\Box \text{diver} \wedge \neg \text{surface} \rightarrow \Diamond^{\leq 30} \text{surface}$

891 Property (1) states that after diving (i.e. not *surface*), the agent must only *surface* with a *diver* on board, and never run *out-of-oxygen* and never get *hit* by an enemy. The size of the DFA for this property is $|\mathcal{D}| = 4$. Property (2) states that once a *diver* is on board the agent must *surface* within 30 timesteps (i.e. rescue the diver).

895 **Hyperparameter settings.** For our shielding approach almost all the hyperparameters are specified
 896 in Appendix E. The only hyperparameter that varies is the model checking horizon H . For property
 897 (1) we use $H = 30$, empirically this seems adequate enough to avoid running *out-of-oxygen* and
 898 begin surfacing in enough time. For property (2) we use $H = 50$, this is to avoid picking up a *diver*
 899 at the bottom of the ocean where it may not be possible to return to the surface in 30 timesteps.

900 **Extended discussion of results.** First we provide slightly larger figures that than provided in the
 901 main paper, see Figure 7

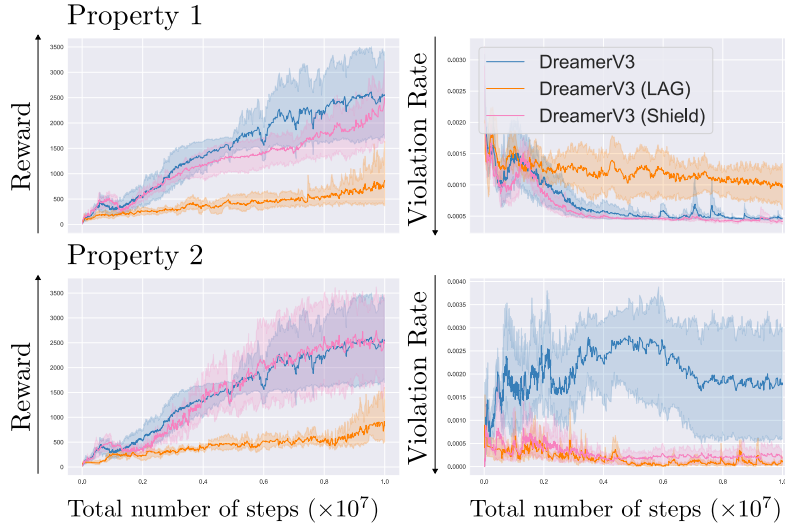


Figure 7: Episode reward and violation rate for deep RL Atari Seaquest.

902 For both safety properties DreamerV3 with shielding obtains comparative performance in terms of
 903 reward with the unmodified DreamerV3 baseline. Of course this baseline entirely ignores the safety
 904 properties and simply maximizes reward. We remark on the differences between the safety properties
 905 themselves, property (1) in particular specifies the natural safety properties of the environment, since
 906 violating property (1) results in a death, the agent only start with 4 lives (and can gain one more ever
 907 10000 points) and so satisfying property (1) is beneficial for long term reward, short the behaviour
 908 satisfying property (1) is correlated with higher reward and we might expect the globally optimal
 909 policy in the environment to never violated property (1). Property (2) specifies that once a diver is
 910 recovered the submarine must return to the surface in 30 timesteps, we would not expect that the
 911 globally optimal policy satisfies this property (2) rather we would expect to converge to a locally
 912 optimal policy satisfying property (2) while still obtaining good reward.

913 With respect to the baseline DreamerV3 (LAG) which has access to the cost function, we see that in
 914 both cases it fails to reliable learn a safe policy that simultaneously maximizes reward. For property
 915 (2) DreamerV3 (LAG) appear to do slightly better in terms of safety, however when qualitatively
 916 inspecting the runs for property (2) we see the DreamerV3 (LAG) agent intentionally get hit by
 917 enemy submarines/sharks to re-spawn on the surface without actually having to navigate there. This
 918 may be a more effective way to satisfy the safety property with high probability but it clearly leads to
 919 worse long term reward.

920 E Hyperparameters & Implementation Details

921 E.1 Access to Code

922 To maintain a high standard of anonymity we provide code for the experiments run on ‘colour’
 923 gridworld as supplementary material, rather than through GitHub. The colour gridworld environment
 924 is implemented with the Gym [14] interface. Tabular Q-learning is implemented with *numpy* in *Python*,
 925 the model checking procedures (both exact and Monte Carlo) are implemented with JAX [12] which
 926 supports vectorized computation on GPU and CPU. The code for the Atari Seaquest experiments

927 are not currently available, although our code base was heavily derived from the code base for
928 *Approximate Model-based Shielding* (AMBS) [30], see <https://github.com/sacktock/AMBS>
929 (MIT License).

930 **Training details.** For collecting both sets of experiments we have access to 2 Nvidia Tesla A30
931 (24GB RAM) GPU and a 24-core/48 thread Intel Xeon CPU each with 32GB RAM. For the ‘colour’
932 gridworld experiments each run can take several minutes up to a day depending on which property is
933 being tested, for example one run for property (3) can take roughly 1.5 days as the product state space
934 is fairly large. For the Atari Seaquest experiments each run can take 8 hours to 1 day depending on
935 the precise configuration of DreamerV3, in general we see a slow down of $\times 2$ when using shielding
936 compared to the unmodified DreamerV3 baseline. Memory requirements may differ depending on
937 the DreamerV3 configuration used, for the *xlarge* DreamerV3 configuration 32GB of GPU memory
938 should suffice.

939 **Statistical significance.** Error bars are provided for each of our experiments. In particular, we report
940 5 random initializations (seeds) for each experiment, the error bars are non-parametric (bootstrap) 95%
941 confidence intervals, provided by `seaborn.lineplot` with default parameters: `errorbar=(‘ci’,`
942 `95)`, `n_boot=1000`. The error bars capture the randomness in the initialization of the DreamerV3
943 world model and policy parameters, the randomness of the environment and any randomness in the
944 batch sampling.

Table 3: Q-learning

Name	Symbol	value
Learning rate	α	0.1
Discount factor	γ	0.95
Exploration type	-	Boltzmann
Temperature	τ	0.05

Table 4: Q-learning with counterfactual experiences [43]

Name	Symbol	value
Learning rate	α	0.1
Discount factor	γ	0.95
Exploration type	-	Boltzmann
Temperature	τ	0.05
Cost coefficient	C	10.0

Table 5: Q-learning with shielding (Algorithm 1)

Name	Symbol	value
Model checking type	-	<i>Monte-Carlo</i>
Approximate model	-	<i>True</i>
Shielding	-	<i>Task</i>
Number of samples	m	varies
Approximation error	ϵ	varies
Failure probability	δ	varies
Model checking horizon	H	varies
Satisfaction prob.	p	varies
Prior	-	<i>uninformative</i>
‘Task policy’ π_{task}		
See Q-learning (Table 3)		
...		
‘Backup policy’ π_{safe}		
See Q-learning with counterfactual experiences (Table 4)		
...		

Table 6: DreamerV3 [34]

Name	Symbol	value
General		
Replay capacity	$ D $	10^6
Batch size	$ B $	16
Batch length	-	64
Number of envs	-	8
Train ratio	-	64
Number of MLP layers	-	5
Number of MLP units	-	1024
Activation	-	LayerNorm + SiLU
World Model		
Configuration size	-	medium
Number of latents	-	32
Classes per latent	-	32
Number of layers	-	3
Number of hidden units	-	640
Number of recurrent units	-	1024
CNN depth	-	48
RSSM loss scales	$\beta_{\text{pred}}, \beta_{\text{dyn}}, \beta_{\text{rep}}$	1.0, 0.5, 0.1
Predictor loss scales	$\beta_o, \beta_r, \beta_c, \beta_\gamma$	1.0, 1.0, 1.0, 1.0
Learning rate	-	10^{-4}
Adam epsilon	ϵ_{adam}	10^{-8}
Gradient clipping	-	1000
Actor Critic		
Imagination horizon	H	15
Discount factor	γ	0.997
TD lambda	λ	0.95
Critic EMA decay	-	0.98
Critic EMA regularizer	-	1
Return norm. scale	S_{reward}	$\text{Per}(R, 95) - \text{Per}(R, 5)$
Return norm. limit	L_{reward}	1
Return norm. decay	-	0.99
Actor entropy scale	η_{actor}	$3 \cdot 10^{-4}$
Learning rate	-	$3 \cdot 10^{-5}$
Adam epsilon	ϵ_{adam}	10^{-5}
Gradient clipping	-	100

Table 7: Augmented Lagrangian [7, 41, 72]

Name	Symbol	value
Augmented Lagrangian		
Penalty multiplier	μ_k	$5 \cdot 10^{-9}$
Initial Lagrange multiplier	λ^k	0.01
Penalty power	σ	10^{-6}
Cost coefficient	C	1.0
Cost threshold	d	1.0
Penalty Critic		
See ‘Actor Critic’ in Table 6		
...		

Table 8: DreamerV3 with Shielding (Algorithm 5)

Name	Symbol	value
Shielding		
Approximation error	ϵ	0.09
Number of samples	m	512
Failure probability	δ	0.01
Look-ahead/shielding horizon	H	varies
Satisfaction prob.	p	0.9
Cost coefficient	C	10
‘Task policy’		
See ‘Actor Critic’ in Table 6		
...		
‘Backup policy’		
See ‘Actor Critic’ in Table 6		
...		

947 **F Ablation Studies**

948 In this section we provide several ablation studies for the ‘colour’ gridworld environment. We test the
 949 most significant hyperparameters and algorithmic components of our method including the baseline
 950 (Q-learning with penalties). In particular we demonstrate the counterfactual experiences is crucial
 951 for learning the safety properties of the environment when the size of the corresponding DFA is non
 952 trivial. We also experiment with using exact model checking – demonstrating that we don’t loose
 953 much by using statistical model checking procedures. Furthermore, we experiment with the cost
 954 coefficient C , the model checking horizon H and the level of stochasticity p .

955 **F.1 Counterfactual experiences**

956 We run our method and the baseline (Q-learning with penalties) without counterfactual experiences
 957 to train the ‘backup policy’ or penalized task policy (baseline).

958 For property (2) and (3) we see a significant
 959 drop in safety performance, since learning to
 960 respect the safety property over the much larger
 961 product state space will require much more expe-
 962 rience and without exploiting the structure of
 963 the DFA (using counterfactual experiences) to
 964 generate synthetic data the task behaviour will
 965 be much more quickly learnt. For property (1),
 966 the invariant property, we observe identical per-
 967 formance as the DFA is trivial (only 2 states),
 968 and so counterfactual experiences is essentially
 969 redundant in this case.

970 **F.2 Exact model checking**

971 We run our method (Shielding) with two differ-
 972 ent configurations: exact model checking with
 973 the ‘approximate’ transition probabilities (learn-
 974 ing from experience) and exact model check-
 975 ing with the ‘true’ transition probabilities. We
 976 compare these two methods to the configura-
 977 tion used in the main paper: Monte Carlo (statisti-
 978 cal) model checking with the learned transition
 979 probabilities.

980 In all cases we see that Shield (MC-Approx)
 981 obtains almost identical performance to Shield
 982 (Exact-True), which demonstrates that we don’t
 983 loose much by statistical model checking with
 984 the learned probabilities, when for example we
 985 don’t have access to the transition probabilities
 986 ahead of time, or the MDP is too large to ex-
 987 act model check. We see some variance with
 988 Shield (Exact-Approx), which can be explained
 989 by sub-optimal convergence in terms of reward,
 990 although note that the safety performance is con-
 991 sistent with the other configurations. Perhaps ex-
 992 act model checking with an inaccurate model of
 993 the transition probabilities restricts exploration
 994 to areas of the state space that are actually safe.

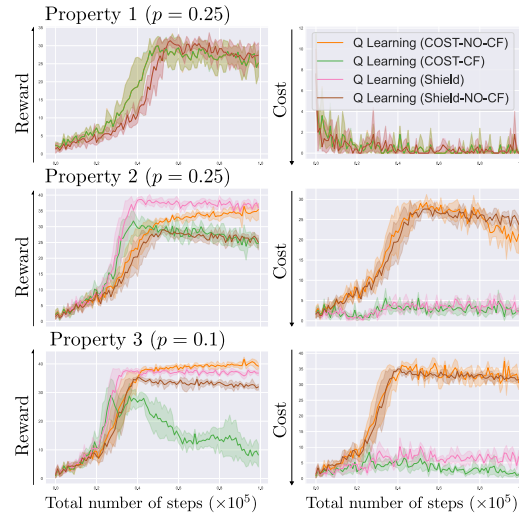


Figure 8: Episode reward and cost for Q-learning (Shield) and Q-learning (COST-CF) with and without counterfactual experiences (CF).

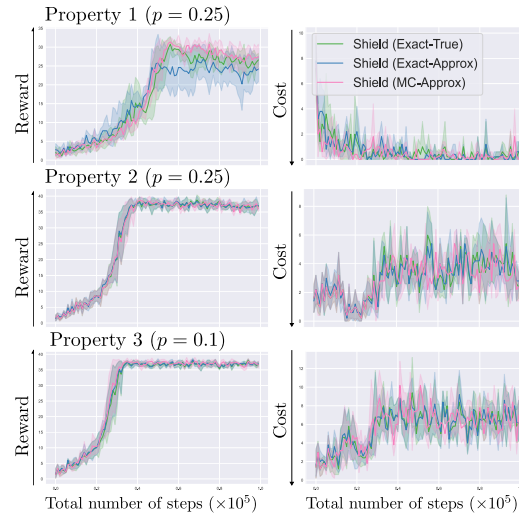


Figure 9: Episode reward and cost for Shield (Exact-True) – exact model checking with the ‘true’ probabilities, Shield (Exact-Approx) - exact model checking with the learning transition probabilities, and Shield (MC-Approx) – from the main paper.

995 **F.3 Cost coefficient C**

996 We experiment with different values for the cost coefficient C used for our baseline (Q-learning with
 997 penalties). In particular, we use $C \in \{0.1, 1.0, 10.0, 100.0\}$, we expect that a larger cost coefficient
 998 will penalize unsafe behaviour more harshly and result in ‘safer’ behaviour (i.e., fewer safety-property
 999 violations).

1000 Unsurprisingly, across the board, by increasing
 1001 the cost coefficient C we obtain a policy that has
 1002 fewer safety-property violations. The improved
 1003 ‘safety performance’ is of course at the expense
 1004 of reward or task performance, this is a trade-off
 1005 we would expect. In particular for $C = 100.0$
 1006 we see that the learned policy essentially avoids
 1007 the goal state (achieving zero reward) all but
 1008 guaranteeing safety (no safety-violations). The
 1009 purpose of this ablation study is to demonstrate
 1010 that while we can achieve any desired level of
 1011 safety by tuning the cost coefficient C , the actual
 1012 value of C offers little to no semantic meaning
 1013 for the probability of violating the safety prop-
 1014 erty.

1015 **F.4 Model checking horizon H**

1016 As was alluded to in the main paper, our method
 1017 can be very sensitive to the model checking hori-
 1018 zon (hyperparameter) H . In particular, if H is
 1019 too large then we might expect the system to
 1020 exhibit overly conservative behaviour. As a rule of thumb we suggest that H should be set to roughly
 1021 the shortest path in the DFA from the initial state to an accepting state – this can easily be computed
 1022 by using Dijkstra’s (shortest-path) algorithm. In this ablation we experiment with much larger H
 1023 than recommended. This significantly impacts the performance of our proposed approach. However,
 1024 we do propose a solution, Q-learning (Shield-Rec) which in short, checks that the action proposed by
 1025 the ‘task policy’ is recoverable with the ‘backup policy’, or in other words by playing with the action
 1026 $a \sim \pi_{task}$ proposed by the ‘task policy’ We can still satisfy $\Pr(\langle s, q \rangle \models \diamond^{\leq H} accept) \leq p_1$ by using
 1027 the ‘backup policy’ after playing a .

1028 In general we observe that when H is too large
 1029 our original method (Shield) is overly conserva-
 1030 tive, sacrificing reward or task performance
 1031 for safety guarantees. Our proposed solution
 1032 (Shield-Rec) is alleviates this issue partly, pro-
 1033 viding reasonable safety performance and compar-
 1034 able task performance. We note that this
 1035 solution is clearly not perfect as it appears
 1036 to be slightly more permissive allowing more
 1037 safety-violations than necessary. More investi-
 1038 gation into this framework would be interesting
 1039 future work, and perhaps more hyperparameter
 1040 tuning, specifically by tuning p_1 , could improve
 1041 this method. The goal would be to obtain an al-
 1042 gorithm that is not overly sensitive to H , and as
 1043 long as H is sufficiently big to guarantee safety
 1044 we don’t see much performance degradation by
 1045 further increasing H .

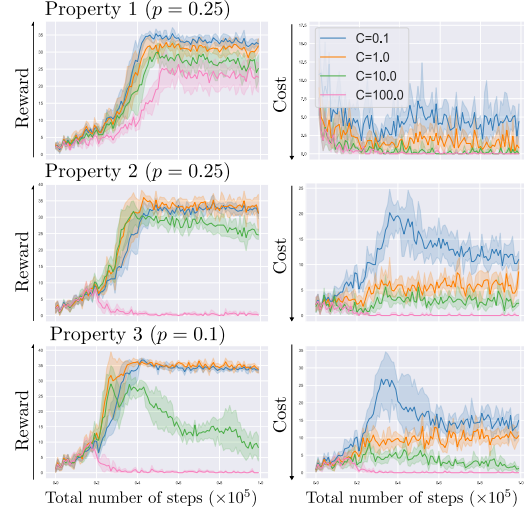


Figure 10: Episode reward and cost for Q-learning (COST-CF) – baseline from the main paper, with different cost coefficients C .

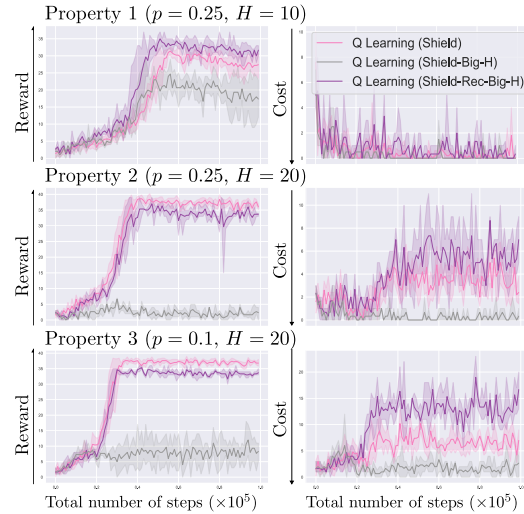


Figure 11: Episode reward and cost for Q-learning (Shield) - from the main paper, Q-learning (Shield) with bigger H and Q-learning (Shield-Rec) with bigger H .

1046 E.5 Level of stochasticity p

1047 Finally we investigate the effect of the level of stochasticity of the environment. Specifically, the value
 1048 p corresponding to the probability that the agent’s action is ignored and another action is chosen
 1049 (uniformly at random) from the action space and played instead. For example, of $p = 0.25$ and the
 1050 agent chooses the action *Right*, there is a 75% chance that the agent goes right and a 25% chance
 1051 the agent goes a different direction. If $p = 0.0$ (deterministic environment) then achieving complete
 1052 safety (zero-violations) becomes easier as the agent has complete control of the environment through
 1053 their actions.

1054 We experiment with the following p values: $p =$
 1055 0.1 for property (1), $p = 0.1$ for property (2)
 1056 and $p = 0.05$ for property (3). For these smaller
 1057 p values we would expect it to be easier for
 1058 our methods including the baseline to achieve
 1059 a higher-rate of safety and possibly complete
 1060 safety in some cases.

1061 We see a similar situation as in the main paper,
 1062 Q-learning (without penalties) simply finds the
 1063 best policy ignoring costs. However, Q-learning
 1064 (with penalties) is able to obtain the same perfor-
 1065 mance now as our method Q-learning (Shield),
 1066 both in terms of reward and cost. With a smaller
 1067 p value the safety-property can be satisfied with
 1068 higher probability while still visiting the goal
 1069 state frequently and obtaining high reward. In
 1070 particular, these p values are chosen such that
 1071 each of the safety properties can be satisfied with
 1072 probability at least 0.9 from the goal state, thus
 1073 penalizing safety-violations with $C = 10.0$ ap-
 1074 pears to be enough to guarantee safety above
 1075 0.9 at each timestep while still achieving high
 1076 reward. For different values of C we might expect the baseline to have a different performance
 1077 profile.

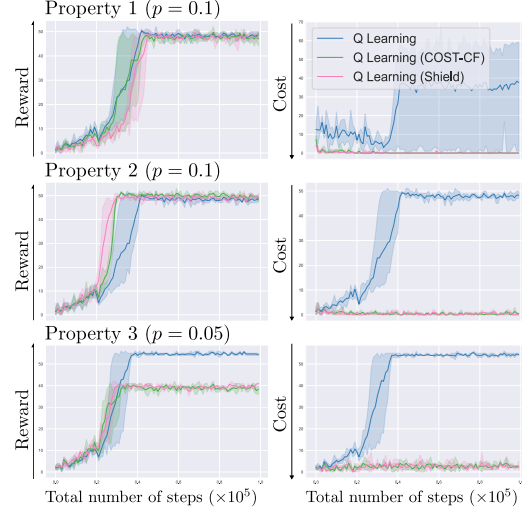


Figure 12: Episode reward and cost for Q-learning, Q-learning (COST-CF) and Q-learning (Shield) – all from the main paper. With smaller levels of stochasticity p

1078 G Comparison to CMDP

1079 In this additional section we analyze the relationships between our problem setup and other common
 1080 CMDP settings, for both the finite horizon and corresponding (discounted) infinite horizon problems.

1081 G.1 Finite Horizon

1082 For reference we restate Problem 4.1 here.

1083 **Problem 4.1 (restated)** (Step-wise bounded regular safety property constraint). Let P_{safe} be a regular
 1084 safety property, \mathcal{D} be the DFA such that $\mathcal{L}(\mathcal{D}) = \text{BadPref}(P_{safe})$ and \mathcal{M} be the MDP;

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T]$$

1085 where all probability is taken under the product Markov Chain $\mathcal{M}_{\pi} \otimes \mathcal{D}$, $p_1 \in [0, 1]$ is a probability
 1086 threshold, H is the model checking horizon and T is the fixed episode length.

1087 G.1.1 Expected Cumulative Constraint

1088 First we restate Problem 4.4.

Problem 4.4 (restated) (Expected cumulative constraint [4, 58]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{E}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} \left[\sum_{t=0}^T C(\langle s_t, q_t \rangle) \right] \leq d_1$$

1089 where $d_1 \in \mathbb{R}_+$ is the cost threshold and T is the fixed episode length.

1090 **Proposition G.1.** *A feasible policy π for Problem 4.1 with parameters $p_1 \in [0, 1]$ is also a feasible*
 1091 *policy for Problem 4.4 with parameter $d_1 \in \mathbb{R}_+$, provided that $d_1 \geq (T + 1) \cdot p_1$.*

1092 *Proof.* For $t \in [0, T]$ we define, the following random variables, X_0, \dots, X_T , where

$$X_t = \mathcal{C}(\langle s_t, q_t \rangle) = 1[\text{accept} \in L'(\langle s_t, q_t \rangle)] \quad (66)$$

1093 where,

$$\mathbb{E}[X_t] = \mathbb{E}[1[\text{accept} \in L'(\langle s_t, q_t \rangle)]] \quad (67)$$

$$= \Pr(\text{accept} \in L'(\langle s_t, q_t \rangle)) \quad (68)$$

$$\leq p_1 \quad (69)$$

1094 The argument is straightforward if at every timestep $t \in [0, T]$ we have $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq$
 1095 p_1 then with probability $\leq p_1$ we have $\text{accept} \in L(\langle s_t, q_t \rangle)$. Then, under mild assumptions
 1096 (i.e. $\mathcal{C}(\langle s_t, q_t \rangle) < \infty$) we consider the following decomposition of the expected cumulative cost,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_\pi \left[\sum_{t=0}^T X_t \right] \quad (70)$$

$$= \mathbb{E}_{s_0 \sim \mathcal{P}_0(\cdot)} [X_0] + \mathbb{E}_{s_1 \sim \mathcal{P}_1(\cdot)} [X_1] + \dots + \mathbb{E}_{s_T \sim \mathcal{P}_T(\cdot)} [X_T] \quad (71)$$

$$= \mathbb{E}_\pi [X_0] + \mathbb{E}_\pi [X_1] + \dots + \mathbb{E}_\pi [X_T] \quad (72)$$

1097 We replace the subscript ' $\langle s_t, q_t \rangle \sim \mathcal{M}_\pi \otimes \mathcal{D}$ ' here for brevity. Clearly by linearity of expectations
 1098 this statement holds. Although it is worth noting that each expectation is taken under a different
 1099 marginal state distribution (i.e. $\mathcal{P}_t(\cdot)$), which depends on π (apart from the initial state distribution
 1100 $\mathcal{P}_0(\cdot)$). From now on we will write this is implicitly (i.e. Eq. 72), rather than writing the marginal
 1101 state distribution (at time t) for each expectation. Using our earlier observations we can now bound
 1102 the expected cumulative cost from above as follows,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_\pi [X_0] + \mathbb{E}_\pi [X_1] + \dots + \mathbb{E}_\pi [X_{T-1}] + \mathbb{E}_\pi [X_T] \quad (73)$$

$$\leq (T + 1) \cdot p_1 \quad (74)$$

1103 \square

1104 **Proposition G.2.** *The converse is not strictly true, since there may be a feasible policy π for Problem*
 1105 *4.4 with threshold $d_1 \leq (T + 1) \cdot p_1$ which does not satisfy the constraints of Problem 4.1.*

1106 *Proof.* We want to prove the following statement, a policy π satisfying,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq (T + 1) \cdot p_1 \quad (75)$$

1107 does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T] \quad (76)$$

1108 To prove this we will show that there may be some policy π that satisfies Eq. 75, but does not satisfy
 1109 Eq. 76 at some timestep t . For simplicity we consider the first timestep (i.e. $t = 0$). First we assume
 1110 π is such that Eq. 75 holds, assuming $H \leq T$ then clearly we have,

$$\mathbb{E}_\pi \left[\sum_{t=0}^H \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq \mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq (T + 1) \cdot p_1 \quad (77)$$

1111 Let $\Pr(\langle s_0, q_0 \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths from the initial state
 1112 $s_0 \sim \mathcal{P}_0(\cdot)$ and automaton state $q_0 = \Delta(\mathcal{Q}_0, L(s_0))$. Suppose π is such that $\Pr(\langle s_0, q_0 \rangle \models$
 1113 $\diamond^{\leq H} \text{accept}) > p_1$. We note that for each path $\rho \in \mathcal{S}^\omega$ and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that
 1114 $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum $\sum_{t=0}^H \mathcal{C}(\langle s_t, q_t \rangle) \geq 1$, and now we have,

$$(T + 1) \cdot p_1 \geq \mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=0}^H \mathcal{C}(\langle s_t, q_t \rangle) \right] > p_1 \quad (78)$$

1115 Now clearly for all $p_1 \in [0, 1]$ and $T \in \mathbb{Z}_+$ the following holds,

$$p_1 < (T + 1) \cdot p_1 \quad (79)$$

1116 This implies that there may exist some π satisfying Eq. 75 and such that $\Pr(\langle s_0, q_0 \rangle \models$
1117 $\diamond^{\leq H} \text{accept}) > p_1$, i.e. does not satisfy Eq. 76 at timestep $t = 0$. \square

1118 **Proposition G.3.** *A feasible policy π for Problem 4.4 with threshold $d_1 \leq p_1$, satisfies $\Pr(\langle s_t, q_t \rangle \models$
1119 $\diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T]$. This bound is tight.*

1120 *Proof.* Firstly, a feasible policy π for Problem 4.4 with threshold $d_1 \leq p_1$ clearly satisfies,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 \quad (80)$$

1121 Assuming $H \leq T$, then this implies that for all $t' \in [0, T - H]$ we have,

$$\mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq \mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 \quad (81)$$

1122 Let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths at timestep t' , where $s_{t'} \sim$
1123 $\mathcal{P}_{t'}(\cdot)$. Here $\mathcal{P}_{t'}(\cdot)$ denotes the marginal state distribution at time t' . Recall that for each path $\rho \in \mathcal{S}^\omega$
1124 and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum $\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1$.
1125 Without loss of generality fix some $t' \in [0, T - H]$ and suppose that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) >$
1126 p_1 . This implies that,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \right] > p_1 \quad (82)$$

1127 Which is a contradiction. Therefore, it must be the case that when Eq. 80 is satisfied then so is
1128 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T - H]$. For the remaining $t' \in [T - H, T]$ a similar
1129 argument can be made, the only detail is to ensure the sum in Eq. 81 is up to T rather than $t' + H$.
1130 To prove that this bound is tight we can again show the possible existence of a counter example. In
1131 particular, we want to prove the following statement, a policy π satisfying,

$$\mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 + c \quad (83)$$

1132 for some constant $c > 0$, does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T] \quad (84)$$

1133 We will show that there may exist some policy π that satisfies Eq. 83 but does not satisfy Eq. 84 at
1134 some timestep t . Firstly, we assume π is such that Eq. 83 holds, this implies that for all $t' \in [0, T - H]$
1135 we have,

$$\mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq \mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 + c \quad (85)$$

1136 Fix some $t' \in [0, T - H]$ and once again let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of
1137 accepting paths at timestep t' . Suppose π is such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. Again recall
1138 that for each path $\rho \in \mathcal{S}^\omega$ and corresponding trace $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$
1139 the sum $\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1$, and so,

$$p_1 + c \geq \mathbb{E}_\pi \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \right] > p_1 \quad (86)$$

1140 Now clearly for all $p_1 \in [0, 1]$ and $c > 0$, the following holds,

$$p_1 < p_1 + c \quad (87)$$

1141 This implies that there may exist some π satisfying Eq. 83 and such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models$
1142 $\diamond^{\leq H} \text{accept}) > p_1$, i.e. does not satisfy Eq. 84 at timestep $t = t'$. \square

1143 **G.1.2 Probabilistic Cumulative Constraint**

1144 First we restate Problem 4.5.

Problem 4.5 (restated) (Probabilistic cumulative constraint [18, 56]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{P}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \leq d_2 \right] \geq 1 - \delta_2$$

1145 where $d_2 \in \mathbb{R}_+$ is the cost threshold, δ_2 is a tolerance parameter and T is the fixed episode length.

1146 **Proposition G.4.** A feasible policy π for Problem 4.1 with parameters $p_1 \in [0, 1]$ is also a
 1147 feasible policy for Problem 4.5 with parameters $d_2 \in \mathbb{R}_+$ and $\delta_2 \in (0, 1]$, provided that,
 1148 $d_2 \geq \sqrt{(T+1)/2 \cdot \log(1/\delta_2)} + (T+1) \cdot p_1$.

1149 *Proof.* For $t \in [0, T]$ we define the following random variables, X_0, \dots, X_T , where,

$$X_t = \mathcal{C}(\langle s_t, q_t \rangle) = 1[\text{accept} \in L'(\langle s_t, q_t \rangle)] \quad (88)$$

1150 and we make the same following observation,

$$\mathbb{E}[X_t] = \mathbb{E}[1[\text{accept} \in L'(\langle s_t, q_t \rangle)]] \quad (89)$$

$$= \Pr(\text{accept} \in L'(\langle s_t, q_t \rangle)) \quad (90)$$

$$\leq p_1 \cdot \delta \quad (91)$$

1151 See the proof of Prop. G.1 for details, the argument is identical. Once again, under mild assumptions
 1152 (i.e. $\mathcal{C}(\langle s_t, q_t \rangle) < \infty$) we consider the following decomposition of the expected cumulative cost,

$$\mathbb{E}_{\pi} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_{\pi}[X_0] + \mathbb{E}_{\pi}[X_1] + \dots + \mathbb{E}_{\pi}[X_T] \quad (92)$$

$$\leq (T+1) \cdot p_1 \quad (93)$$

1153 Again we replace the subscript ' $\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}$ ' here for brevity, see the proof of Prop. G.1 for the
 1154 full details. Before we proceed we must first deal with the dependence between the random variables
 1155 X_0, \dots, X_T . Strictly speaking it is not the case that $\Pr(X_t = 1 \mid X_{t-1}, \dots, X_0) = \Pr(X_t = 1)$.
 1156 However, we have already established that $\Pr(X_t = 1) \leq p_1$, as such we can simulate X_0, \dots, X_T
 1157 as a sequence of independent coin flips Y_0, \dots, Y_T with probability p_1 , it is then the case that
 1158 $\mathbb{P}[\sum_{t=0}^T X_t > d_2] \leq \mathbb{P}[\sum_{t=0}^T Y_t > d_2]$. We can now continue by bounding the probability we care
 1159 about,

$$1 - \mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \leq d_2 \right] = \mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) > d_2 \right] \quad (94)$$

$$= \mathbb{P} \left[\sum_{t=0}^T X_t > d_2 \right] \quad (95)$$

$$\leq \mathbb{P} \left[\sum_{t=0}^T Y_t > d_2 \right] \quad (96)$$

$$= \mathbb{P} \left[\sum_{t=0}^T Y_t > (T+1) \cdot p_1 + d_2 - (T+1) \cdot p_1 \right] \quad (97)$$

$$= \mathbb{P} \left[\sum_{t=0}^T Y_t > \mathbb{E} \left[\sum_{t=0}^T Y_t \right] + d_2 - (T+1) \cdot p_1 \right] \quad (98)$$

$$\leq \exp \left(- \frac{2 \cdot (d_2 - (T+1) \cdot p_1)^2}{\sum_{t=0}^T (\max\{Y_i\} - \min\{Y_i\})^2} \right) \quad (99)$$

$$= \exp \left(- \frac{2 \cdot (d_2 - (T+1) \cdot p_1)^2}{(T+1)} \right) \quad (100)$$

1160 The first inequality (Eq. 96) comes from our earlier construction and the second (Eq. 99) is obtained
 1161 from Hoeffding's inequality [40] for bounded random variables. Finally, bounding the final expression
 1162 from above by δ_2 and rearranging gives the desired result. \square

1163 **Proposition G.5.** A feasible policy π for Problem 4.5 with parameters $\delta_2 \leq p_1$ and $d_2 < 1$, satisfies
 1164 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T]$. This bound is tight.

1165 *Proof.* A feasible policy π for Problem 4.5 with parameters $\delta_2 \leq p_1$ and $d_2 < 1$ clearly implies that,

$$\mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq 1 - p_1 \quad (101)$$

1166 Assuming $H \leq T$, then this implies that for all $t' \in [0, T - H]$ we have,

$$\mathbb{P} \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq \mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq 1 - p_1 \quad (102)$$

1167 Let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths at timestep t' , where $s_{t'} \sim$
 1168 $\mathcal{P}_{t'}(\cdot)$. Again $\mathcal{P}_{t'}(\cdot)$ denotes the marginal state distribution at time t' . Recall that for each path $\rho \in \mathcal{S}^\omega$
 1169 and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum $\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1$.
 1170 Without loss of generality fix some $t' \in [0, T - H]$ and suppose that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) >$
 1171 p_1 . This implies that,

$$\mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \geq 1 \right] \geq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1 \right] > p_1 \quad (103)$$

1172 Which is a contradiction. Therefore, it must be the case that when Eq. 101 is satisfied then so is
 1173 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T - H]$. For the remaining $t' \in [T - H, T]$ a similar
 1174 argument can be made, the only detail is to ensure the sum in Eq. 102 is up to T rather than $t' + H$. To
 1175 prove that this bound is tight we can show the possible existence of a counter example. In particular,
 1176 we want to prove the following statement, a policy π satisfying,

$$\mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq 1 - (p_1 + c) \quad (104)$$

1177 for some constant $c > 0$ does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T] \quad (105)$$

1178 We will show that there may exist some policy π that satisfies Eq. 104 but does not satisfy Eq. 105
 1179 at some timestep t . Firstly, we assume π is such that Eq. 104 holds, this implies that for all
 1180 $t' \in [0, T - H]$ we have,

$$\mathbb{P} \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq \mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) < 1 \right] \geq 1 - (p_1 + c) \quad (106)$$

1181 Fix some $t' \in [0, T - H]$ and let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting
 1182 paths at timestep t' . Suppose that π is such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. Again recall that
 1183 for each path $\rho \in \mathcal{S}^\omega$ and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum
 1184 $\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1$, and so,

$$p_1 + c \geq \mathbb{P} \left[\sum_{t=0}^T \mathcal{C}(\langle s_t, q_t \rangle) \geq 1 \right] \geq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \mathcal{C}(\langle s_t, q_t \rangle) \geq 1 \right] > p_1 \quad (107)$$

1185 Now clearly for all $p_1 \in [0, 1]$ and $c > 0$, the following holds,

$$p_1 < p_1 + c \quad (108)$$

1186 This implies that there may exist some π satisfying Eq. 104 and such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models$
 1187 $\diamond^{\leq H} \text{accept}) > p_1$, i.e. does not satisfy Eq. 105 at timestep $t = t'$. \square

1188 **G.1.3 Instantaneous constraint**

1189 First we restate Problem 4.6.

Problem 4.6 (restated) (Instantaneous constraint [23, 60, 69]).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{P}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} [\mathcal{C}(\langle s_t, q_t \rangle) \leq d_3] = 1 \quad \forall t \in [0, T]$$

1190 **Proposition G.6.** *A feasible policy π for Problem 4.6 with threshold $d_3 < 1$ (otherwise the problem*
 1191 *is trivial) is a feasible policy for Problem 4.1 if and only if $p_1 = 0$.*

1192 *Proof.* We start by proving the 4.6 \Rightarrow 4.1 direction. A feasible policy π for Problem 4.6 with $d_3 < 1$
 1193 satisfies,

$$\Pr(\mathcal{C}(\langle s_t, q_t \rangle) < 1) = 1 \quad \forall t \in [0, T] \quad (109)$$

1194 which implies that,

$$\Pr(\mathcal{C}(\langle s_t, q_t \rangle) = 0) = 1 \quad \forall t \in [0, T] \quad (110)$$

1195 and by Defn. 4.3,

$$\Pr(\text{accept} \notin L'(\langle s_t, q_t \rangle)) = 1 \quad \forall t \in [0, T] \quad (111)$$

1196 Then if for all $t \in [0, T]$, $\text{accept} \notin L'(\langle s_t, q_t \rangle)$ then we have $\Pr(\langle s_0, q_0 \rangle \not\models \diamond \text{accept}) = 1$, where
 1197 $q_0 = \Delta(\mathcal{Q}_0, L(s_0))$ and by extension we have $\Pr(\langle s_t, q_t \rangle \not\models \diamond \text{accept}^{\leq H}) = 1$ for all $t \in [0, T]$.
 1198 This completes the proof of this direction.

1199 Now we prove the 4.1 \Rightarrow 4.6 direction. A policy π satisfying $\Pr(\langle s_t, q_t \rangle \models \diamond \text{accept}^{\leq H}) = 0$ for all
 1200 $t \in [0, T]$ implies that $\Pr(\langle s_t, q_t \rangle \not\models \diamond \text{accept}^{\leq H}) = 1$ for all $t \in [0, T]$ which implies the following,

$$\Pr(\text{accept} \notin L'(\langle s_t, q_t \rangle)) = 1 \quad \forall t \in [0, T] \quad (112)$$

1201 and by Defn. 4.3,

$$\Pr[\mathcal{C}(\langle s_t, q_t \rangle) = 0] = 1 \quad \forall t \in [0, T] \quad (113)$$

1202 which implies that,

$$\Pr[\mathcal{C}(\langle s_t, q_t \rangle) < 1] = 1 \quad \forall t \in [0, T] \quad (114)$$

1203 which concludes the proof. \square

1204 **G.2 Infinite Horizon**

1205 While in this paper we only consider finite horizon problems with a fixed episode length T , we note
 1206 that we can also make a set of similar statements for the infinite horizon (discounted) setting. In this
 1207 section we provide the corresponding statements and proofs for the infinite horizon setting. Firstly,
 1208 we consider the following infinite horizon problem.

1209 **Problem G.7** (Step-wise bounded regular safety property constraint). *Let P_{safe} be a regular safety*
 1210 *property, \mathcal{D} be the DFA such that $\mathcal{L}(\mathcal{D}) = \text{BadPref}(P_{\text{safe}})$ and \mathcal{M} be the MDP;*

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t = 0, 1, 2, \dots$$

1211 *where all probability is taken under the product Markov chain $\mathcal{M}_{\pi} \otimes \mathcal{D}$, $p_1 \in [0, 1]$ is a probability*
 1212 *threshold H is the model checking horizon .*

1213 **G.2.1 Expected Cumulative Constraint**

Problem G.8 (Expected cumulative constraint).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{E}_{\langle s_t, q_t \rangle \sim \mathcal{M}_{\pi} \otimes \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq d_1$$

1214 *where $d_1 \in \mathbb{R}_+$ is the cost threshold and $\gamma \in [0, 1)$ is the discount factor.*

1215 **Proposition G.9.** *A feasible policy π for Problem G.7 with parameters $p_1 \in [0, 1]$, is also a feasible*
 1216 *policy for Problem G.8 with parameter $d_1 \in \mathbb{R}_+$, provided that $d_1 \geq T \cdot p_1$, where $T = 1/(1 - \gamma)$ is*
 1217 *the effective horizon.*

1218 *Proof.* For $t = 0, 1, 2, \dots$ we define, the following random variables, X_0, X_1, X_2, \dots , where,

$$X_t = \mathcal{C}(\langle s_t, q_t \rangle) = 1 [\text{accept} \in L'(\langle s_t, q_t \rangle)] \quad (115)$$

1219 where,

$$\mathbb{E}[X_t] = \mathbb{E}[1 [\text{accept} \in L'(\langle s_t, q_t \rangle)]] \quad (116)$$

$$= \Pr(\text{accept} \in L'(\langle s_t, q_t \rangle)) \quad (117)$$

$$\leq p_1 \quad (118)$$

1220 The argument for this is straightforward. If at every timestep $t = 0, 1, 2, \dots$ we have $\Pr(\langle s_t, q_t \rangle \models$
 1221 $\diamond^{\leq H} \text{accept}) \leq p_1$ then with probability $\leq p_1$ we have $\text{accept} \in L'(\langle s_t, q_t \rangle)$. Let $T = 1/(1 - \gamma)$
 1222 be the effective horizon, then under mild assumptions (i.e. $\mathcal{C}(\langle s_t, q_t \rangle) < \infty$) we can consider the
 1223 following decomposition of the expected cumulative cost,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t X_t \right] \quad (119)$$

$$= \mathbb{E}_{s_0 \sim \mathcal{P}_0(\cdot)} [X_0] + \gamma \cdot \mathbb{E}_{s_1 \sim \mathcal{P}_1(\cdot)} [X_1] + \dots \quad (120)$$

$$+ \gamma^T \cdot \mathbb{E}_{s_T \sim \mathcal{P}_T(\cdot)} [X_T] + \dots$$

$$= \mathbb{E}_\pi [X_0] + \gamma \cdot \mathbb{E}_\pi [X_1] + \dots + \gamma^T \cdot \mathbb{E}_\pi [X_T] + \dots \quad (121)$$

1224 We replace the subscript ' $\langle s_t, q_t \rangle \sim \mathcal{M}_\pi \otimes \mathcal{D}$ ' here for brevity. Clearly by linearity of expectations
 1225 this statement holds. Although it is worth noting that each expectation is taken under a different
 1226 marginal state distribution (i.e. $\mathcal{P}_t(\cdot)$), which depends on π (apart from the initial state distribution
 1227 $\mathcal{P}_0(\cdot)$). From now on we will write this is implicitly (i.e. Eq. 121), rather than writing the marginal
 1228 state distribution (at time t) for each expectation. Using our earlier observations we can now bound
 1229 the expected cumulative cost from above as follows,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_\pi [X_0] + \gamma \cdot \mathbb{E}_\pi [X_1] + \dots + \gamma^T \cdot \mathbb{E}_\pi [X_T] + \dots \quad (122)$$

$$\leq p_1 + \gamma \cdot p_1 + \dots + \gamma^{T-1} \cdot p_1 + \gamma^T \cdot p_1 + \dots \quad (123)$$

$$= p_1 \cdot \sum_{t=0}^{\infty} \gamma^t = p_1 \cdot (1/(1 - \gamma)) = T \cdot p_1 \quad (124)$$

1230

□

1231 **Proposition G.10.** *The converse is not strictly true, since there may be a feasible policy π for*
 1232 *Problem G.8 with threshold $d_1 \leq T \cdot p_1$ which does not satisfy the constraints of Problem G.7*

1233 We want to prove the following statement, a policy π satisfying,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq T \cdot p_1 \quad (125)$$

1234 does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t = 0, 1, 2, \dots \quad (126)$$

1235 *Proof.* To prove this we will show that there may be some policy π that satisfies Eq. 125, but does
 1236 not satisfy Eq. 126 at some timestep t . For simplicity we consider the first timestep (i.e. $t = 0$). First
 1237 we assume π is such that Eq. 125 holds, then clearly we have,

$$\mathbb{E}_\pi \left[\sum_{t=0}^H \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq T \cdot p_1 \quad (127)$$

1238 Let $\Pr(\langle s_0, q_0 \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths from the initial state $s_0 \sim$
 1239 $\mathcal{P}_0(\cdot)$. Suppose π is such that $\Pr(\langle s_0, q_0 \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. We note that for each path $\rho \in \mathcal{S}^\omega$

1240 and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum $\sum_{t=0}^H \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \geq$
 1241 γ^H , and so,

$$T \cdot p_1 \geq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=0}^H \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] > p_1 \cdot \gamma^H \quad (128)$$

1242 Now clearly for all $p_1 \in [0, 1]$, $\gamma \in [0, 1]$, $H \in \mathbb{Z}_+$ and $T = 1/(1 - \gamma)$ the following holds,

$$p_1 \cdot \gamma^H < T \cdot p_1 \quad (129)$$

1243 This implies that there may exist some π satisfying Eq. 125 and such that $\Pr(\langle s_0, q_0 \rangle \models$
 1244 $\diamond^{\leq H} \text{accept}) > p_1$, i.e. does not satisfy Eq. 126 at timestep $t = 0$. \square

1245 **Proposition G.11.** *A feasible policy π for Problem 4.4 with threshold $d_1 \leq p_1 \cdot \gamma^{T+H}$ satisfies*
 1246 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ up to the effective horizon $T = 1/(1 - \gamma)$. *This bound is tight.*

1247 *Proof.* Let $T = 1/(1 - \gamma)$ be the effective horizon. A feasible policy π for Problem 4.4 with threshold
 1248 $d_1 \leq p_1 \cdot \gamma^{T+H}$ clearly satisfies,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 \cdot \gamma^{T+H} \quad (130)$$

1249 which implies that for all $t' \in [0, T]$ we have,

$$p_1 \cdot \gamma^{T+H} \geq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (131)$$

$$= \mathbb{E}_\pi \left[\gamma^{t'} \sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (132)$$

$$= \gamma^{t'} \cdot \mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (133)$$

1250 Let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths at timestep t' , where
 1251 $s_{t'} \sim \mathcal{P}_{t'}(\cdot)$. Here $\mathcal{P}_{t'}(\cdot)$ denotes the marginal state distribution at time t' . Recall that for
 1252 each path $\rho \in \mathcal{S}^\omega$ and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum
 1253 $\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^H$. Without loss of generality fix some $t' \in [0, T]$ and suppose that
 1254 $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. This implies that,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \gamma^{t'} \cdot \mathbb{E}_\pi \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (134)$$

$$> p_1 \cdot \gamma^H \cdot \gamma^{t'} \geq p_1 \cdot \gamma^{T+H} \quad (135)$$

1255 Which is a contradiction. Therefore, it must be the case that when Eq. 130 is satisfied then so is
 1256 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T]$. To prove that this bound is tight we can again
 1257 show the possible existence of a counter example. In particular, we want to prove the following
 1258 statement, a policy π satisfying,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \leq p_1 \cdot \gamma^{T+H} + c \quad (136)$$

1259 for some constant $c > 0$, does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T] \quad (137)$$

1260 We will show that there may exist some policy π that satisfies Eq. 136 but does not satisfy Eq. 137 at
 1261 some timestep t . For simplicity we consider timestep $t = T$, although we note that with a little extra

1262 work we could come up with a proof for any $t \in [0, T]$. Firstly, we assume π is such that Eq. 136
 1263 holds, then we have,

$$p_1 \cdot \gamma^{T+H} + c \geq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \geq \mathbb{E}_\pi \left[\sum_{t=T}^{T+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (138)$$

1264 Let $\Pr(\langle s_T, q_T \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths at timestep T . Suppose π
 1265 is such that $\Pr(\langle s_T, q_T \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. We note that for each path $\rho \in \mathcal{S}^\omega$ and corresponding
 1266 $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum $\sum_{t=T}^{T+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^{T+H}$, and so,

$$p_1 \cdot \gamma^{T+H} + c \geq \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (139)$$

$$\geq \mathbb{E}_\pi \left[\sum_{t=T}^{T+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \right] \quad (140)$$

$$> p_1 \cdot \gamma^{T+H} \quad (141)$$

1267 Now clearly for all $p_1 \in [0, 1]$, $\gamma \in [0, 1]$, $c > 0$, $H \in \mathbb{Z}_+$ and $T = 1/(1 - \gamma)$, the following holds,

$$p_1 \cdot \gamma^{T+H} < p_1 \cdot \gamma^{T+H} + c \quad (142)$$

1268 This implies that there may exist some π satisfying Eq. 136 and such that $\Pr(\langle s_T, q_T \rangle \models$
 1269 $\diamond^{\leq H} \text{accept}) > p_1$, i.e. does not satisfy Eq. 137 at timestep $t = T$. \square

1270 G.3 Probabilistic Cumulative Constraint

Problem G.12 (Probabilistic cumulative constraint).

$$\max_{\pi} V_{\pi} \quad \text{subject to} \quad \mathbb{P}_{\langle s_t, q_t \rangle \sim \mathcal{M}_\pi \otimes \mathcal{D}} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \leq d_2 \right] \geq 1 - \delta_2$$

1271 where $d_2 \in \mathbb{R}_+$ is the cost threshold, δ_2 is a tolerance parameter and $\gamma \in [0, 1)$ is the discount factor.
 1272

1273 **Proposition G.13.** *A feasible policy π for Problem G.7 with parameters $p_1 \in [0, 1]$, is also a*
 1274 *feasible policy for Problem G.12 with parameters $d_2 \in \mathbb{R}_+$ and $\delta_2 \in (0, 1]$, provided that, $d_2 \geq$*
 1275 *$\sqrt{(\lceil \log(T) \rceil \cdot T)/2 \cdot \log(1/\delta_2)} + \lceil \log(T) \rceil \cdot T \cdot p_1 + 1$, where $T = 1/(1 - \gamma)$ is the effective horizon.*
 1276

1277 *Proof.* Again $t = 0, 1, 2, \dots$ we define the following random variables, X_0, X_1, X_2, \dots , where,

$$X_t = \mathcal{C}(\langle s_t, q_t \rangle) = 1 [\text{accept} \in L'(\langle s_t, q_t \rangle)] \quad (143)$$

1278 and we make the following observation,

$$\mathbb{E}[X_t] = \mathbb{E}[1 [\text{accept} \in L'(\langle s_t, q_t \rangle)]] \quad (144)$$

$$= \Pr(\text{accept} \in L'(\langle s_t, q_t \rangle)) \quad (145)$$

$$\leq p_1 \quad (146)$$

1279 See the proof of Prop. G.9, the argument is identical. Under mild assumptions (i.e. $\mathcal{C}(\langle s_t, q_t \rangle) < \infty$)
 1280 we consider the following decomposition of the (undiscounted) expected cumulative cost up to
 1281 timestep $\lceil \log(T) \rceil \cdot T - 1$,

$$\mathbb{E}_\pi \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} \mathcal{C}(\langle s_t, q_t \rangle) \right] = \mathbb{E}_\pi[X_0] + \mathbb{E}_\pi[X_1] + \dots + \mathbb{E}_\pi[X_{\lceil \log(T) \rceil \cdot T - 1}] \quad (147)$$

$$\leq \lceil \log(T) \rceil \cdot T \cdot p_1 \quad (148)$$

1282 Again we replace the subscript ' $\langle s_t, q_t \rangle \sim \mathcal{M}_\pi \otimes \mathcal{D}$ ' here for brevity, see the proof of Prop. G.9 for
 1283 more details. Before we proceed we must first deal with the dependence between the random variables
 1284 $X_0, \dots, X_{\lceil \log(T) \rceil \cdot T - 1}$. Strictly speaking it is not the case that $\Pr(X_t = 1 \mid X_{t-1}, \dots, X_0) =$

1285 $\Pr(X_t = 1)$. However, we have already established that $\Pr(X_t = 1) \leq p_1$, as such we can simulate
1286 $X_0, \dots, X_{\lceil \log(T) \rceil \cdot T - 1}$ as a sequence of independent coin flips $Y_0, \dots, Y_{\lceil \log(T) \rceil \cdot T - 1}$ with probability
1287 p_1 , it is then the case that $\mathbb{P}[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} X_t > d_2] \leq \mathbb{P}[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} Y_t > d_2]$. Now we can
1288 bound the probability that we care about,

$$1 - \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \leq d_2 \right] = \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) > d_2 \right] \quad (149)$$

$$= \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t X_t > d_2 \right] \quad (150)$$

$$= \mathbb{P} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} \gamma^t X_t + \sum_{t=\lceil \log(T) \rceil \cdot T}^{\infty} \gamma^t X_t > d_2 \right] \quad (151)$$

$$\leq \mathbb{P} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} X_t + 1 > d_2 \right] \quad (152)$$

$$\leq \mathbb{P} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} Y_t + 1 > d_2 \right] \quad (153)$$

$$= \mathbb{P} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} Y_t > \lceil \log(T) \rceil \cdot T \cdot p_1 + d_2 - \lceil \log(T) \rceil \cdot T \cdot p_1 - 1 \right] \quad (154)$$

$$= \mathbb{P} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} Y_t > \mathbb{E} \left[\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} Y_t \right] + d_2 - \lceil \log(T) \rceil \cdot T \cdot p_1 - 1 \right] \quad (155)$$

$$\leq \exp \left(- \frac{2 \cdot (d_2 - \lceil \log(T) \rceil \cdot T \cdot p_1 - 1)^2}{\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} (\max\{Y_i\} - \min\{Y_i\})^2} \right) \quad (156)$$

$$= \exp \left(- \frac{2 \cdot (d_2 - \lceil \log(T) \rceil \cdot T \cdot p_1 - 1)^2}{\lceil \log(T) \rceil \cdot T} \right) \quad (157)$$

1289 Here the first inequality (Eq. 152) comes from the following two facts, certainly
1290 $\sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} \gamma^t X_t \leq \sum_{t=0}^{\lceil \log(T) \rceil \cdot T - 1} X_t$ and we have that $\sum_{t=\lceil \log(T) \rceil \cdot T}^{\infty} \gamma^t X_t \leq 1$. The second
1291 fact is a little harder to see, first we note that $\lim_{\gamma \rightarrow 1} \gamma^T = 1/e$, where $T = 1/(1 - \gamma)$ is the
1292 effective horizon. Then we can rewrite,

$$\sum_{t=\lceil \log(T) \rceil \cdot T}^{\infty} \gamma^t X_t = (\gamma^{\lceil \log(T) \rceil \cdot T}) \cdot \left(\sum_{t=\lceil \log(T) \rceil \cdot T}^{\infty} \gamma^{t - \lceil \log(T) \rceil \cdot T} X_t \right) \quad (158)$$

$$= ((\gamma^T)^{\lceil \log(T) \rceil}) \cdot \left(\sum_{t=\lceil \log(T) \rceil \cdot T}^{\infty} \gamma^{t - \lceil \log(T) \rceil \cdot T} X_t \right) \quad (159)$$

$$\leq \left(\frac{1}{e} \right)^{\lceil \log(T) \rceil} \cdot \left(\frac{1}{1 - \gamma} \right) \leq \left(\frac{1}{e} \right)^{\log(T)} \cdot T = \frac{1}{T} \cdot T = 1 \quad (160)$$

1293 The second inequality (Eq. 153) comes from our earlier construction. The final inequality (Eq. 156)
1294 is obtained from Hoeffding's inequality [40] for bounded random variables. Finally, by bounding the
1295 final expression (Eq. 157) from above by δ_2 and rearranging gives the desired result. \square

1296 **Proposition G.14.** *A feasible policy π for Problem G.12 with parameters $\delta_2 \leq p_1$ and $d_2 < \gamma^{T+H}$,
1297 satisfies $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ up to the effective horizon $T = 1/(1 - \gamma)$. This bound is
1298 tight.*

1299 *Proof.* A feasible policy π for Problem G.12 with parameters $\delta_2 \leq p_1$ and $d_2 < \gamma^{T+H}$ clearly
 1300 implies that,

$$\mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \geq 1 - p_1 \quad (161)$$

1301 and certainly for all $t' \in [0, T]$ we have that,

$$1 - p_1 \leq \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \quad (162)$$

$$\leq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \quad (163)$$

$$= \mathbb{P} \left[\gamma^{t'} \sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \quad (164)$$

$$= \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) < (\gamma^{T+H} / \gamma^{t'}) \right] \quad (165)$$

$$\leq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^H \right] \quad (166)$$

1302 Let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths at timestep t' , where
 1303 $s_{t'} \sim \mathcal{P}_{t'}(\cdot)$. Here $\mathcal{P}_{t'}(\cdot)$ denotes the marginal state distribution at time t' . Recall that for
 1304 each path $\rho \in \mathcal{S}^\omega$ and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum
 1305 $\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^H$. Without loss of generality fix some $t' \in [0, T]$ and suppose that
 1306 $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. This implies that,

$$\mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^{T+H} \right] \geq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^H \right] > p_1 \quad (167)$$

1307 Which is a contradiction. Therefore, it must be the case that when Eq. 161 is satisfied then so is
 1308 $\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1$ for all $t \in [0, T]$. To prove that this bound is tight we can show
 1309 the possible existence of a counter example. In particular, we want to prove the following statement,
 1310 a policy π satisfying,

$$\mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \geq 1 - (p_1 + c) \quad (168)$$

1311 for some constant $c > 0$ does not imply that,

$$\Pr(\langle s_t, q_t \rangle \models \diamond^{\leq H} \text{accept}) \leq p_1 \quad \forall t \in [0, T] \quad (169)$$

1312 We will show that there may exist some policy π that satisfies Eq. 168 but does not satisfy Eq. 169 at
 1313 some timestep t . Firstly, we assume π is such that Eq. 168 holds, this implies that for all $t' \in [0, T]$
 1314 we have,

$$1 - (p_1 + c) \leq \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \quad (170)$$

$$\leq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^{T+H} \right] \quad (171)$$

$$\leq \mathbb{P} \left[\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) < \gamma^H \right] \quad (172)$$

1315 Fix some $t' \in [0, T]$ and let $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept})$ denote the proportion of accepting paths
1316 at timestep t' . Suppose that π is such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) > p_1$. Again recall that
1317 for each path $\rho \in \mathcal{S}^\omega$ and corresponding $\text{trace}(\rho) \in \Sigma^\omega$ such that $\text{trace}(\rho) \models \diamond^{\leq H} \text{accept}$ the sum
1318 $\sum_{t=t'}^{t'+H} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^H$, and so,

$$p_1 + c \geq \mathbb{P} \left[\sum_{t=0}^{\infty} \gamma^t \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^{T+H} \right] \quad (173)$$

$$\geq \mathbb{P} \left[\sum_{t=t'}^{t'+b} \gamma^{t-t'} \mathcal{C}(\langle s_t, q_t \rangle) \geq \gamma^H \right] > p_1 \quad (174)$$

1319 Now clearly for all $p_1 \in [0, 1]$, and $c > 0$, the following holds,

$$p_1 < p_1 + c \quad (175)$$

1320 This implies that there may exist some π satisfying Eq. 168 such that $\Pr(\langle s_{t'}, q_{t'} \rangle \models \diamond^{\leq H} \text{accept}) >$
1321 p_1 , i.e. does not satisfy Eq. 169 at timestep $t = t'$. \square

1322 **NeurIPS Paper Checklist**

1323 **1. Claims**

1324 Question: Do the main claims made in the abstract and introduction accurately reflect the
1325 paper's contributions and scope?

1326 Answer: [\[Yes\]](#)

1327 Justification: The proof of safety guarantees (Theorem 6.5) is provided in Appendix C.5.
1328 Furthermore, we provide experimental results demonstrating the scalability of our approach
1329 in Section 7.

1330 Guidelines:

- 1331 • The answer NA means that the abstract and introduction do not include the claims
1332 made in the paper.
- 1333 • The abstract and/or introduction should clearly state the claims made, including the
1334 contributions made in the paper and important assumptions and limitations. A No or
1335 NA answer to this question will not be perceived well by the reviewers.
- 1336 • The claims made should match theoretical and experimental results, and reflect how
1337 much the results can be expected to generalize to other settings.
- 1338 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1339 are not attained by the paper.

1340 **2. Limitations**

1341 Question: Does the paper discuss the limitations of the work performed by the authors?

1342 Answer: [\[Yes\]](#)

1343 Justification: The limitations of our framework are discussed under in Section 7 in the
1344 paragraph **Separating Reward and Safety**, furthermore the limitations of our proposed
1345 method are also discussed and explored in Appendix F.

1346 Guidelines:

- 1347 • The answer NA means that the paper has no limitation while the answer No means that
1348 the paper has limitations, but those are not discussed in the paper.
- 1349 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1350 • The paper should point out any strong assumptions and how robust the results are to
1351 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1352 model well-specification, asymptotic approximations only holding locally). The authors
1353 should reflect on how these assumptions might be violated in practice and what the
1354 implications would be.
- 1355 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1356 only tested on a few datasets or with a few runs. In general, empirical results often
1357 depend on implicit assumptions, which should be articulated.
- 1358 • The authors should reflect on the factors that influence the performance of the approach.
1359 For example, a facial recognition algorithm may perform poorly when image resolution
1360 is low or images are taken in low lighting. Or a speech-to-text system might not be
1361 used reliably to provide closed captions for online lectures because it fails to handle
1362 technical jargon.
- 1363 • The authors should discuss the computational efficiency of the proposed algorithms
1364 and how they scale with dataset size.
- 1365 • If applicable, the authors should discuss possible limitations of their approach to
1366 address problems of privacy and fairness.
- 1367 • While the authors might fear that complete honesty about limitations might be used by
1368 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1369 limitations that aren't acknowledged in the paper. The authors should use their best
1370 judgment and recognize that individual actions in favor of transparency play an impor-
1371 tant role in developing norms that preserve the integrity of the community. Reviewers
1372 will be specifically instructed to not penalize honesty concerning limitations.

1373 **3. Theory Assumptions and Proofs**

1374 Question: For each theoretical result, does the paper provide the full set of assumptions and
1375 a complete (and correct) proof?

1376 Answer: [Yes]

1377 Justification: For the key proofs in the main paper we explicitly provide the assumptions
1378 used, the proofs of each theoretical result from the main paper can also be found in Appendix
1379 C.

1380 Guidelines:

- 1381 • The answer NA means that the paper does not include theoretical results.
- 1382 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
1383 referenced.
- 1384 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1385 • The proofs can either appear in the main paper or the supplemental material, but if
1386 they appear in the supplemental material, the authors are encouraged to provide a short
1387 proof sketch to provide intuition.
- 1388 • Inversely, any informal proof provided in the core of the paper should be complemented
1389 by formal proofs provided in appendix or supplemental material.
- 1390 • Theorems and Lemmas that the proof relies upon should be properly referenced.

1391 4. Experimental Result Reproducibility

1392 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
1393 perimental results of the paper to the extent that it affects the main claims and/or conclusions
1394 of the paper (regardless of whether the code and data are provided or not)?

1395 Answer: [Yes]

1396 Justification: other than access to code, we provide pseudo-code for the algorithms used in
1397 our experiments (see Appendix A)

1398 Guidelines:

- 1399 • The answer NA means that the paper does not include experiments.
- 1400 • If the paper includes experiments, a No answer to this question will not be perceived
1401 well by the reviewers: Making the paper reproducible is important, regardless of
1402 whether the code and data are provided or not.
- 1403 • If the contribution is a dataset and/or model, the authors should describe the steps taken
1404 to make their results reproducible or verifiable.
- 1405 • Depending on the contribution, reproducibility can be accomplished in various ways.
1406 For example, if the contribution is a novel architecture, describing the architecture fully
1407 might suffice, or if the contribution is a specific model and empirical evaluation, it may
1408 be necessary to either make it possible for others to replicate the model with the same
1409 dataset, or provide access to the model. In general, releasing code and data is often
1410 one good way to accomplish this, but reproducibility can also be provided via detailed
1411 instructions for how to replicate the results, access to a hosted model (e.g., in the case
1412 of a large language model), releasing of a model checkpoint, or other means that are
1413 appropriate to the research performed.
- 1414 • While NeurIPS does not require releasing code, the conference does require all submis-
1415 sions to provide some reasonable avenue for reproducibility, which may depend on the
1416 nature of the contribution. For example
 - 1417 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
1418 to reproduce that algorithm.
 - 1419 (b) If the contribution is primarily a new model architecture, the paper should describe
1420 the architecture clearly and fully.
 - 1421 (c) If the contribution is a new model (e.g., a large language model), then there should
1422 either be a way to access this model for reproducing the results or a way to reproduce
1423 the model (e.g., with an open-source dataset or instructions for how to construct
1424 the dataset).
 - 1425 (d) We recognize that reproducibility may be tricky in some cases, in which case
1426 authors are welcome to describe the particular way they provide for reproducibility.
1427 In the case of closed-source models, it may be that access to the model is limited in

1428 some way (e.g., to registered users), but it should be possible for other researchers
1429 to have some path to reproducing or verifying the results.

1430 5. Open access to data and code

1431 Question: Does the paper provide open access to the data and code, with sufficient instruc-
1432 tions to faithfully reproduce the main experimental results, as described in supplemental
1433 material?

1434 Answer: [Yes]

1435 Justification: We provide access to the code for our first set of experiments in the supple-
1436 mentary material, with a corresponding script to reproduce the results in the main paper. For
1437 the second set of experiments we provide directions to the code base that we adapted and
1438 throughout the paper and appendices we provide sufficient details to reproduce these results
1439 without too much difficulty.

1440 Guidelines:

- 1441 • The answer NA means that paper does not include experiments requiring code.
- 1442 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
1443 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1444 • While we encourage the release of code and data, we understand that this might not be
1445 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
1446 including code, unless this is central to the contribution (e.g., for a new open-source
1447 benchmark).
- 1448 • The instructions should contain the exact command and environment needed to run to
1449 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
1450 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1451 • The authors should provide instructions on data access and preparation, including how
1452 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1453 • The authors should provide scripts to reproduce all experimental results for the new
1454 proposed method and baselines. If only a subset of experiments are reproducible, they
1455 should state which ones are omitted from the script and why.
- 1456 • At submission time, to preserve anonymity, the authors should release anonymized
1457 versions (if applicable).
- 1458 • Providing as much information as possible in supplemental material (appended to the
1459 paper) is recommended, but including URLs to data and code is permitted.

1460 6. Experimental Setting/Details

1461 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
1462 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
1463 results?

1464 Answer: [Yes]

1465 Justification: We provide a thorough description of the environmental settings in Appendix
1466 D.1 and D.2, furthermore, hyperparameters and details with regards to access to the code
1467 are provided in Appendix E.

1468 Guidelines:

- 1469 • The answer NA means that the paper does not include experiments.
- 1470 • The experimental setting should be presented in the core of the paper to a level of detail
1471 that is necessary to appreciate the results and make sense of them.
- 1472 • The full details can be provided either with the code, in appendix, or as supplemental
1473 material.

1474 7. Experiment Statistical Significance

1475 Question: Does the paper report error bars suitably and correctly defined or other appropriate
1476 information about the statistical significance of the experiments?

1477 Answer: [Yes] ,

1478 Justification: We provide error bars for all of our experiments, over 5 random initializations
1479 (seeds), provided by `seaborn.lineplot`, see Appendix E for details.

1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide details the compute resources used for our experiments in Appendix E and the expected time and memory requirements for an individual run.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: We have reviewed the NeurIPS Code of Ethics and we do not think the research in this paper has any particular ethical concerns, to the best of our ability we have tried to maintain anonymity.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

1531 Question: Does the paper discuss both potential positive societal impacts and negative
1532 societal impacts of the work performed?

1533 Answer: [NA] .

1534 Justification: As our paper is mostly foundational we do not foresee any immediate positive
1535 or negative societal impact of this research.

1536 Guidelines:

- 1537 • The answer NA means that there is no societal impact of the work performed.
- 1538 • If the authors answer NA or No, they should explain why their work has no societal
1539 impact or why the paper does not address societal impact.
- 1540 • Examples of negative societal impacts include potential malicious or unintended uses
1541 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1542 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1543 groups), privacy considerations, and security considerations.
- 1544 • The conference expects that many papers will be foundational research and not tied
1545 to particular applications, let alone deployments. However, if there is a direct path to
1546 any negative applications, the authors should point it out. For example, it is legitimate
1547 to point out that an improvement in the quality of generative models could be used to
1548 generate deepfakes for disinformation. On the other hand, it is not needed to point out
1549 that a generic algorithm for optimizing neural networks could enable people to train
1550 models that generate Deepfakes faster.
- 1551 • The authors should consider possible harms that could arise when the technology is
1552 being used as intended and functioning correctly, harms that could arise when the
1553 technology is being used as intended but gives incorrect results, and harms following
1554 from (intentional or unintentional) misuse of the technology.
- 1555 • If there are negative societal impacts, the authors could also discuss possible mitigation
1556 strategies (e.g., gated release of models, providing defenses in addition to attacks,
1557 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
1558 feedback over time, improving the efficiency and accessibility of ML).

1559 11. Safeguards

1560 Question: Does the paper describe safeguards that have been put in place for responsible
1561 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1562 image generators, or scraped datasets)?

1563 Answer: [NA] .

1564 Justification: We do not believe that the new assets provided in the paper pose any such
1565 risks.

1566 Guidelines:

- 1567 • The answer NA means that the paper poses no such risks.
- 1568 • Released models that have a high risk for misuse or dual-use should be released with
1569 necessary safeguards to allow for controlled use of the model, for example by requiring
1570 that users adhere to usage guidelines or restrictions to access the model or implementing
1571 safety filters.
- 1572 • Datasets that have been scraped from the Internet could pose safety risks. The authors
1573 should describe how they avoided releasing unsafe images.
- 1574 • We recognize that providing effective safeguards is challenging, and many papers do
1575 not require this, but we encourage authors to take this into account and make a best
1576 faith effort.

1577 12. Licenses for existing assets

1578 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1579 the paper, properly credited and are the license and terms of use explicitly mentioned and
1580 properly respected?

1581 Answer: [Yes]

1582 Justification: We are the original creators/owners of the code used for the first set of
1583 experiments, for the second set of experiments we explicitly cite the paper and provide the
1584 URL for the code that we have adapted in this paper, which is available under the MIT
1585 License as stated in Appendix E.

1586 Guidelines:

- 1587 • The answer NA means that the paper does not use existing assets.
- 1588 • The authors should cite the original paper that produced the code package or dataset.
- 1589 • The authors should state which version of the asset is used and, if possible, include a
1590 URL.
- 1591 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1592 • For scraped data from a particular source (e.g., website), the copyright and terms of
1593 service of that source should be provided.
- 1594 • If assets are released, the license, copyright information, and terms of use in the
1595 package should be provided. For popular datasets, `paperswithcode.com/datasets`
1596 has curated licenses for some datasets. Their licensing guide can help determine the
1597 license of a dataset.
- 1598 • For existing datasets that are re-packaged, both the original license and the license of
1599 the derived asset (if it has changed) should be provided.
- 1600 • If this information is not available online, the authors are encouraged to reach out to
1601 the asset's creators.

1602 13. New Assets

1603 Question: Are new assets introduced in the paper well documented and is the documentation
1604 provided alongside the assets?

1605 Answer: [Yes]

1606 Justification: We provide

1607 Guidelines:

- 1608 • The answer NA means that the paper does not release new assets.
- 1609 • Researchers should communicate the details of the dataset/code/model as part of their
1610 submissions via structured templates. This includes details about training, license,
1611 limitations, etc.
- 1612 • The paper should discuss whether and how consent was obtained from people whose
1613 asset is used.
- 1614 • At submission time, remember to anonymize your assets (if applicable). You can either
1615 create an anonymized URL or include an anonymized zip file.

1616 14. Crowdsourcing and Research with Human Subjects

1617 Question: For crowdsourcing experiments and research with human subjects, does the paper
1618 include the full text of instructions given to participants and screenshots, if applicable, as
1619 well as details about compensation (if any)?

1620 Answer: [NA] .

1621 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1622 Guidelines:

- 1623 • The answer NA means that the paper does not involve crowdsourcing nor research with
1624 human subjects.
- 1625 • Including this information in the supplemental material is fine, but if the main contribu-
1626 tion of the paper involves human subjects, then as much detail as possible should be
1627 included in the main paper.
- 1628 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1629 or other labor should be paid at least the minimum wage in the country of the data
1630 collector.

1631 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1632 Subjects

1633 Question: Does the paper describe potential risks incurred by study participants, whether
1634 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1635 approvals (or an equivalent approval/review based on the requirements of your country or
1636 institution) were obtained?

1637 Answer: [NA] .

1638 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1639 Guidelines:

- 1640 • The answer NA means that the paper does not involve crowdsourcing nor research with
1641 human subjects.
- 1642 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1643 may be required for any human subjects research. If you obtained IRB approval, you
1644 should clearly state this in the paper.
- 1645 • We recognize that the procedures for this may vary significantly between institutions
1646 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1647 guidelines for their institution.
- 1648 • For initial submissions, do not include any information that would break anonymity (if
1649 applicable), such as the institution conducting the review.