

1 Safe Reinforcement Learning 2 using Probabilistic Shields

3 **Nils Jansen**

4 Radboud University, Nijmegen, The Netherlands

5 **Bettina Könighofer**

6 Graz University of Technology, Graz, Austria

7 **Sebastian Junges**

8 University of California, Berkeley, CA, USA

9 **Alex Serban**

10 Radboud University, Nijmegen, The Netherlands

11 **Roderick Bloem**

12 Graz University of Technology, Graz, Austria

13 — Abstract —

14 This paper concerns the efficient construction of a safety shield for reinforcement learning. We
15 specifically target scenarios that incorporate uncertainty and use Markov decision processes (MDPs)
16 as the underlying model to capture such problems. Reinforcement learning (RL) is a machine
17 learning technique that can determine near-optimal policies in MDPs that may be unknown before
18 exploring the model. However, during exploration, RL is prone to induce behavior that is undesirable
19 or not allowed in safety- or mission-critical contexts. We introduce the concept of a probabilistic
20 shield that enables RL decision-making to adhere to safety constraints with high probability. We
21 employ formal verification to efficiently compute the probabilities of critical decisions within a
22 safety-relevant fragment of the MDP. These results help to realize a shield that, when applied to
23 an RL algorithm, restricts the agent from taking unsafe actions, while optimizing the performance
24 objective. We discuss tradeoffs between sufficient progress in the exploration of the environment
25 and ensuring safety. In our experiments, we demonstrate on the arcade game PAC-MAN and on a
26 case study involving service robots that the learning efficiency increases as the learning needs orders
27 of magnitude fewer episodes.

28 **2012 ACM Subject Classification** Computing methodologies → Reinforcement learning; Theory of
29 computation → Verification by model checking; Theory of computation → Reinforcement learning;
30 Computing methodologies → Markov decision processes

31 **Keywords and phrases** Safe Reinforcement Learning, Formal Verification, Safe Exploration, Model
32 Checking, Markov Decision Process

33 **Digital Object Identifier** 10.4230/LIPIcs.CONCUR.2020.3

34 **Category** Invited Paper

35 **Supplementary Material** <http://shieldrl.nilsjansen.org>

36 **1** Introduction

37 Recent years showed increased use of reinforcement learning (RL) in solving tasks such as
38 complex games [60] or robotic manipulation [67]. In RL, an agent perceives the surrounding
39 environment and acts towards maximizing a long-term reward. A major open challenge
40 for systems employing RL is the *safety* of decision-making [61, 30]. In particular during
41 the exploration phase – when an agent chooses random actions in order to examine its
42 surroundings – it is important to avoid actions that may cause unsafe outcomes. The area of



© Nils Jansen, Bettina Könighofer, Sebastian Junges, Alex Serban, Roderick Bloem;
licensed under Creative Commons License CC-BY

31st International Conference on Concurrency Theory (CONCUR 2020).

Editors: Igor Konnov and Laura Kovács; Article No. 3; pp. 3:1–3:16

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 *safe exploration* investigates how RL agents may be forced to adhere to safety requirements
 44 during this phase [54, 4].

45 A suite of methods that deliver theoretical guarantees are so-called *safety-shields* [15].
 46 Shields prevent an agent from taking unsafe actions at runtime. To this end, the performance
 47 objective is extended with a constraint specifying that unsafe states should *never* be visited.
 48 This new safety objective ensures there are no violations during the exploration phase.

49 We propose to incorporate constraints that enforce safety violations to occur *only with*
 50 *small probability*. If an action increases the probability of a safety violation by more than
 51 a threshold δ with respect to the optimal safety probability, the shield blocks the action.
 52 Consequently, an agent augmented with a shield is *guided* to satisfy the safety objective
 53 during exploration (or as long as the shield is used). The shield is *adaptive* with respect
 54 to δ , as a high value for δ yields a stricter shield, a smaller value a more permissive shield.
 55 The value for δ can be changed on-the-fly, and may depend on the individual minimal safety
 56 probabilities at each state. Moreover, in case there is not suitable safe action with respect to
 57 δ , the shield can always pick the optimal action as a fallback.

58 We base our formal notion of a probabilistic shield on MDPs, which constitute a popular
 59 modeling formalism for decision-making under uncertainty [69] and is widely used in model-
 60 based RL. We assess safety by means of probabilistic *temporal logic constraints* [7] that limit,
 61 for example, the probability to reach a set of critical states in the MDP.

62 In order to assess the risk of one action, we (1) construct a behavior model for the
 63 environment using model-based RL [26]. By plugging this model into any concrete scenario,
 64 we obtain an MDP. To construct the shield, we (2) use a model-based verification technique
 65 known as *model checking* [23, 7] that assesses whether a system model satisfies a specification.
 66 In particular, we obtain precise *safety probabilities of any possible decision* within the MDP.
 67 These probabilities can be looked up efficiently and compared to the threshold δ . The shield
 68 then readily (3) augments either model-free or model-based RL.

69 We identify three key challenges. Firstly, model checking – as any model-based technique –
 70 is susceptible to scalability issues. A key advantage of using a separate safety objective is
 71 that we may analyze safety on just a fraction of the system, the *safety-critical MDP*. In
 72 our experiments, these MDP fragments are at least ten orders of magnitude smaller than
 73 a full model of the system, rendering model checking applicable to realistic scenarios. We
 74 introduce further optimizations based on problem-specific abstraction techniques.

75 Secondly, without randomness, all states are either absolutely safe or unsafe. However, in
 76 the presence of randomness, safety may be seen as a quantitative measure: in some states
 77 certain actions induce a large risk, while others may be considered *relatively* safe. Therefore,
 78 we propose an *adaptive* notion of shielding, in which the pre-selection of actions is not based
 79 on absolute thresholds.

80 Lastly, shielding may *restrict* exploration and lead to suboptimal policies. Therefore, it
 81 should not be considered in isolation. The trade-off between optimizing the performance
 82 objective and the achieved safety is intricate. Intuitively, accepting small short term risks
 83 may allow for efficient exploration and limit the risk long-term. To this end, we provide and
 84 discuss mechanisms that allow to adjust the shield based on such observations.

85 We apply shielding to two distinct use cases: the arcade game PAC-MAN and a new case
 86 study involving service robots in a warehouse. Shielded RL leads to improved policies for
 87 both case studies with fewer safety violations and performance superior to unshielded RL.

88 Related Work

89 Safety in RL

90 Safe RL [33, 54] is concerned with providing safety guarantees for learned agents. Our
91 work focusses on the safe exploration [52], we refer to [33] for other types of safe RL. Using
92 their taxonomy, shielding is an instance of “teacher provides advice” [24], where a teacher
93 with additional information about the system guides the RL agent to pick the right actions.
94 Apprenticeship learning [1] is a closely related variant where the teacher gives (positive)
95 examples and has been used in the context of verification [72]. Some of the work does not
96 assume a model for the environment, making the problem intrinsically harder—and often
97 limiting the safety during exploration. We refer to [52, 31, 21, 32, 37, 36, 38] for some
98 interesting approaches.

99 Shielding and non-probabilistic specifications

100 In the remainder of the work, we focus on related work that assumes that the (relevant)
101 dynamics are known. Let us focus on discrete systems first. Exploiting the progress in
102 reactive synthesis [13] provides a shield [44], a maximally permissive policy that contains all
103 actions that will not violate the safety specification. This approach has been shown to be
104 successful in combination with RL [3], and has some noticeable variants. In [68], the temporal
105 specification is extended with an unknown reward function. To enforce complex timing
106 behavior, shields from timed safety properties given as timed automata were considered in [14].
107 Finally, shielding multi-agent systems is considered in [11], and a shield for almost-sure
108 specifications in partially observable MDPs is introduced in [41].

109 Shielding and probabilistic specifications

110 The difference and novel contribution of this paper with the aforementioned papers is rooted
111 in the consideration of stochastic behavior, which is natural to RL. Intuitively, without
112 stochasticities, a learning agent does not take any risk, which is unrealistic in most scenarios.
113 Moreover, often one cannot assume that a 100% (or almost-sure) safety is realizable. A
114 permissive policy for these cases is provided in [29]. Such policies can also be computed
115 from abstract environments [51]. However, as there is no notion of maximally-permissive
116 policies for probabilistic guarantees, multiple permissive policies need to be deployed [42].
117 Furthermore, the computation of these permissive policies is even more expansive than
118 probabilistic model checking itself, harming scalability. A similar approach to ours was
119 developed independently in [16], but targets a different application area and does not consider
120 scalability issues of formal verification.

121 Continuous domains

122 In [71], shielding with qualitative guarantees for continuous domains is discussed. For
123 probabilistic properties and continuous MDPs, additional assumptions are necessary to
124 provide guarantees. Often, these assumption help make some ranking or barrier function [10,
125 53, 22, 34, 2], and have been extended to work with uncertain specifications [64]. UPPAAL
126 STRATEGO provides a number of algorithms combining safety synthesis with optimizing RL
127 for continuous space MDPs [25, 40].

128 **Reinforcement learning in verification**

129 The recent area of programmatic reinforcement learning aims to find (simple) programs
 130 rather than policies represented by deep neural networks, such that these programs can
 131 be verified [8, 66, 65]. These approaches can be seen as an extension to guided policy
 132 synthesis [50]. Such programs have also been used as shields [73]. More generally, Ashok et
 133 al. [6, 5] consider *post-processing* controllers into decision trees. Recently, also the verification
 134 of recurrent neural network controllers has been investigated [18]. These approaches do not
 135 aim to generate permissive shields, but do post-processing to provide guarantees about the
 136 final result. Similarly, methods from reinforcement learning have been successfully employed
 137 to improve the scalability of verification methods for MDPs [17, 57, 45] or other areas of
 138 formal methods [20, 49].

139 **2 Problem Statement**140 **2.1 Foundations**

141 A *probability distribution* over a countable set X is a function $\mu: X \rightarrow [0, 1]$ with $\sum_{x \in X} \mu(x) =$
 142 1. $Distr(X)$ denotes all distributions on X . The support of $\mu \in Distr(X)$ is $supp(\mu) = \{x \in$
 143 $X \mid \mu(x) > 0\}$.

144 **► Definition 1 (MDP).** A Markov decision process (MDP) $\mathcal{M} = (S, Act, \mathcal{P}, r)$ has a set S of
 145 states, a finite set Act of actions, a (partial) probabilistic transition function $\mathcal{P}: S \times Act \rightarrow$
 146 $Distr(S)$, and an immediate reward function $r: S \times Act \rightarrow \mathbb{R}_{\geq 0}$. For all $s \in S$ the available
 147 actions are $Act(s) = \{\alpha \in Act \mid \mathcal{P}(s, \alpha) \neq \perp\}$ and we assume $|Act(s)| \geq 1$.

148 MDPs operate by means of nondeterministic *choices* of actions at each state. A *policy* for
 149 an MDP is a function $\sigma: S^* \rightarrow Distr(Act)$ with $supp(\sigma(s_1 \dots s_n)) \subseteq Act(s_n)$ and S^* a finite
 150 sequence of states.

151 In formal methods, safety properties are often specified as *linear temporal logic* (LTL)
 152 properties [55]. For an MDP \mathcal{M} , probabilistic model checking [43, 47] employs value iteration
 153 or linear programming to compute the probabilities of *all states and actions of the MDP*
 154 to satisfy an LTL property φ . Specifically, we compute $\eta_{\varphi, \mathcal{M}}^{\max}: S \rightarrow [0, 1]$ or $\eta_{\varphi, \mathcal{M}}^{\min}: S \rightarrow [0, 1]$,
 155 which yields for all states the minimal (or maximal) probability over all possible policies
 156 to satisfy φ . For instance, for φ encoding to reach a set of states T , $\eta_{\varphi, \mathcal{M}}^{\max}(s)$ describes the
 157 maximal probability to “eventually” reach a state in T .

158 **2.2 Setting**

159 We define a setting where one controllable agent (the *avatar*) and a number of uncontrollable
 160 agents (the *adversaries*) operate within an *arena*. The arena is a compact, high-level
 161 description of the underlying model. From this arena, the potential states and actions of all
 162 agents may be inferred. For safety considerations, the reward structure can be neglected,
 163 effectively reducing the state space for our model-based safety computations. Formally, an
 164 *arena* is a directed graph $G = (V, E)$ with a finite sets V of nodes and $E \subseteq V \times V$ of edges.
 165 The agent’s *position* is defined via the current node $v \in V$. The agent *decides* on a new edge
 166 $(v, v') \in E$ and determines its next position v' . Some (combinations of) agent positions are
 167 safety-critical, as they e.g., correspond to collisions or falling off a cliff. A safety property
 168 may describe reaching such positions, or use any other property expressible in (the safety
 169 fragment of) temporal logic.

170 While the underlying model for the arena suffices to specify the safe behavior, it is
 171 not sufficiently succinct to model the performance via rewards. Consider an edge that is
 172 safety-relevant, but the agent is only rewarded the first time taking this edge. In a flat model
 173 with rewards, two different edges are necessary to model this behavior. However, the reward
 174 (and thus the difference between these edges) is not needed to assess the safety, and the
 175 safety-relevant model may be pruned to an exponentially smaller model. We use a *token*
 176 function that implicitly extends the underlying model by a reward structure, enabling a
 177 separation of concerns between safety and performance.

178 Technically, we associate edges with a token function $\circ: E \rightarrow \{0, 1\}$, indicating the status
 179 of an edge. Tokens can be (de-) activated and have an associated *reward* earned upon taking
 180 edges with an active token.

181 ► **Example 2 (Autonomous driving).** An autonomous taxi (the avatar) operates within a
 182 road network encoded by an arena. The taxi has to visit several points to pick up or drop
 183 off passengers [28, 35]. Upon visiting such a point, a corresponding token activates and a
 184 reward is earned, afterwards the token is deactivated permanently. Meanwhile, the taxi has
 185 to account for other traffic participants or further environmental factors (the adversaries). A
 186 sensible safety specification may restrict the probability for collision with other cars to 0.5%.
 187 Note that the token structure is not relevant for such a specification.

188 ► **Example 3 (Robot logistics in a smart factory).** Take a factory floor plan with several
 189 corridors with machines. The adversaries are (possibly autonomous) transporters moving
 190 parts within the factory. The avatar models a specific service unit moving around and
 191 inspecting machines where an issue has been raised (as indicated by a token), while accounting
 192 for the behavior of the adversaries. Corridors might be too narrow for multiple (facing) robots,
 193 which poses a safety critical situation. The tokens allow to have a *state-dependent* cost, either
 194 as long as they are present (indicating the costs of a broken machine) or for removing the
 195 tokens (indicating costs for inspecting the machine). A similar scenario has been investigated
 196 in [12].

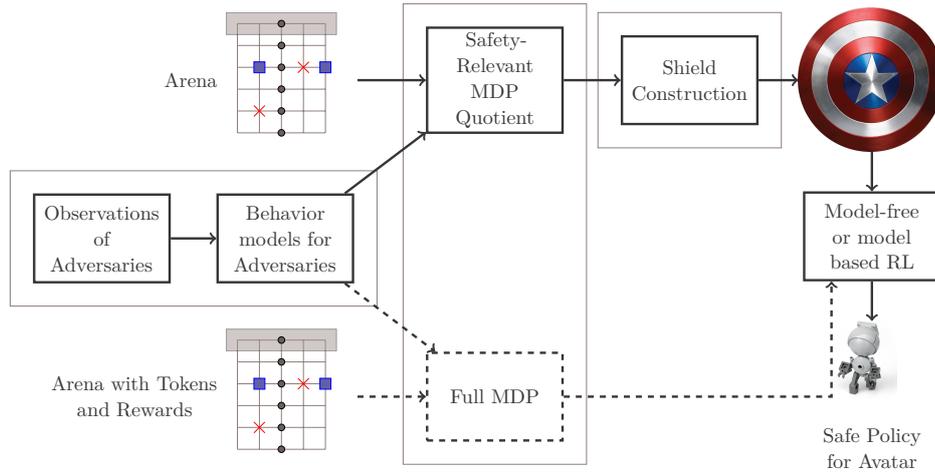
197 2.3 Problem

198 Consider an environment described by an arena as above and a safety specification. We
 199 assume stochastic behaviors for the adversaries, e.g. obtained using RL [58, 59] in a training
 200 environment. In fact, this stochastic behavior determines all actions of the adversaries via
 201 probabilities. The underlying model is then an MDP: the avatar executes an action, and upon
 202 this execution the next exact positions (the state of the system) are determined stochastically.

203 We compute a δ -shield that prevents avatar decisions that violate this specification by
 204 more than a threshold δ with respect to the optimal safety probability. We evaluate the
 205 shield using a model-based or model-free RL avatar that aims to optimize the performance.
 206 The shield therefore has to handle an intricate tradeoff between strictly focussing on (short
 207 and midterm) safety and performance.

208 3 Constructing Shields for MDPs

209 We outline the workflow of our approach in Fig. 1, and describe it below. We employ a
 210 separation of concerns between the model-based shield construction and potentially model-free
 211 reinforcement learning (RL). First, we construct a *behavior model* for each adversary. Based
 212 on this model and a concrete arena, we construct a compact MDP model: the *safety-relevant*



■ **Figure 1** Workflow of the Shield Construction

213 *MDP quotient*. In this MDP, we compute the *shield* which enables safe RL for the full MDP.
 214 We now detail the individual technical steps to realize our proposed method.

215 3.1 Behavior Models for Adversaries

216 We learn an adversary model by observing behavior in a set of similar, but smaller, arenas,
 217 until we gain sufficient confidence that more training data would not change the behavior
 218 significantly [58]. To reduce the size of the training set, we devise a data augmentation
 219 technique using domain knowledge of the arenas [46, 70]. In particular, we abstract away
 220 from the precise configuration of the arena by partitioning the graph into zones that are
 221 relative to the view-point of the adversary (e.g., near or far, north or south, east or west).
 222 The intuitive assumption is that the specific position of an adversary is not important, but
 223 some key information is (e.g., the relation to the position of the avatar). This approach
 224 (1) speeds up the learning process and (2) renders the resulting behavior model applicable
 225 for varying the concrete instance of the same setting.

226 Zones are uniquely identified by a coloring with a finite set C of colors. Formally, for
 227 an arena $G = (V, E)$, *zones relative to a node* $v \in V$ are given by a function $z_v: V \rightarrow C$.
 228 For nodes $x, y \in V$, with $z_v(x) = z_v(y)$, the assumption is that the adversary in v behaves
 229 similarly regardless whether the avatar is in x or y . From our observations, we extract a
 230 *histogram* $h: E \times C \rightarrow \mathbb{N}$, where $h(e, c)$ describes how often the adversary takes an edge
 231 $e = (v, v') \in E$ while the avatar is in a node u with $z_v(u) = c$. We translate these likelihoods
 232 into distributions over possible edges in the arena.

233 ► **Definition 4** (Adversary Behavior). *For an arena $G = (V, E)$, zones $z_u: V \rightarrow C$ for every*
 234 *$u \in V$, and a histogram $h: E \times C \rightarrow \mathbb{N}$, the adversary behavior is a function $B: V \times C \rightarrow$*
 235 *$\text{Distr}(E)$ with*

$$236 \quad B(v, c) = \frac{h((v, v'), c)}{\sum_{(v, v') \in E} h((v, v'), c)} .$$

237 While we employ a simple normalization of likelihoods, alternatively one may also utilize,
 238 e.g., a softmax function which is adjustable to favor more or less likely decisions [62].

3.2 Safety-Relevant Quotient MDP

The construction of the MDP $\mathcal{M} = (S, Act, \mathcal{P})$ augments an arena by behavior models B_i . First, the *states* $S = V^{m+1} \times \{0, \dots, m\}$ encode the positions for all agents and whose turn it is. The *decision states* of the safety-relevant MDP \mathcal{M} are $S_d = \{s_d \in S \mid s_d = (\dots, 0)\}$, i.e., it's the turn of the avatar. The *actions* $Act = \{\alpha_0\} \cup Act_E$ with $Act_E = \{\alpha_e \mid e \in E\}$ determine the movements of the avatar and the adversaries. For $(v, \dots, 0) = s_d \in S_d$ (the avatar moves next), the available actions are $\alpha_e \in Act(s_d) \subseteq Act_e$, where α_e corresponds to an outgoing edge of v . For $(v, \dots, 0) = s_d \in S_d$, α_e with $e = (v, v')$ leads with probability one to a state $s_e = (v', \dots, 1)$. For $(v, \dots, v_i, \dots, i > 0)$ (an adversary moves next), there is a unique action α_0 where v_i is changed to v'_i , randomly determined according to the behavior B_i , which also updates i to $i + 1$ modulo m . These transitions induce the only probabilistic choices in the MDP. A policy only has to choose an action at decision states. At all other states, only the unique action α_0 emanates. Consequently, a policy for \mathcal{M} is a policy for the avatar.

Under the assumption that the reward function is known and not discovered during the exploration of the MDP, one can build the full MDP for the arena (V, E) and the token function $\circ : E \rightarrow \{0, 1\}$. Then, one can compute the reward-optimal and safe policy without need for further learning techniques. As there are 2^E token configurations, the state space blows up exponentially, which prevents the successful application of model checking or planning techniques for anything but very small examples.

3.3 Shield Construction

For the safety-relevant MDP \mathcal{M} , a set of unsafe states $T \subseteq S$ should preferably not be reached from any state. The property $\varphi = \Diamond T$ encodes the **violation** of this safety constraint, that is, eventually reaching T within \mathcal{M} . The shield needs to limit the probability to *satisfy* φ . We evaluate all decision states $s_d \in S_d$ with respect to this probability: We compute $\eta_{\varphi, \mathcal{M}}^{\min}(s_e)$, i.e., the minimal probability to satisfy φ from s_e , which is the state reached after taking action $\alpha_e \in Act_e$ in s_d .

► **Definition 5** (Action-valuation). *An action-valuation for each available action $\alpha_e \in Act_e$ at each decision state $s_d \in S_d$ is given by*

$$val_{s_d}^{\mathcal{M}} : Act(s_d) \rightarrow [0, 1], \text{ with } val_{s_d}^{\mathcal{M}}(\alpha_e) = \eta_{\varphi, \mathcal{M}}^{\min}(s_e) .$$

The optimal action-value for s_d is $optval_{s_d}^{\mathcal{M}} = \min_{\alpha' \in Act} val_{s_d}^{\mathcal{M}}(\alpha')$, the set of all action-valuations at s_d is $ActVals_{s_d}$.

We now define a shield for the safety-relevant MDP \mathcal{M} using the action values. Specifically, a δ -shield for $\delta \in [0, 1]$ determines a set of actions at each decision state s_d that are δ -optimal for the specification φ . All other actions are “shielded” or “blocked”.

► **Definition 6** (Shield). *For action-valuation $val_{s_d}^{\mathcal{M}}$ and $\delta \in [0, 1]$, a δ -shield for state $s_d \in S_d$ is given by*

$$shield_{\delta}^{s_d} : ActVals_{s_d} \rightarrow 2^{Act(s_d)}$$

$$\text{with } shield_{\delta}^{s_d} \mapsto \{\alpha \in Act(s_d) \mid \delta \cdot val_{s_d}^{\mathcal{M}}(\alpha) \leq optval_{s_d}^{\mathcal{M}}\}.$$

Intuitively, δ enforces a constraint on actions that are acceptable with respect to the optimal probability. The shield is *adaptive* with respect to δ , as a high value for δ yields a stricter

282 shield, a smaller value a more permissive shield. The shield is stored using a lookup-table,
 283 and the value for δ can then be changed on-the-fly. In particularly critical situations, the
 284 shield can enforce the decision-maker to resort to (only) the optimal actions w.r.t. the safety
 285 objective. A δ -shield for the MDP \mathcal{M} is built by constructing and applying δ -shields to all
 286 decision states.

287 ► **Definition 7** (Shielded MDP). *The shielded MDP $\mathcal{M}_\square = (S, Act, \mathcal{P}_\square)$ for a safety-relevant
 288 quotient MDP $\mathcal{M} = (S, Act, \mathcal{P})$ and a δ -shield for all $s_d \in S_d$ is given by the transition
 289 probability \mathcal{P}_\square with $\mathcal{P}_\square(s, \alpha) = \mathcal{P}(s, \alpha)$ if $\alpha \in \text{shield}_\delta^s(\text{val}_s^{\mathcal{M}})$ and $\mathcal{P}_\square(s, \alpha) = \perp$ otherwise.*

290 ► **Lemma 8.** *If MDP \mathcal{M} is deadlock-free if and only if the shielded MDP \mathcal{M}_\square is deadlock-free.*

291 We compute the shield relative to optimal values $\text{optval}_{s_d}^{\mathcal{M}}$. Consequently, for $\delta = 1$, only
 292 optimal actions are preserved, and for $\delta = 0$ no actions are blocked.

293 ► **Theorem 9.** *For an MDP \mathcal{M} and a δ -shield, it holds for any state s that $\text{val}_s^{\mathcal{M}} = \text{val}_s^{\mathcal{M}_\square}$.*

294 As optimal actions for the safety objective are not removed, optimality w.r.t. safety is
 295 preserved in the shielded MDP. Thus, during construction of the shield, we compute the
 296 action-valuations in fact *for the shielded MDP*. Observe that computing a shield for a state
 297 is done *independently* from the application of the shield to other states.

298 3.4 Guaranteed Safety

299 A δ -shield ensures that only actions that are δ -optimal with respect to an LTL property φ are
 300 allowed. In particular, for each action $\alpha \in Act_e$ at state s_e , we use the *minimal* probability
 301 $\eta_{\varphi, \mathcal{M}}^{\min}(s_e)$ to satisfy φ , see Def. 5. Under *optimal* (subsequent) choices, the value $\eta_{\varphi, \mathcal{M}}^{\min}(s_e)$
 302 will be achieved. In contrast, a sequence of bad choices may violate φ with high probability.
 303 A more conservative notion would be to use the minimal action value while assuming that in
 304 all subsequent states the worst-case decisions corresponding to the maximal probabilities are
 305 taken. These values are computable by model checking. Regardless of subsequent choices, at
 306 least $\text{val}_{s_d}^{\mathcal{M}}(\alpha_e)$ is then guaranteed. A sensible notion to construct a shield would then be to
 307 impose a threshold $\lambda \in [0, 1]$ such that only actions with $\text{val}_{s_d}^{\mathcal{M}}(\alpha_e) \leq \lambda$ are allowed. A shield
 308 with such a guaranteed safety probability may induce a shielded MDP (Def. 7) that is *not*
 309 *deadlock free*. Moreover, the shield may become too restrictive for the agent.

310 3.5 Scalable Shield Construction

311 Although we apply model checking only in the safety-relevant MDP, scalability issues for large
 312 applications remain. We employ several optimizations towards computational tractability.

313 Finite Horizon

314 For infinite horizon properties, the probability to violate safety (in the long run) is often one
 315 in our examples. Furthermore, our learned MDP model is inherently an approximation of
 316 the real world. Errors originating from this approximation accumulate for growing horizons.
 317 Thus, we focus on a finite horizon such that the action values (and consequently, a policy
 318 for the avatar) carry only guarantees for the next steps. This assumption also allows us to
 319 prune the safety-relevant MDP (see below), increasing the scalability.

320 Piecewise Construction

321 Computing a shield for all states in an MDP concurrently yields a large memory footprint.
 322 To alleviate this footprint, we compute the shield states independently, in accordance with
 323 Theorem 9. The independent computation prunes the relevant part of the MDP, as the
 324 number of states reachable within the horizon is drastically reduced. Additionally, the
 325 independent computation allows for parallelizing the computation.

326 Independent Agents

327 The explosion of state spaces stems mostly from the number of agents. Here, an important
 328 observation is that we can consider agents independently. For instance, the probability for the
 329 avatar to crash with an adversary is stochastically independent from crashing with the others.
 330 Instead of determining the shield for all adversaries at once, we perform computations for
 331 each agent individually, and combine them via the inclusion-exclusion principle. Afterwards,
 332 the shield is composed from the shields dedicated to individual adversaries.

333 Abstractions

334 We observe that for finite horizon properties and piecewise construction, adversaries may
 335 be far away—beyond the horizon—without a chance to reach the avatar. We do not need
 336 to consider such (positions for) adversaries, as in these states, the shield will not block any
 337 actions.

338 **Fewer Decision States.** Depending on the setting, there might be only a few critical
 339 situations in which the agent requires shielding to ensure safety. The shield can be computed
 340 for this critical states only. Consequently, the agent makes shielded decisions in the adapted
 341 decision states, and unshielded decisions in all other ones.

342 **Shielding versus Performance.** A shield which is *minimally invasive* gives the RL
 343 agent the most freedom to optimize the performance objective. We propose two methods
 344 to alleviate invasiveness, all of them assume *domain knowledge* of the rationale behind the
 345 decision procedure.

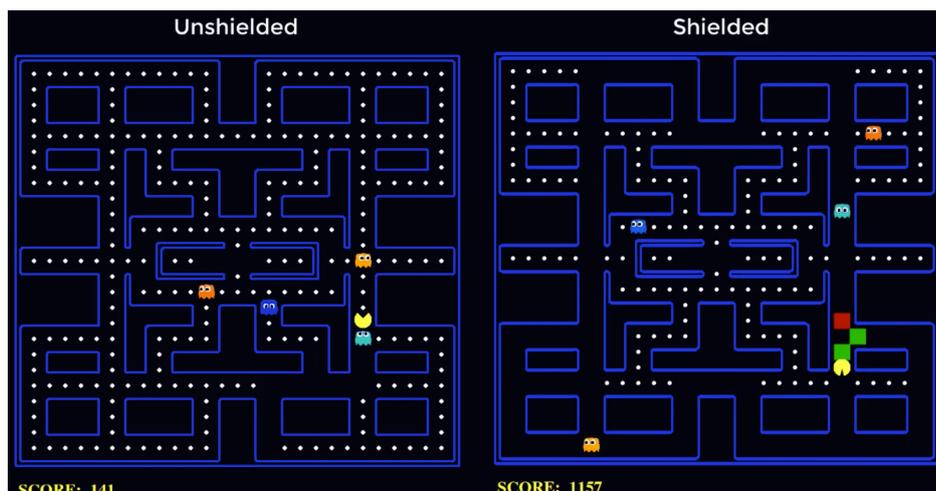
346 **Iterative Weakening.** During runtime, we may observe that the progress of the avatar
 347 regarding the performance objective is not increasing anymore. Then, we weaken the shield by
 348 $\delta - \varepsilon$, allowing additional actions. As soon as progress is made, we reset δ to its former value.
 349 The adaption of $shield_{\delta}^s$ to $shield_{\delta-\varepsilon}^s$ can be done on the fly, without new computations.

350 **Adapted Specifications.** If the goal of the decision maker is known *and* can be
 351 captured in temporal logic, we may adapt the original specification accordingly. There are
 352 often natural trade-offs between safety and performance. These trade-offs might be resolved
 353 via weights, but this process is often undesirable [56] and similar to reward engineering.
 354 Instead, optimizing the conditional performance while assuming to stay sufficiently safe [63],
 355 avoids side-effects of attaching some weights to the safety specification.

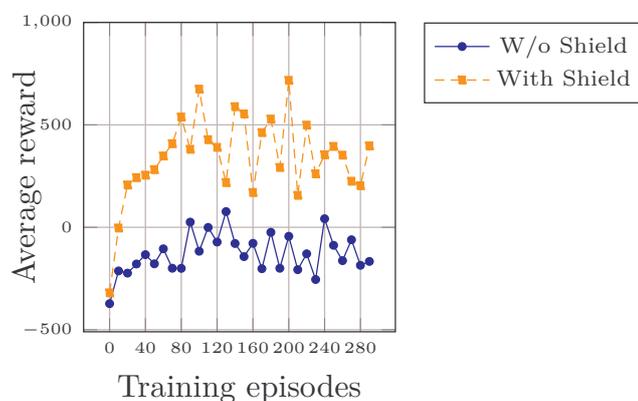
356 4 Implementation and Numerical Experiments

357 4.1 Set-up

358 We run experiments using an Intel Core i7-4790K CPU with 16 GB of RAM using 4 cores. We
 359 give the timing results for a single CPU. Since the shield may be computed in a multi-threaded
 360 architecture, this time can be divided by the number of cores available.



■ **Figure 2** Video still for classic PAC-MAN



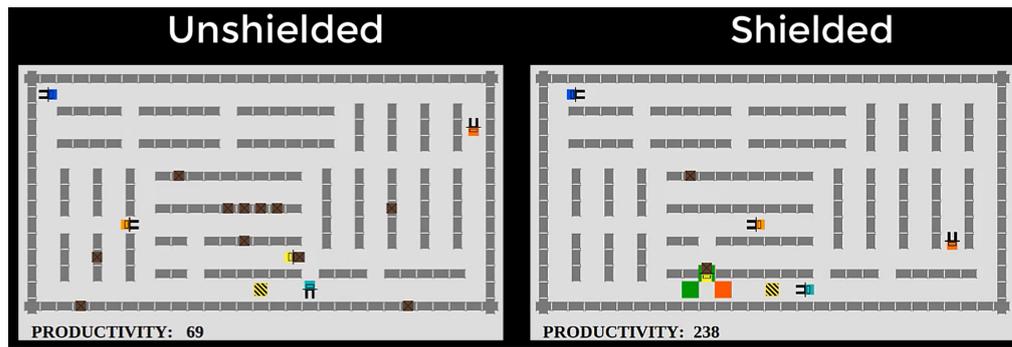
■ **Figure 3** Training scores classic PAC-MAN

361 The supplementary materials, namely the source code and videos are available online¹.

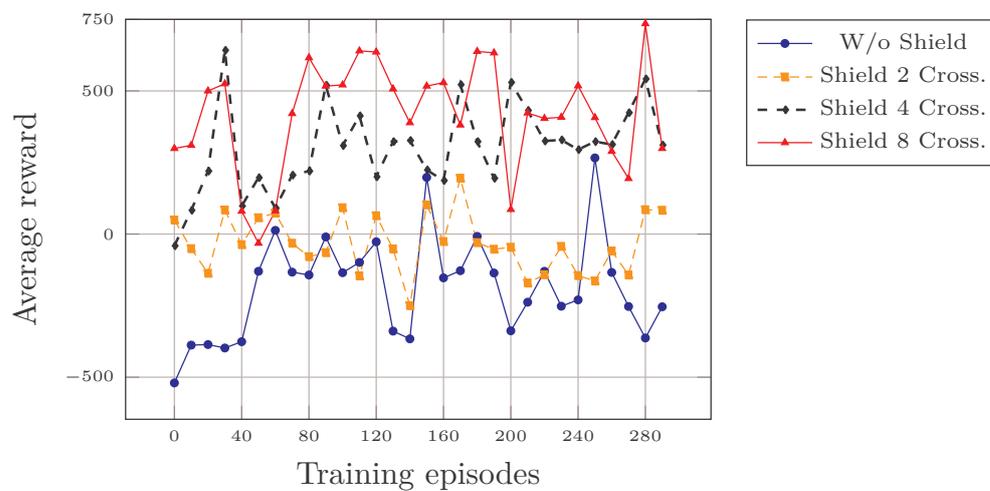
362 We demonstrate the applicability of our approach by means of two case studies. For both
 363 case studies, we learn the adversary behavior in small arenas, individually for each adversary.
 364 These behavior models are applicable to any benchmark instance, as they are independent of
 365 concrete positions.

366 For the arcade game PAC-MAN, PM (the avatar) aims to collect *PAC-dots* in a *maze*
 367 and not get caught by *ghosts* (the adversaries). We model various instance of the game (with
 368 different sizes) as an arena, where tokens represent the dots at each position in the maze,
 369 such that a dot is either present or collected. The score (reward, performance) is positively
 370 affected (+10) by collecting a dot and negatively by time (each step: -1). If PM either
 371 collects all dots (+500) or is caught (-500), the game is restarted. RL approaches exist [9],
 372 but they suffer from the fact that during the exploration phase PM is often caught by the
 373 ghosts, achieving very poor scores. The safety specification places a lower bound on the
 374 probability of reaching states in the underlying MDP that correspond to being caught.

¹ <http://shieldrl.nilsjansen.org>



■ **Figure 4** Video still for warehouse



■ **Figure 5** Training scores warehouse

375 We also consider a warehouse floor plan with several corridors. A similar scenario has
 376 been investigated in [12]. In the arena, nodes describe crossings, the edges the corridors
 377 with shelves, and the distances the corridor length. The agents are fork-lift units picking up
 378 packages from the shelves and delivering them to the exit; tokens represent the presence of a
 379 package at its position. The avatar is a specific (yellow) fork-lift unit that has to account for
 380 other units, the adversaries. The performance (reward) is positively affected by loading and
 381 delivering packages (+20, respectively) and negatively by time (each step: -1). Delivering all
 382 packages yields a large bonus (+500) and a collision leads to a large punishment (-500), both
 383 cases end the scenario. Corridors might be too narrow for multiple (facing) units, which
 384 poses a safety-critical situation. Most crucial is the crowded area near the exit, since all
 385 units have to deliver the packages to the exit.

386 Transferring the stochastic adversary behavior to any arena (without tokens) yields a
 387 concrete safety-relevant MDP. In particular, we specify an arena with the positions of the
 388 avatar and the adversaries as well as their behavior in the high-level PRISM-language [48].
 389 We employ a script that automatically generates arenas to enable a broad set of benchmarks.
 390 Taking, e.g., the PAC-MAN arena from Fig. 2, the considered MDP has roughly 10^{12} states
 391 (compared to 10^{50} for the full MDP). For a safety-relevant MDP, we compute a δ -shield
 392 (with iterative weakening) via the model checker Storm [27], using a horizon of 10 steps. The

immense size even of safety-relevant MDPs requires optimizations such as a piecewise and independent shield construction. Moreover, a multi-threaded architecture lets us construct shields for very large examples. In particular, we perform model checking for (many) MDPs of roughly 10^6 states. The computation time for the largest PM instance takes about 6 hours (single-threaded), while memory is not an issue due to the piecewise shield construction.

We compare RL to shielded RL on different instances. The key comparison criterion is the performance (detailed above) during learning. Our implementation is based on an existing PAC-MAN environment² using an approximate Q-learning agent [62] with the following feature vectors:

- for PAC-MAN: (1) distance to the closest dot, (2) whether a ghost collision is imminent, and (3) whether a ghost is one step away.
- for Warehouse: (1) has the unit loaded or unloaded, (2) the distance to the next package and (3) to the exit, (4) whether another unit is three steps away and (5) one step away.

The results are basic reflex controllers. The Q-learning uses the learning rate $\alpha = 0.2$ and the discount factor $\gamma = 0.8$ for the Q-update and an ϵ -greedy exploration policy with $\epsilon = 0.05$. One episode lasts until either the game is restarted. We describe results for the training phase of RL (300 episodes).

4.2 Results

Figures 2 and 4 show screenshots of a series of recommended videos (available in the supplementary material). Each video compares how RL performs either shielded or unshielded on a instance of the case study. In the shielded version, at each decision state in the underlying MDP, we indicate the risk of decisions from low to high by the colors green, orange, red.

Consider PAC-MAN in detail: Figure 3 depicts the scores obtained during RL. The curves (blue, solid: unshielded, orange, dashed: shielded) show the average scores for every ten training episodes. Table 1 shows results for instances in increasing size. We list the number of model checking calls and the time to construct the shield. We list the scores with and without shield, and the *winning rate* capturing the ratio of successfully ended episodes. For all instances, we see a large difference in scores due to the fact that PM is often rescued by the shield. The winning rates differ for most benchmarks, favoring shielded RL. For three or four ghosts, a shield with a ten-step horizon cannot guide PM to avoid being encircled by the ghosts long enough to successfully end the game. Nevertheless, the shield often safes PM, leading to superior scores. Moreover, the shield helps learning an optimal policy much faster as fewer restarts are needed.

For the warehouse case study, we choose to vary the decision states, i.e., the positions of the avatar for which we compute a shield. We present results for shielding the 2–8 crossings closest to the exit. Figure 5 shows the average score for the different variants, Table 2 summarizes average score and win rate. Unsurprisingly, the score gets better the more states are shielded. Furthermore, we have seen that shielding even more states has only a very limited effect.

5 Conclusion and Future Work

We developed the concept of shields for MDPs. Utilizing probabilistic model checking, we maintained probabilistic safety measures during reinforcement learning. We addressed

² http://ai.berkeley.edu/project_overview.html

■ **Table 1** Average scores and win rates for PM

Size, #Ghosts	#Model Checking	time (s)	Score w/o Shield	Score w. Shield	Win Rate w/o Shield	Win Rate w. Shield
9x7,1	5912	584	-359,6	535,3	0,04	0,84
17x6,2	5841	1072	-195,6	253,9	0,04	0,4
17x10,3	51732	3681	-220,79	-40,52	0,01	0,07
27x25,4	269426	19941	-129,25	339,89	0,00	0,00

■ **Table 2** Average scores and win rates for warehouse

Crossings shielded	0	2	4	8
Score	-186	-27.6	303	420
Win Rate	0.16	0.31	0.59	0.71

435 inherent scalability issues and provided means to deal with typical trade-off between safety
 436 and performance. Our experiments showed that we improved the state-of-the-art in safe
 437 reinforcement learning.

438 For future work, we will extend the applications to more arcade games and employ deep
 439 recurrent neural networks as means of decision-making [39, 19, 18]. Another interesting
 440 direction is to explore (possibly model-free) learning of shields, instead of employing model-
 441 based model checking.

442 ——— References ———

- 443 1 Pieter Abbeel and Andrew Y. Ng. Exploration and apprenticeship learning in reinforcement
 444 learning. In *ICML*, volume 119 of *ACM International Conference Proceeding Series*, pages
 445 1–8. ACM, 2005.
- 446 2 Anayo K. Akametalu, Shahab Kaynama, Jaime F. Fisac, Melanie Nicole Zeilinger, Jeremy H.
 447 Gillula, and Claire J. Tomlin. Reachability-based safe learning with gaussian processes. In
 448 *CDC*, pages 1424–1431. IEEE, 2014.
- 449 3 Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum,
 450 and Ufuk Topcu. Safe reinforcement learning via shielding. In *AAAI*, pages 2669–2678. AAAI
 451 Press, 2018.
- 452 4 Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané.
 453 Concrete problems in AI safety. *CoRR*, abs/1606.06565, 2016.
- 454 5 Pranav Ashok, Mathias Jackermeier, Pushpak Jagtap, Jan Kretínský, Maximilian Weininger,
 455 and Majid Zamani. dtcontrol: decision tree learning algorithms for controller representation.
 456 In *HSCC*, pages 30:1–30:2. ACM, 2020.
- 457 6 Pranav Ashok, Jan Kretínský, Kim Guldstrand Larsen, Adrien Le Coënt, Jakob Haahr
 458 Taankvist, and Maximilian Weininger. SOS: safe, optimal and small strategies for hybrid
 459 markov decision processes. In *QEST*, volume 11785 of *Lecture Notes in Computer Science*,
 460 pages 147–164. Springer, 2019.
- 461 7 Christel Baier and J.P. Katoen. *Principles of Model Checking*. MIT Press, 2008.
- 462 8 Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via
 463 policy extraction. In *NeurIPS*, pages 2499–2509, 2018.
- 464 9 UC Berkeley. Intro to AI – Reinforcement Learning , 2018.
 465 <http://ai.berkeley.edu/reinforcement.html>.
- 466 10 Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. Safe model-based
 467 reinforcement learning with stability guarantees. In *NIPS*, pages 908–919, 2017.

- 468 11 Suda Bharadwaj, Roderick Bloem, Rayna Dimitrova, Bettina Könighofer, and Ufuk Topcu.
469 Synthesis of minimum-cost shields for multi-agent systems. In *ACC*, pages 1048–1055. IEEE,
470 2019.
- 471 12 Arthur Bit-Monnot, Francesco Leofante, Luca Pulina, Erika Ábrahám, and Armando Tacchella.
472 Smartplan: a task planner for smart factories. *CoRR*, abs/1806.07135, 2018.
- 473 13 Roderick Bloem, Krishnendu Chatterjee, and Barbara Jobstmann. Graph games and reactive
474 synthesis. In *Handbook of Model Checking*, pages 921–962. Springer, 2018.
- 475 14 Roderick Bloem, Peter Gjøøl Jensen, Bettina Könighofer, Kim Guldstrand Larsen, Florian
476 Lorber, and Alexander Palmisano. It’s time to play safe: Shield synthesis for timed systems.
477 *CoRR*, abs/2006.16688, 2020.
- 478 15 Roderick Bloem, Bettina Könighofer, Robert Könighofer, and Chao Wang. Shield synthesis: -
479 runtime enforcement for reactive systems. In *TACAS*, volume 9035 of *LNCS*, pages 533–548.
480 Springer, 2015.
- 481 16 Maxime Bouton, Jesper Karlsson, Alireza Nakhaei, Kikuo Fujimura, Mykel J. Kochenderfer,
482 and Jana Tumova. Reinforcement learning with probabilistic guarantees for autonomous
483 driving. *CoRR*, abs/1904.07189, 2019.
- 484 17 Tomáš Brázdil, Krishnendu Chatterjee, Martin Chmelík, Vojtěch Forejt, Jan Křetínský,
485 Marta Z. Kwiatkowska, David Parker, and Mateusz Ujma. Verification of Markov decision
486 processes using learning algorithms. In *ATVA*, 2014.
- 487 18 Steven Carr, Nils Jansen, and Ufuk Topcu. Verifiable rnn-based policies for pomdps under
488 temporal logic constraints. *CoRR*, abs/2002.05615, 2020.
- 489 19 Steven Carr, Nils Jansen, Ralf Wimmer, Alexandru Constantin Serban, Bernd Becker, and
490 Ufuk Topcu. Counterexample-guided strategy improvement for pomdps using recurrent neural
491 networks. In *IJCAI*, pages 5532–5539. ijcai.org, 2019.
- 492 20 Jia Chen, Jiayi Wei, Yu Feng, Osbert Bastani, and Isil Dillig. Relational verification using
493 reinforcement learning. *Proc. ACM Program. Lang.*, 3(OOPSLA):141:1–141:30, 2019.
- 494 21 Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. End-to-end safe
495 reinforcement learning through barrier functions for safety-critical continuous control tasks.
496 *AAAI*, 2019.
- 497 22 Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh. A Lyapunov-based approach
498 to safe reinforcement learning. In *NIPS*, pages 8103–8112, 2018.
- 499 23 Edmund M. Clarke, Orna Grumberg, and Doron A. Peled. *Model checking*. MIT Press, 2001.
- 500 24 Jeffery A. Clouse and Paul E. Utgoff. A teaching method for reinforcement learning. In *ML*,
501 pages 92–110. Morgan Kaufmann, 1992.
- 502 25 Alexandre David, Peter Gjøøl Jensen, Kim Guldstrand Larsen, Marius Mikucionis, and
503 Jakob Haahr Taankvist. Uppaal stratego. In *TACAS*, volume 9035 of *LNCS*, pages 206–
504 211. Springer, 2015.
- 505 26 Peter Dayan and Yael Niv. Reinforcement learning: the good, the bad and the ugly. *Current
506 opinion in neurobiology*, 18(2):185–196, 2008.
- 507 27 Christian Dehnert, Sebastian Junges, Joost-Pieter Katoen, and Matthias Volk. A storm is
508 coming: A modern probabilistic model checker. In *CAV (2)*, volume 10427 of *LNCS*, pages
509 592–600. Springer, 2017.
- 510 28 Thomas G Dietterich. Hierarchical reinforcement learning with the maxq value function
511 decomposition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- 512 29 Klaus Dräger, Vojtech Forejt, Marta Z. Kwiatkowska, David Parker, and Mateusz Ujma.
513 Permissive controller synthesis for probabilistic systems. *Log. Methods Comput. Sci.*, 11(2),
514 2015.
- 515 30 Richard G Freedman and Shlomo Zilberstein. Safety in AI-HRI: Challenges complementing
516 user experience quality. In *AAAI Fall Symposium Series*, 2016.
- 517 31 Jie Fu and Ufuk Topcu. Probably approximately correct mdp learning and control with
518 temporal logic constraints. In *RSS*, 2014.

- 519 32 Nathan Fulton and André Platzer. Verifiably safe off-model reinforcement learning. 11427:413–
520 430, 2019.
- 521 33 Javier Garcia and Fernando Fernández. A comprehensive survey on safe reinforcement learning.
522 *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- 523 34 Javier García and Fernando Fernández. Probabilistic policy reuse for safe reinforcement
524 learning. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, 13(3):14, 2019.
- 525 35 OpenAI Gym. Taxi-v2, 2018. <https://gym.openai.com/envs/Taxi-v2/>.
- 526 36 Ernst Moritz Hahn, Mateo Perez, Sven Schewe, Fabio Somenzi, Ashutosh Trivedi, and Dominik
527 Wojtczak. Omega-regular objectives in model-free reinforcement learning. In *TACAS (1)*,
528 volume 11427 of *LNCS*, pages 395–412. Springer, 2019.
- 529 37 Mohammadhosein Hasanbeig, Alessandro Abate, and Daniel Kroening. Logically-correct
530 reinforcement learning. *CoRR*, abs/1801.08099, 2018.
- 531 38 Mohammadhosein Hasanbeig, Yiannis Kantaros, Alessandro Abate, Daniel Kroening, George J.
532 Pappas, and Insup Lee. Reinforcement learning for temporal logic control synthesis with
533 probabilistic satisfaction guarantees. In *CDC*, pages 5338–5343. IEEE, 2019.
- 534 39 Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable
535 mdps. *CoRR*, abs/1507.06527, 2015.
- 536 40 Manfred Jaeger, Peter Gjør Jensen, Kim Guldstrand Larsen, Axel Legay, Sean Sedwards, and
537 Jakob Haahr Taankvist. Teaching stratego to play ball: Optimal synthesis for continuous
538 space mdps. In *ATVA*, volume 11781 of *Lecture Notes in Computer Science*, pages 81–97.
539 Springer, 2019.
- 540 41 Sebastian Junges, Nils Jansen, and Sanjit A. Seshia. Enforcing almost-sure reachability in
541 pomdps. *CoRR*, abs/2007.00085, 2020.
- 542 42 Sebastian Junges, Nils Jansen, Christian Dehnert, Ufuk Topcu, and Joost-Pieter Katoen.
543 Safety-constrained reinforcement learning for MDPs. In *TACAS*, volume 9636 of *LNCS*, pages
544 130–146. Springer, 2016.
- 545 43 Joost-Pieter Katoen. The probabilistic model checking landscape. In *LICS*, pages 31–45. ACM,
546 2016.
- 547 44 Bettina Könighofer, Mohammed Alshiekh, Roderick Bloem, Laura R. Humphrey, Robert
548 Könighofer, Ufuk Topcu, and Chao Wang. Shield synthesis. *Formal Methods Syst. Des.*,
549 51(2):332–361, 2017.
- 550 45 Jan Kretínský, Guillermo A. Pérez, and Jean-François Raskin. Learning-based mean-payoff
551 optimization in an unknown MDP under omega-regular constraints. In *CONCUR*, volume 118
552 of *LIPICs*, pages 8:1–8:18. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2018.
- 553 46 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
554 convolutional neural networks. In *NIPS*, pages 1097–1105, 2012.
- 555 47 Marta Z. Kwiatkowska. Model checking for probability and time: from theory to practice. In
556 *LICS*, page 351. IEEE CS, 2003.
- 557 48 Marta Z. Kwiatkowska, Gethin Norman, and David Parker. PRISM 4.0: Verification of
558 probabilistic real-time systems. In *CAV*, volume 6806 of *LNCS*, pages 585–591. Springer, 2011.
- 559 49 Gil Lederman, Markus N. Rabe, Sanjit Seshia, and Edward A. Lee. Learning heuristics for
560 quantified boolean formulas through reinforcement learning. In *ICLR*. OpenReview.net, 2020.
- 561 50 Sergey Levine and Vladlen Koltun. Guided policy search. In *ICML (3)*, volume 28 of *JMLR
562 Workshop and Conference Proceedings*, pages 1–9. JMLR.org, 2013.
- 563 51 George Mason, Radu Calinescu, Daniel Kudenko, and Alec Banks. Assured reinforcement
564 learning with formally verified abstract policies. In *ICAART (2)*, pages 105–117. SciTePress,
565 2017.
- 566 52 Teodor Mihai Moldovan and Pieter Abbeel. Safe exploration in markov decision processes. In
567 *ICML*. icml.cc / Omnipress, 2012.
- 568 53 M. Ohnishi, L. Wang, G. Notomista, and M. Egerstedt. Barrier-certified adaptive reinforcement
569 learning with applications to brushbot navigation. *IEEE Transactions on Robotics*, pages
570 1–20, 2019.

- 571 **54** Martin Pecka and Tomas Svoboda. Safe exploration techniques for reinforcement learning—an
572 overview. In *MESAS*, pages 357–375. Springer, 2014.
- 573 **55** Amir Pnueli. The temporal logic of programs. In *Foundations of Computer Science*, pages
574 46–57. IEEE, 1977.
- 575 **56** Diederik M. Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of
576 multi-objective sequential decision-making. *J. Artif. Intell. Res.*, 48:67–113, 2013.
- 577 **57** Dorsa Sadigh, Eric S Kim, Samuel Coogan, S Shankar Sastry, and Sanjit A Seshia. A learning
578 based approach to control synthesis of markov decision processes for linear temporal logic
579 specifications. In *CDC*, pages 1091–1096. IEEE, 2014.
- 580 **58** Dorsa Sadigh, Nick Landolfi, Shankar S Sastry, Sanjit A Seshia, and Anca D Dragan. Planning
581 for cars that coordinate with people: leveraging effects on human actions for planning and
582 active information gathering over human internal state. *Autonomous Robots*, 42(7):1405–1426,
583 2018.
- 584 **59** Dorsa Sadigh, Shankar Sastry, Sanjit A Seshia, and Anca D Dragan. Planning for autonomous
585 cars that leverage effects on human actions. In *Robotics: Science and Systems*, 2016.
- 586 **60** David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den
587 Driessche, Julian Schrittwieser, Ioannis Antonoglou, Vedavyas Panneshelvam, Marc Lanctot,
588 Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy P.
589 Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering
590 the game of go with deep neural networks and tree search. *Nat.*, 529(7587):484–489, 2016.
- 591 **61** Ion Stoica, Dawn Song, Raluca Ada Popa, David Patterson, Michael W Mahoney, Randy
592 Katz, Anthony D Joseph, Michael Jordan, Joseph M Hellerstein, Joseph E Gonzalez, et al. A
593 berkeley view of systems challenges for AI. *CoRR*, abs/1712.05855, 2017.
- 594 **62** Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT Press,
595 1998.
- 596 **63** Florent Teichteil-Königsbuch. Stochastic safest and shortest path problems. In *AAAI*. AAAI
597 Press, 2012.
- 598 **64** Matteo Turchetta, Felix Berkenkamp, and Andreas Krause. Safe exploration for interactive
599 machine learning. In *NeurIPS*, pages 2887–2897, 2019.
- 600 **65** Abhinav Verma, Hoang Minh Le, Yisong Yue, and Swarat Chaudhuri. Imitation-projected
601 programmatic reinforcement learning. In *NeurIPS*, pages 15726–15737, 2019.
- 602 **66** Abhinav Verma, Vijayaraghavan Murali, Rishabh Singh, Pushmeet Kohli, and Swarat Chaud-
603 huri. Programmatically interpretable reinforcement learning. In *ICML*, volume 80 of *Proceedings*
604 *of Machine Learning Research*, pages 5052–5061. PMLR, 2018.
- 605 **67** Angelina Wang, Thanard Kurutach, Kara Liu, Pieter Abbeel, and Aviv Tamar. Learning
606 robotic manipulation through visual planning and acting. arXiv preprint arXiv:1905.04411,
607 2019.
- 608 **68** Min Wen, Rüdiger Ehlers, and Ufuk Topcu. Correct-by-synthesis reinforcement learning with
609 temporal logic constraints. In *IROS*, 2015.
- 610 **69** Douglas J White. Real applications of Markov decision processes. *Interfaces*, 15(6):73–83,
611 1985.
- 612 **70** Ian H Witten, Eibe Frank, Mark A Hall, and Christopher J Pal. *Data Mining: Practical*
613 *machine learning tools and techniques*. Morgan Kaufmann, 2016.
- 614 **71** Meng Wu, Jingbo Wang, Jyotirmoy Deshmukh, and Chao Wang. Shield synthesis for real:
615 Enforcing safety in cyber-physical systems. In *FMCAD*, pages 129–137. IEEE, 2019.
- 616 **72** Weichao Zhou and Wenchao Li. Safety-aware apprenticeship learning. In *CAV (1)*, volume
617 10981 of *Lecture Notes in Computer Science*, pages 662–680. Springer, 2018.
- 618 **73** He Zhu, Zikang Xiong, Stephen Magill, and Suresh Jagannathan. An inductive synthesis
619 framework for verifiable reinforcement learning. In *PLDI*, pages 686–701. ACM, 2019.