Value Iteration for Simple Stochastic Games: Stopping Criterion and Learning Algorithm^{*}

Edon Kelmendi, Julia Krämer, Jan Křetínský, and Maximilian Weininger



Technical University of Munich

Abstract. Simple stochastic games can be solved by value iteration (VI), which yields a sequence of under-approximations of the value of the game. This sequence is guaranteed to converge to the value only in the limit. Since no stopping criterion is known, this technique does not provide any guarantees on its results. We provide the first stopping criterion for VI on simple stochastic games. It is achieved by additionally computing a convergent sequence of *over-approximations* of the value, relying on an analysis of the game graph. Consequently, VI becomes an anytime algorithm returning the approximation of the value and the current error bound. As another consequence, we can provide a simulation-based asynchronous VI algorithm, which yields the same guarantees, but without necessarily exploring the whole game graph.

1 Introduction

Simple stochastic game (SG) [Con92] is a zero-sum two-player game played on a graph by Maximizer and Minimizer, who choose actions in their respective vertices (also called states). Each action is associated with a probability distribution determining the next state to move to. The objective of Maximizer is to maximize the probability of reaching a given target state; the objective of Minimizer is the opposite.

Stochastic games constitute a fundamental problem for several reasons. From the theoretical point of view, the complexity of this problem¹ is known to be in $\mathbf{UP} \cap \mathbf{coUP}$ [HK66], but no polynomial-time algorithm is known. Further, several other important problems can be reduced to SG, for instance parity games, mean-payoff games, discounted-payoff games and their stochastic extensions [CF11]. The task of solving SG is also polynomial-time equivalent to solving perfect information Shapley, Everett and Gillette games [AM09]. Besides, the problem is practically relevant in verification and synthesis. SG can model reactive systems, with players corresponding to the controller of the system and

^{*} This research was funded in part by the German Excellence Initiative and the European Union Seventh Framework Programme under grant agreement No. 291763 for TUM – IAS, the Studienstiftung des deutschen Volkes project "Formal methods for analysis of attack-defence diagrams", the Czech Science Foundation grant No. 18-11193S, TUM IGSSE Grant 10.06 (PARSEC), and the German Research Foundation (DFG) project KR 4890/2-1 "Statistical Unbounded Verification".

¹ Formally, the problem is to decide, for a given $p \in [0, 1]$ whether Maximizer has a strategy ensuring probability at least p to reach the target.

to its environment, where quantified uncertainty is explicitly modelled. This is useful in many application domains, ranging from smart energy management [CFK⁺13b] to autonomous urban driving [CKSW13], robot motion planning [LaV00] to self-adaptive systems [CMG14]; for various recent case studies, see e.g. [SK16]. Finally, since Markov decision processes (MDP) [Put14] are a special case with only one player, SG can serve as abstractions of large MDP [KKNP10].

Solution techniques There are several classes of algorithms for solving SG, most importantly strategy iteration (SI) algorithms [HK66] and value iteration (VI) algorithms [Con92]. Since the repetitive evaluation of strategies in SI is often slow in practice, VI is usually preferred, similarly to the special case of MDPs [KM17]. For instance, the most used probabilistic model checker PRISM [KNP11] and its branch PRISM-Games [CFK⁺13a] use VI for MDP and SG as the default option, respectively. However, while SI is in principle a precise method, VI is an approximative method, which converges only in the limit. Unfortunately, there is no known stopping criterion for VI applied to SG. Consequently, there are no guarantees on the results returned in finite time. Therefore, current tools stop when the difference between the two most recent approximations is low, and thus may return arbitrarily imprecise results [HM17].

Value iteration with guarantees In the special case of MDP, in order to obtain bounds on the imprecision of the result, one can employ a bounded variant of VI [MLG05,BCC⁺14] (also called *interval iteration* [HM17]). Here one computes not only an under-approximation, but also an over-approximation of the actual value as follows. On the one hand, iterative computation of the least fixpoint of Bellman equations yields an under-approximating sequence converging to the value. On the other hand, iterative computation of the greatest fixpoint yields an over-approximation, which, however, does not converge to the value. Moreover, it often results in the trivial bound of 1. A solution suggested for MDPs [BCC⁺14,HM17] is to modify the underlying graph, namely to collapse end components. In the resulting MDP there is only one fixpoint, thus the least and greatest fixpoint coincide and both approximating sequences converge to the actual value. In contrast, for general SG no procedure where the greatest fixpoint converges to the value is known. In this paper we provide one, yielding a stopping criterion. We show that the pre-processing approach of collapsing is not applicable in general and provide a solution on the original graph. We also characterize SG where the fixpoints coincide and no processing is needed. The main technical challenge is that states in an end component in SG can have different values, in contrast to the case of MDP.

Practical efficiency using guarantees We further utilize the obtained guarantees to practically improve our algorithm. Similar to the MDP case [BCC⁺14], the quantification of the error allows for ignoring parts of the state space, and thus a speed up without jeopardizing the correctness of the result. Indeed, we provide a technique where some states are not explored and processed at all, but their potential effect is still taken into accountThe information is further used to decide the states to be explored next and to be analyzed in more detail. To this end, simulations and learning are used as tools. While for MDP this idea

has already demonstrated speed ups in orders of magnitude [BCC⁺14,ACD⁺17], this paper provides the first technique of this kind for SG. **Our contribution** is summarized as follows

- We introduce a VI algorithm yielding both under- and over-approximation sequences, both of which converge to the value of the game. Thus we present the first stopping criterion for VI on SG and the first anytime algorithm with guaranteed precision. We also characterize when a simpler solution is sufficient.
- We provide a learning-based algorithm, which preserves the guarantees, but is in some cases more efficient since it avoids exploring the whole state space.
- We evaluate the running times of the algorithms experimentally, concluding that obtaining guarantees requires an overhead that is either negligible or mitigated by the learning-based approach.

Related work The works closest to ours are the following. As mentioned above, [BCC⁺14,HM17] describe the solution to the special case of MDP. While [BCC⁺14] also provides a learning-based algorithm, [HM17] discusses the convergence rate and the exact solution. The basic algorithm of [HM17] is implemented in PRISM [BKL⁺17] and the learning approach of [BCC⁺14] in STORM [DJKV17a]. The extension for SG where the interleaving of players is severely limited (every end component belongs to one player only) is discussed in [Ujm15].

Further, in the area of probabilistic planning, bounded real-time dynamic programming [MLG05] is related to our learning-based approach. However, it is limited to the setting of stopping MDP where the target sink or the non-target sink is reached almost surely under any pair of strategies and thus the fixpoints coincide. Our algorithm works for general SG, not only for stopping ones, without any blowup.

For SG, the tools implementing the standard SI and/or VI algorithms are PRISM-games [CFK⁺13a], GAVS+ [CKLB11] and GIST [CHJR10]. The latter two are, however, neither maintained nor accessible via the links provided in their publications any more.

Apart from fundamental algorithms to solve SG, there are various practically efficient heuristics that, however, provide none or weak guarantees, often based on some form of learning [BT00,LL08,WT16,TT16,AY17,BBS08]. Finally, the only currently available way to obtain any guarantees through VI is to perform γ^2 iterations and then round to the nearest multiple of $1/\gamma$, yielding the value of the game with precision $1/\gamma$ [CH08]; here γ cannot be freely chosen, but it is a fixed number, exponential in the number of states and the used probability denominators. However, since the precision cannot be chosen and the number of iterations is always exponential, this approach is infeasible even for small games. **Organization of the paper** Section 2 introduces the basic notions and revises value iteration. Section 3 explains the idea of our approach on an example. Section 4 provides a full technical treatment of the method as well as the learningbased variation. Section 5 discusses experimental results and Section 6 concludes. The appendix (available in [KKKW18]) gives technical details on the pseudocode as well as the conducted experiments and provides more extensive proofs to the theorems and lemmata; in this paper, there are only proof sketches and ideas.

2 Preliminaries

2.1 Basic definitions

A probability distribution on a finite set X is a mapping $\delta : X \to [0, 1]$, such that $\sum_{x \in X} \delta(x) = 1$. The set of all probability distributions on X is denoted by $\mathcal{D}(X)$. Now we define stochastic games, in literature often referred as simple stochastic games or stochastic two-player games with a reachability objective.

Definition 1 (SG). A stochastic game (SG) is a tuple $(S, S_{\Box}, S_{\bigcirc}, \mathbf{s}_0, \mathbf{A}, \mathbf{Av}, \delta, \mathbf{1}, \mathbf{0})$, where S is a finite set of states partitioned into the sets S_{\Box} and S_{\bigcirc} of states of the player Maximizer and Minimizer, respectively, $\mathbf{s}_0, \mathbf{1}, \mathbf{0} \in S$ is the initial state, target state, and sink state, respectively, \mathbf{A} is a finite set of actions, $\mathbf{Av} : S \to 2^{\mathbf{A}}$ assigns to every state a set of available actions, and $\delta : S \times \mathbf{A} \to \mathcal{D}(S)$ is a transition function that given a state \mathbf{s} and an action $\mathbf{a} \in \mathbf{Av}(\mathbf{s})$ yields a probability distribution over successor states.

A Markov decision process (MDP) is a special case of SG where $S_{\bigcirc} = \emptyset$. We assume that SGs are non-blocking, so for all states s we have $Av(s) \neq \emptyset$. Further, 1 and o only have one action and it is a self-loop with probability 1. Additionally, we can assume that the SG is preprocessed so that all states with no path to 1 are merged with o.

For a state s and an available action $\mathbf{a} \in \mathsf{Av}(\mathbf{s})$, we denote the set of successors by $\mathsf{Post}(\mathbf{s}, \mathbf{a}) := \{\mathbf{s}' \mid \delta(\mathbf{s}, \mathbf{a}, \mathbf{s}') > 0\}$. Finally, for any set of states $T \subseteq S$, we use T_{\Box} and T_{\bigcirc} to denote the states in T that belong to Maximizer and Minimizer, whose states are drawn in the figures as \Box and \bigcirc , respectively.

The semantics of SG is given in the usual way by means of strategies and the induced Markov chain and the respective probability space, as follows. An *infinite path* ρ is an infinite sequence $\rho = s_0 a_0 s_1 a_1 \ldots \in (S \times A)^{\omega}$, such that for every $i \in \mathbb{N}$, $\mathbf{a}_i \in \mathsf{Av}(\mathbf{s}_i)$ and $\mathbf{s}_{i+1} \in \mathsf{Post}(\mathbf{s}_i, \mathbf{a}_i)$. Finite paths are defined analogously as elements of $(S \times A)^* \times S$. Since this paper deals with the reachability objective, we can restrict our attention to memoryless strategies, which are optimal for this objective. We still allow randomizing strategies, because they are needed for the learning-based algorithm later on. A strategy of Maximizer or Minimizer is a function $\sigma : S_{\Box} \to \mathcal{D}(A)$ or $S_{\bigcirc} \to \mathcal{D}(A)$, respectively, such that $\sigma(\mathbf{s}) \in \mathcal{D}(\mathsf{Av}(\mathbf{s}))$ for all \mathbf{s} . We call a strategy deterministic if it maps to Dirac distributions only. Note that there are finitely many deterministic strategies. A pair (σ, τ) of strategies of Maximizer and Minimizer induces a Markov chain $\mathsf{G}^{\sigma,\tau}$ where the transition probabilities are defined as $\delta(\mathbf{s}, \mathbf{s}') = \sum_{\mathbf{a} \in \mathsf{Av}(\mathbf{s})} \sigma(\mathbf{s}, \mathbf{a}) \cdot \delta(\mathbf{s}, \mathbf{a}, \mathbf{s}')$ for states of Maximizer and analogously for states of Minimizer, with σ replaced by τ . The Markov chain induces a unique probability distribution $\mathbb{P}_{\mathbf{s}}^{\sigma,\tau}$ over measurable sets of infinite paths [BK08, Ch. 10].

We write $\Diamond \mathbf{1} := \{ \rho \mid \exists i \in \mathbb{N}. \ \rho(i) = \mathbf{1} \}$ to denote the (measurable) set of all paths which eventually reach $\mathbf{1}$. For each $\mathbf{s} \in S$, we define the *value* in \mathbf{s} as

$$\mathsf{V}(\mathsf{s}) := \sup_{\sigma} \inf_{\tau} \mathbb{P}_{s}^{\sigma,\tau}(\Diamond \mathfrak{l}) = \inf_{\tau} \sup_{\sigma} \mathbb{P}_{s}^{\sigma,\tau}(\Diamond \mathfrak{l}),$$

where the equality follows from [Mar75]. We are interested not only in $V(s_0)$, but also its ε -approximations and the corresponding (ε -)optimal strategies for both players.

Now we recall a fundamental tool for analysis of MDP called end components. We introduce the following notation. Given a set of states $T \subseteq S$, a state $s \in T$ and an action $a \in Av(s)$, we say that (s, a) exits T if $Post(s, a) \not\subseteq T$. We define an end component of a SG as the end component of the underlying MDP with both players unified.

Definition 2 (EC). A non-empty set $T \subseteq S$ of states is an end component (EC) if there is a non-empty set $B \subseteq \bigcup_{s \in T} Av(s)$ of actions such that

- 1. for each $s \in T$, $a \in B \cap Av(s)$ we do not have (s, a) exits T,
- 2. for each $s, s' \in T$ there is a finite path $w = sa_0 \dots a_n s' \in (T \times B)^* \times T$, i.e. the path stays inside T and only uses actions in B.

Intuitively, ECs correspond to bottom strongly connected components of the Markov chains induced by possible strategies, so for some pair of strategies all possible paths starting in the EC remain there. An end component T is a maximal end component (MEC) if there is no other end component T' such that $T \subseteq T'$. Given an SG G, the set of its MECs is denoted by MEC(G) and can be computed in polynomial time [CY95].

2.2 (Bounded) value iteration

The value function V satisfies the following system of equations, which is referred to as the *Bellman equations*:

$$V(\mathbf{s}) = \begin{cases} \max_{\mathbf{a} \in A_{\mathbf{v}}(\mathbf{s})} V(\mathbf{s}, \mathbf{a}) & \text{if } \mathbf{s} \in S_{\Box} \\ \min_{\mathbf{a} \in A_{\mathbf{v}}(\mathbf{s})} V(\mathbf{s}, \mathbf{a}) & \text{if } \mathbf{s} \in S_{\Box} \\ 1 & \text{if } \mathbf{s} = \mathbf{1} \\ 0 & \text{if } \mathbf{s} = \mathbf{o} \end{cases}$$
(1)

 $where^2$

$$\mathsf{V}(\mathsf{s},\mathsf{a}) := \sum_{s' \in S} \delta(\mathsf{s},\mathsf{a},\mathsf{s}') \cdot \mathsf{V}(\mathsf{s}')$$
⁽²⁾

Moreover, V is the *least* solution to the Bellman equations, see e.g. [CH08]. To compute the value of V for all states in an SG, one can thus utilize the iterative approximation method *value iteration* (VI) as follows. We start with a lower bound function $L_0: S \to [0, 1]$ such that $L_0(1) = 1$ and, for all other $s \in S$, $L_0(s) = 0$. Then we repetitively apply Bellman updates (3) and (4)

$$\mathsf{L}_{n}(\mathsf{s},\mathsf{a}) \coloneqq \sum_{s' \in S} \delta(\mathsf{s},\mathsf{a},\mathsf{s}') \cdot \mathsf{L}_{n-1}(\mathsf{s}')$$
(3)

$$\mathsf{L}_{n}(\mathsf{s}) := \begin{cases} \max_{\mathsf{a}\in\mathsf{Av}(\mathsf{s})}\mathsf{L}_{n}(\mathsf{s},\mathsf{a}) & \text{if } \mathsf{s}\in S_{\Box} \\ \min_{\mathsf{a}\in\mathsf{Av}(\mathsf{s})}\mathsf{L}_{n}(\mathsf{s},\mathsf{a}) & \text{if } \mathsf{s}\in S_{\Box} \end{cases}$$
(4)

² Throughout the paper, for any function $f: S \to [0, 1]$ we overload the notation and also write $f(\mathbf{s}, \mathbf{a})$ meaning $\sum_{s' \in S} \delta(\mathbf{s}, \mathbf{a}, \mathbf{s}') \cdot f(\mathbf{s}')$.

until convergence. Note that convergence may happen only in the limit even for such a simple game as in Figure 1 on the left. The sequence is monotonic, at all times a *lower* bound on V, i.e. $L_i(s) \leq V(s)$ for all $s \in S$, and the least fixpoint satisfies $L^* := \lim_{n\to\infty} L_n = V$.

Unfortunately, there is no known stopping criterion, i.e. no guarantees how close the current under-approximation is to the value [HM17]. The current tools stop when the difference between two successive approximations is smaller than a certain threshold, which can lead to arbitrarily wrong results [HM17].

For the special case of MDP, it has been suggested to also compute the greatest fixpoint [MLG05] and thus an *upper* bound as follows. The function $G: S \to [0, 1]$ is initialized for all states $s \in S$ as $G_0(s) = 1$ except for $G_0(o) = 0$. Then we repetitively apply updates (3) and (4), where L is replaced by G. The resulting sequence G_n is monotonic, provides an upper bound on V and the greatest fixpoint $G^* := \lim_n G_n$ is the greatest solution to the Bellman equations on $[0, 1]^S$.

This approach is called *bounded value iteration (BVI)* (or *bounded realtime dynamic programming (BRTDP)* [MLG05,BCC⁺14] or *interval iteration* [HM17]). If $L^* = G^*$ then they are both equal to V and we say that *BVI converges.* BVI is guaranteed to converge in MDP if the only ECs are those of 1 and o [BCC⁺14]. Otherwise, if there are non-trivial ECs they have to be "collapsed"³. Computing the greatest fixpoint on the modified MDP results in another sequence U_i of upper bounds on V, converging to U^{*} := $\lim_{n} U_{n}$. Then BVI converges even for general MDPs, U^{*} = V [BCC⁺14], when transformed this way. The next section illustrates this difficulty and the solution through collapsing on an example.

In summary, all versions of BVI discussed so far and later on in the paper follow the pattern of Algorithm 1. In the naive version, UPDATE just performs the Bellman update on L and U according to Equations (3) and (4).⁴ For a general MDP, U does not converge to V, but to G^* , and thus the termination criterion may never be met if $G^*(s_0) - V(s_0) > 0$. If the ECs are collapsed in pre-processing then U converges to V.

For the general case of SG, the collapsing approach fails and this paper provides another version of BVI where U converges to V, based on a more detailed structural analysis of the game.

3 Example

In this section, we illustrate the issues preventing BVI convergence and our solution on a few examples. Recall that G is the sequence converging to the greatest solution of the Bellman equations, while U is in general any sequence over-approximating V that one or another BVI algorithm suggests.

Firstly, we illustrate the issue that arises already for the special case of MDP. Consider the MPD of Figure 1 on the left. Although V(s) = V(t) = 0.5, we have

³ All states of an EC are merged into one, all leaving actions are preserved and all other actions are discarded. For more detail see [KKKW18, Appendix A.1.]

⁴ For the straightforward pseudocode, see [KKKW18, Appendix A.2.].

Algorithm 1 Bounded value iteration algorithm

1:	: procedure BVI(precision $\epsilon > 0$)						
2:	for $s \in S$ do	* Initialization $*$					
3:	L(s) = 0	$\ \ \ \ \ \ \ \ \ \ \ \ \ $					
4:	U(s) = 1	$\ \ \ \ \ \ \ \ \ \ \ \ \ $					
5:	$L(\mathtt{1}) = 1$	$\$ Value of sinks is determined a priori * $\$					
6:	$U(\mathfrak{o})=0$						
7:	repeat						
8:	UPDATE(L,	U) $\$ Bellman updates or their modification * $\$					
9:	until $U(s_0) - L$	$(\mathbf{s}_0) < \epsilon$ * Guaranteed error bound * \					

 $G_i(s) = G_i(t) = 1$ for all *i*. Indeed, the upper bound for t is always updated as the maximum of $G_i(t, c)$ and $G_i(t, b)$. Although $G_i(t, c)$ decreases over time, $G_i(t, b)$ remains the same, namely equal to $G_i(s)$, which in turn remains equal to $G_i(s, a) = G_i(t)$. This cyclic dependency lets both s and t remain in an "illusion" that the value of the other one is 1.

The solution for MDP is to remove this cyclic dependency by collapsing all MECs into singletons and removing the resulting purely self-looping actions. Figure 1 in the middle shows the MDP after collapsing the EC $\{s, t\}$. This turns the MDP into a stopping one, where 1 or o is under any strategy reached with probability 1. In such MDP, there is a unique solution to the Bellman equations. Therefore, the greatest fixpoint is equal to the least one and thus to V.

Secondly, we illustrate the issues that additionally arise for general SG. It turns out that the collapsing approach can be extended only to games where all states of each EC belong to one player only [Ujm15]. In this case, both Maximizer's and Minimizer's ECs are collapsed the same way as in MDP.

However, when both players are present in an EC, then collapsing may not solve the issue. Consider the SG of Figure 2. Here α and β represent the values of the respective actions.⁵ There are three cases:

First, let $\alpha < \beta$. If the bounds converge to these values we eventually observe $G_i(q, e) < L_i(r, f)$ and learn the induced inequality. Since p is a Minimizer's state it will never pick the action leading to the greater value of β . Therefore, we can safely merge p and q, and remove the action leading to r, as shown in the second subfigure.

Second, if $\alpha > \beta$, **p** and **r** can be merged in an analogous way, as shown in the third subfigure.

Third, if $\alpha = \beta$, both previous solutions as well as collapsing all three states as in the fourth subfigure is possible. However, since the approximants may only converge to α and β in the limit, we may not know in finite time which of these cases applies and thus cannot decide for any of the collapses.

Consequently, the approach of collapsing is not applicable in general. In order to ensure BVI convergence, we suggest a different method, which we call

⁵ Precisely, we consider them to stand for a probabilistic branching with probability α (or β) to 1 and with the remaining probability to \mathfrak{o} . To avoid clutter in the figure, we omit this branching and depict only the value.

deflating. It does not involve changing the state space, but rather decreasing the upper bound U_i to the least value that is currently provable (and thus still correct). To this end, we analyze the exiting actions, i.e. with successors outside of the EC, for the following reason. If the play stays in the EC forever, the target is never reached and Minimizer wins. Therefore, Maximizer needs to pick some exiting action to avoid staying in the EC.



Fig. 1: Left: An MDP (as special case of SG) where BVI does not converge due to the grayed EC. Middle: The same MDP where the EC is collapsed, making BVI converge. Right: The approximations illustrating the convergence of the MDP in the middle.



Fig. 2: Left: Collapsing ECs in SG may lead to incorrect results. The Greek letters on the leaving arrows denote the values of the exiting actions. Right three figures: Correct collapsing in different cases, depending on the relationship of α and β . In contrast to MDP, some actions of the EC exiting the collapsed part have to be removed.

For the EC with the states s and t in Figure 1, the only exiting action is c. In this example, since c is the only exiting action, $U_i(t, c)$ is the highest possible upper bound that the EC can achieve. Thus, by decreasing the upper bound of all states in the EC to that number⁶, we still have a safe upper bound. Moreover, with this modification BVI converges in this example, intuitively because now the upper bound of t depends on action c as it should.

For the example in Figure 2, it is correct to decrease the upper bound to the maximal exiting one, i.e. $\max{\{\hat{\alpha}, \hat{\beta}\}}$, where $\hat{\alpha} := \mathsf{U}_i(\mathsf{a}), \hat{\beta} := \mathsf{U}_i(\mathsf{b})$ are the

 $^{^6}$ We choose the name "deflating" to evoke decreasing the overly high "pressure" in the EC until it equalizes with the actual "pressure" outside.

current approximations of α and of β . However, this itself does not ensure BVI convergence. Indeed, if for instance $\hat{\alpha} < \hat{\beta}$ then deflating all states to $\hat{\beta}$ is not tight enough, as values of **p** and **q** can even be bounded by $\hat{\alpha}$. In fact, we have to find a certain sub-EC that corresponds to $\hat{\alpha}$, in this case {**p**, **q**} and set all its upper bounds to $\hat{\alpha}$. We define and compute these sub-ECs in the next section.

In summary, the general structure of our convergent BVI algorithm is to produce the sequence U by application of Bellman updates and occasionally find the relevant sub-ECs and deflate them. The main technical challenge is that states in an EC in SG can have different values, in contrast to the case of MDP.

4 Convergent Over-approximation

In Section 4.1, we characterize SGs where Bellman equations have more solutions. Based on the analysis, subsequent sections show how to alter the procedure computing the sequence G_i over-approximating V so that the resulting tighter sequence U_i still over-approximates V, but also converges to V. This ensures that thus modified BVI converges. Section 4.4 presents the learning-based variant of our BVI.

4.1 Bloated end components cause non-convergence

As we have seen in the example of Fig. 2, BVI generally does not converge due to ECs with a particular structure of the exiting actions. The analysis of ECs relies on the extremal values that can be achieved by exiting actions (in the example, α and β). Given the value function V or just its current overapproximation U_i , we define the most profitable exiting action for Maximizer (denoted by \Box) and Minimizer (denoted by \bigcirc) as follows.

Definition 3 (bestExit). Given a set of states $T \subseteq S$ and a function $f : S \rightarrow [0,1]$ (see footnote 2), the f-value of the best T-exiting action of Maximizer and Minimizer, respectively, is defined as

$$bestExit_{f}^{\square}(T) = \max_{\substack{\mathbf{s} \in T_{\square} \\ (\mathbf{s}, \mathbf{a}) \text{ exits } T}} f(\mathbf{s}, \mathbf{a})$$
$$bestExit_{f}^{\bigcirc}(T) = \min_{\substack{\mathbf{s} \in T_{\bigcirc} \\ (\mathbf{s}, \mathbf{a}) \text{ exits } T}} f(\mathbf{s}, \mathbf{a})$$

with the convention that $\max_{\emptyset} = 0$ and $\min_{\emptyset} = 1$.

Example 1. In the example of Fig. 2 on the left with $T = \{\mathsf{p}, \mathsf{q}, \mathsf{r}\}$ and $\alpha < \beta$, we have $\mathsf{bestExit}_V^{\Box}(T) = \beta$, $\mathsf{bestExit}_V^{\Box}(T) = 1$. It is due to $\beta < 1$ that BVI does not converge here. We generalize this in the following lemma.

Lemma 1. Let T be an EC. For every m satisfying $\text{bestExit}_V^{\square}(T) \leq m \leq \text{bestExit}_V^{\square}(T)$, there is a solution $f: S \to [0, 1]$ to the Bellman equations, which on T is constant and equal to m.

Proof (Idea). Intuitively, such a constant m is a solution to the Bellman equations on T for the following reasons. As both players prefer getting m to exiting and getting "only" the values of their respective **bestExit**, they both choose to stay in the EC (and the extrema in the Bellman equations are realized on non-exiting actions). On the one hand, Maximizer (Bellman equations with max) is hoping for the promised m, which is however not backed up by any actions actually exiting towards the target. On the other hand, Minimizer (Bellman equations with min) does not realize that staying forever results in her optimal value 0 instead of m.

Corollary 1. If $\text{bestExit}_{V}^{\bigcirc}(T) > \text{bestExit}_{V}^{\bigcirc}(T)$ for some ECT, then $G^* \neq V$.

Proof. Since there are m_1, m_2 such that $\mathsf{bestExit}_{\mathsf{V}}^{\Box}(T) < m_1 < m_2 < \mathsf{bestExit}_{\mathsf{V}}^{\bigcirc}(T)$, by Lemma 1 there are two different solutions to the Bellman equations. In particular, $\mathsf{G}^* > \mathsf{L}^* = \mathsf{V}$, and BVI does not converge.

In accordance with our intuition that ECs satisfying the above inequality should be deflated, we call them bloated.

Definition 4 (BEC). An ECT is called a bloated end component (BEC), if $\text{bestExit}_{V}^{\bigcirc}(T) > \text{bestExit}_{V}^{\bigcirc}(T)$.

Example 2. In the example of Fig. 2 on the left with $\alpha < \beta$, the ECs {p,q} and {p,q,r} are BECs.

Example 3. If an EC T has no exiting actions of Minimizer (or no Minimizer's states at all, as in an MDP), then $\mathsf{bestExit}_V^{\bigcirc}(T) = 1$ (the case with \min_{\emptyset}). Hence all numbers between $\mathsf{bestExit}_V^{\square}(T)$ and 1 are a solution to the Bellman equations and $\mathsf{G}^*(s) = 1$ for all states $s \in T$.

Analogously, if Maximizer does not have any exiting action in T, then it holds that $\mathsf{bestExit}_V^{\Box}(T) = 0$ (the case with \max_{\emptyset}), T is a BEC and all numbers between 0 and $\mathsf{bestExit}_V^{\bigcirc}(T)$ are a solution to the Bellman equations.

Note that in MDP all ECs belong to one player, namely Maximizer. Consequently, all ECs are BECs except for ECs where Maximizer has an exiting action with value 1; all other ECs thus have to be collapsed (or deflated) to ensure BVI convergence in MDPs. Interestingly, all non-trivial ECs in MDPs are a problem, while in SGs through the presence of the other player some ECs can converge, namely if both players want to exit (See e.g. [KKKW18, Appendix A.3.]). \triangle

We show that BECs are indeed the only obstacle for BVI convergence.

Theorem 1. If the SG contains no BECs except for $\{0\}$ and $\{1\}$, then $G^* = V$.

Proof (Sketch). Assume, towards a contradiction, that there is some state s with a positive difference $G^*(s) - V(s) > 0$. Consider the set D of states with the maximal difference. D can be shown to be an EC. Since it is not a BEC there has to be an action exiting D and realizing the optimum in that state. Consequently, this action also has the maximal difference, and all its successors, too. Since some of the successors are outside of D, we get a contradiction with the maximality of D.

In Section 4.2, we show how to eliminate BECs by collapsing their "core" parts, called below MSECs (maximal simple end components). Since MSECs can only be identified with enough information about V, Section 4.3 shows how to avoid direct *a priori* collapsing and instead dynamically deflate candidates for MSECs in a conservative way.

4.2 Static MSEC decomposition

Now we turn our attention to SG with BECs. Intuitively, since in a BEC all Minimizer's exiting actions have a higher value than what Maximizer can achieve, Minimizer does not want to use any of his own exiting actions and prefers staying in the EC (or steering Maximizer towards his worse exiting actions). Consequently, only Maximizer wants to take an exiting action. In the MDP case he can pick any desirable one. Indeed, he can wait until he reaches a state where it is available. As a result, in MDP all states of an EC have the *same value* and can all be collapsed into one state. In the SG case, he may be restricted by Minimizer's behaviour or even not given any chance to exit the EC at all. As a result, a BEC may contain several parts (below denoted MSECs), each with different value, intuitively corresponding to different exits. Thus instead of MECs, we have to decompose into finer MSECs and only collapse these.

Definition 5 (Simple EC). An ECT is called simple (SEC), if for all $s \in T$ we have $V(s) = bestExit_{V}^{\Box}(T)$.

A SEC C is maximal (MSEC) if there is no SEC C' such that $C \subsetneq C'$.

Intuitively, an EC is simple, if Minimizer cannot keep Maximizer away from his bestExit. Independently of Minimizer's decisions, Maximizer can reach the bestExit almost surely, unless Minimizer decides to leave, in which case Maximizer could achieve an even higher value.

Example 4. Assume $\alpha < \beta$ in the example of Figure 2. Then $\{\mathbf{p}, \mathbf{q}\}$ is a SEC and an MSEC. Further observe that action **c** is sub-optimal for Minimizer and removing it does not affect the value of any state, but simplifies the graph structure. Namely, it destructs the whole EC into several (here only one) SECs and some non-EC states (here **r**).

Algorithm 2, called FIND_MSEC, shows how to compute MSECs. It returns the set of all MSECs if called with parameter V. However, later we also call this function with other parameters $f: S \to [0, 1]$. The idea of the algorithm is the following. The set X consists of Minimizer's sub-optimal actions, leading to a higher value. As such they cannot be a part of any SEC and thus should be ignored when identifying SECs. (The previous example illustrates that ignoring X is indeed safe as it does not change the value of the game.) We denote the game G where the available actions Av are changed to the new available actions Av' (ignoring the Minimizer's sub-optimal ones) as $G_{[Av/Av']}$. Once removed, Minimizer has no choices to affect the value and thus each EC is simple.

Algorithm 2 FIND_MSEC

1:	function	F	IND.	MS	EC	C(j	f: S	$5 \rightarrow$	[0, 1])	
		~ /	~				- /		- () -) -	-

2:

3:

4:

Lemma 2 (Correctness of Algorithm 2). $T \in FIND_{MSEC}(V)$ if and only if T is a MSEC.

Proof (Sketch). "If": As T is an MSEC, all states in T have the value $\mathsf{bestExit}_V^{\square}(T)$, and hence also all actions that stay inside T have this value. Thus, no action that stays in T is removed by Line 3 and it is still a MEC in the modified game.

"Only if": If $T \in \mathsf{FIND}_{\mathsf{MSEC}}(\mathsf{V})$, then T is a MEC of the game where the suboptimal available actions (those in X) of Minimizer have been removed. Hence for all $s \in T$: $V(s) = bestExit_V^{\square}(T)$, because intuitively Minimizer has no possibility to influence the value any further, since all actions that could do so were in X and have been removed. Since T is a MEC in the modified game, it certainly is an EC in the original game. Hence T is a SEC. The inclusion maximality follows from the fact that we compute MECs in the modified game. Thus T is an MSEC.

Remark 1 (Algorithm with an oracle). In Section 3, we have seen that collapsing MECs does not ensure BVI convergence. Collapsing does not preserve the values. since in BECs we would be collapsing states with different values. Hence we want to collapse only MSECs, where the values are the same. If, moreover, we remove X in such a collapsed SG, then there are no (non-sink) ECs and BVI converges on this SG to the original value.

The difficulty with this algorithm is that it requires an oracle to compare values, for instance a sufficiently precise approximation of V. Consequently, we cannot pre-compute the MSECs, but have to find them while running BVI. Moreover, since the approximations converge only in the limit we may never be able to conclude on simplicity of some ECs. For instance, if $\alpha = \beta$ in Figure 2, and if the approximations converge at different speeds, then Algorithm 2 always outputs only a part of the EC, although the whole EC on $\{p, q, r\}$ is simple.

In MDPs, all ECs are simple, because there is no second player to be resolved and all states in an EC have the same value. Thus for MDPs it suffices to collapse all MECs, in contrast to SG.

4.3**Dynamic MSEC decomposition**

Since MSECs cannot be identified from approximants of V for sure, we refrain from collapsing⁷ and instead only decrease the over-approximation in the corresponding way. We call the method *deflating*, by which we mean decreasing the upper bound of all states in an EC to its $bestExit_{U}^{\Box}$, see Algorithm 3.

 $^{^{7}}$ Our subsequent method can be combined with local collapsing whenever the lower and upper bounds on V are conclusive.

The procedure DEFLATE (called on the current upper bound U_i) decreases this upper bound to the minimum possible value according to the current approximation and thus prevents states from only depending on each other, as in SECs. Intuitively, it gradually approximates SECs and performs the corresponding adjustments, but does not commit to any of the approximations.

Algorithm 3 DEFLATE

1: function DEFLATE(EC $T, f : S \to [0, 1]$) 2: for $\mathbf{s} \in T$ do 3: $f(\mathbf{s}) \leftarrow \min(f(\mathbf{s}), \mathsf{bestExit}_f^{\Box}(T))$ * Decrease the upper bound * \ 4: return f

Lemma 3 (DEFLATE is sound). For any $f : S \to [0, 1]$ such that $f \ge V$ and any ECT, $\mathsf{DEFLATE}(T, f) \ge V$.

This allows us to define our BVI algorithm as the naive BVI with only the additional lines 3-4, see Algorithm 4.

Algo	rithm 4 UPDATE procedure for	or bounded value iteration on SG				
1: procedure UPDATE(L : $S \rightarrow [0, 1], U : S \rightarrow [0, 1])$						
2:	L,U get updated according to E	Eq. (3) and (4) \setminus * Bellman updates * \setminus				
3:	for $T \in FIND_MSEC(L)$ do	$\$ Use lower bound to find ECs * $\$				
4:	$U \leftarrow DEFLATE(T,U)$	$\backslash *$ and deflate the upper bound there * \backslash				

Theorem 2 (Soundness and completeness). Algorithm 1 (calling Algorithm 4) produces monotonic sequences L under- and U over-approximating V, and terminates.

Proof (Sketch). The crux is to show that U converges to V. We assume towards a contradiction, that there exists a state s with $\lim_{n\to\infty} U_n(s) - V(s) > 0$. Then there exists a nonempty set of states X where the difference between $\lim_{n\to\infty} U_n$ and V is maximal. If the upper bound of states in X depends on states outside of X, this yields a contradiction, because then the difference between upper bound and value would decrease in the next Bellman update. So X must be an EC where all states depend on each other. However, if that is the case, calling DEFLATE decreases the upper bound to something depending on the states outside of X, thus also yielding a contradiction. □

Summary of our approach:

- 1. We cannot collapse MECs, because we cannot collapse BECs with nonconstant values.
- 2. If we remove X (the sub-optimal actions of Minimizer) we can collapse MECs (now actually MSECs with constant values).
- 3. Since we know neither X nor SECs we gradually deflate SEC approximations.

4.4 Learning-based algorithm

Asynchronous value iteration selects in each round a subset $T \subseteq S$ of states and performs the Bellman update in that round only on T. Consequently, it may speed up computation if "important" states are selected. However, using the standard VI it is even more difficult to determine the current error bound. Moreover, if some states are not selected infinitely often the lower bound may not even converge.

In the setting of bounded value iteration, the current error bound is known for each state and thus convergence can easily be enforced. This gave rise to asynchronous VI, such as BRTDP (bounded real time dynamic programing) in the setting of stopping MDPs [MLG05], where the states are selected as those that appear on a simulation run. Very similar is the adaptation for general MDP [BCC⁺14]. In order to simulate a run, the transition probabilities determine how to resolve the probabilistic choice. In order to resolve the non-deterministic choice of Maximizer, the "most promising action" is taken, i.e., with the highest U. This choice is derived from a reinforcement algorithm called delayed Q-learning and ensures convergence while practically performing well [BCC⁺14].

In this section, we harvest our convergence results and BVI algorithm for SG, which allow us to trivially extend the asynchronous learning-based approach of BRTDP to SGs. On the one hand, the only difference to the MDP algorithm is how to resolve the choice for Minimizer. Since the situation is dual, we again pick the "most promising action", in this case with the lowest L. On the other hand, the only difference to Algorithm 1 calling Algorithm 4 is that the Bellman updates of U and L are performed on the states of the simulation run only, see lines 2-3 of Algorithm 5.

Al	$\mathbf{gorithm} \ 5 \ \mathrm{Update} \ \mathrm{procedure} \ \mathrm{for} \ \mathrm{the} \ \mathrm{learning}/\mathrm{BRTDP} \ \mathrm{version} \ \mathrm{of} \ \mathrm{BVI} \ \mathrm{on} \ \mathrm{SG}$
1:	procedure UPDATE(L : $S \rightarrow [0, 1], U : S \rightarrow [0, 1])$
2:	$\rho \leftarrow \text{path } \mathbf{s}_0, \mathbf{s}_1, \dots, \mathbf{s}_\ell$ of length $\ell \leq k$, obtained by simulation where the
	successor of s is s' with probability $\delta(s, a, s')$ and a is sampled randomly from
	$\arg \max_{a} U(s, a)$ and $\arg \min_{a} L(s, a)$ for $s \in S_{\Box}$ and $s \in S_{\bigcirc}$, respectively
3:	L, U get updated by Eq. (3) and (4) on states $s_{\ell}, s_{\ell-1}, \ldots, s_0$ * all $s \in \rho * \setminus$
4:	for $T \in FIND_{-}MSEC(L)$ do
5:	DEFLATE(T,U)

If $\mathbf{1}$ or $\mathbf{0}$ is reached in a simulation, we can terminate it. It can happen that the simulation cycles in an EC. To that end, we have a bound k on the maximum number of steps. The choice of k is discussed in detail in [BCC⁺14] and we use $2 \cdot |S|$ to guarantee the possibility of reaching sinks as well as exploring new states. If the simulation cycles in an EC, the subsequent call of DEFLATE ensures that next time there is a positive probability to exit this EC. Further details can be found in [KKKW18, Appendix A.4.].

5 Experimental results

We implemented both our algorithms as an extension of PRISM-games [CFK⁺13a], a branch of PRISM [KNP11] that allows for modelling SGs, utilizing previous work of [BCC⁺14,Ujm15] for MDP and SG with single-player ECs. We tested the implementation on the SGs from the PRISM-games case studies [gam] that have reachability properties and one additional model from [CKJ12] that was also used in [Ujm15]. We compared the results with both the explicit and the hybrid engine of PRISM-games, but since the models are small both of them performed similar and we only display the results of the hybrid engine in Table 1.

Furthermore we ran experiments on MDPs from the PRISM benchmark suite [KNP12]. We compared our results there to the hybrid and explicit engine of PRISM, the interval iteration implemented in PRISM [HM17], the hybrid engine of STORM [DJKV17b] and the BRTDP implementation of [BCC⁺14].

Recall that the aim of the paper is not to provide a faster VI algorithm, but rather the first guaranteed one. Consequently, the aim of the experiments is not to show any speed ups, but to experimentally estimate the overhead needed for computing the guarantees.

For information on the technical details of the experiments, all the models and the tables for the experiments on MDPs we refer to [KKKW18, Appendix B]. Note that although some of the SG models are parametrized they could only be scaled by manually changing the model file, which complicates extensive benchmarking.

Although our approaches compute the additional upper bound to give the convergence guarantees, for each of the experiments one of our algorithms performed similar to PRISM-games. Table 1 shows this result for three of the four SG models in the benchmarking set. On the fourth model, PRISM's precomputations already solve the problem and hence it cannot be used to compare the approaches. For completeness, the results are displayed in [KKKW18, Appendix B.5].

Table 1: Experimental results for the experiments on SGs. The left two columns denote the model and the given parameters, if present. Columns 3 to 5 display the verification time in seconds for each of the solvers, namely PRISM-games (referred as PRISM), our BVI algorithm (BVI) and our learning-based algorithm (BRTDP). The next two columns compare the number of states that BRTDP explored (#States_B) to the total number of states in the model. The rightmost column shows the number of MSECs in the model.

Model	Parameters	PRISM	BVI	BRTDP	#States_B	#States	#MSECs
mdem	prop=1	8	8	17	767	62,245	1
musm	prop=2	4	4	29	407	62,245	1
cdmsn		2	2	3	1,212	1,240	1
cloud	N=5	3	7	15	1,302	8,842	4,421
	N=6	6	59	4	570	34,954	17,477

Whenever there are few MSECs, as in mdsm and cdmsn, BVI performs like PRISM-games, because only little time is used for deflating. Apparently the additional upper bound computation takes very little time in comparison to the other tasks (e.g. parsing, generating the model, pre-computation) and does not slow down the verification significantly. For cloud, BVI is slower than PRISMgames, because there are thousands of MSECs and deflating them takes over 80% of the time. This comes from the fact that we need to compute the expensive end component decomposition for each deflating step. BRTDP performs well for cloud, because in this model, as well as generally often if there are many MECs [BCC⁺14], only a small part of the state space is relevant for convergence. For the other models, BRTDP is slower than the deterministic approaches, because the models are so small that it is faster to first construct them completely than to explore them by simulation.

Our more extensive experiments on MDPs compare the guaranteed approaches based on collapsing (i.e. learning-based from [BCC⁺14] and deterministic from [HM17]) to our guaranteed approaches based on deflating (so BRTDP and BVI). Since both learning-based approaches as well as both deterministic approaches perform similarly (see Table 2 in [KKKW18, Appendix B]), we conclude that collapsing and deflating are both useful for practical purposes, while the latter is also applicable to SGs. Furthermore we compared the usual unguaranteed value iteration of PRISM's explicit engine to BVI and saw that our guaranteed approach did not take significantly more time in most cases. This strengthens the point that the overhead for the computation of the guarantees is negligible

6 Conclusions

We have provided the first stopping criterion for value iteration on simple stochastic games and an anytime algorithm with bounds on the current error (guarantees on the precision of the result). The main technical challenge was that states in end components in SG can have different values, in contrast to the case of MDP. We have shown that collapsing is in general not possible, but we utilized the analysis to obtain the procedure of *deflating*, a solution on the original graph. Besides, whenever a SEC is identified for sure it can be collapsed and the two techniques of collapsing and deflating can thus be combined.

The experiments indicate that the price to pay for the overhead to compute the error bound is often negligible. For each of the available models, at least one of our two implementations has performed similar to or better than the standard approach that yields no guarantees. Further, the obtained guarantees open the door to (e.g. learning-based) heuristics which treat only a part of the state space and can thus potentially lead to huge improvements. Surprisingly, already our straightforward adaptation of such an algorithm for MDP to SG yields interesting results, palliating the overhead of our non-learning method, despite the most naive implementation of deflating. Future work could reveal whether other heuristics or more efficient implementation can lead to huge savings as in the case of MDP [BCC⁺14].

References

- [ACD⁺17] Pranav Ashok, Krishnendu Chatterjee, Przemyslaw Daca, Jan Kretínský, and Tobias Meggendorfer. Value iteration for long-run average reward in markov decision processes. In CAV, pages 201–221, 2017.
- [AM09] Daniel Andersson and Peter Bro Miltersen. The complexity of solving stochastic games on graphs. In *ISAAC*, pages 112–121, 2009.
- [AY17] Gürdal Arslan and Serdar Yüksel. Decentralized q-learning for stochastic teams and games. *IEEE Trans. Automat. Contr.*, 62(4):1545–1558, 2017.
- [BBS08] Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Trans. Systems, Man,* and Cybernetics, Part C, 38(2):156–172, 2008.
- [BCC⁺14] Tomás Brázdil, Krishnendu Chatterjee, Martin Chmelik, Vojtech Forejt, Jan Kretínský, Marta Z. Kwiatkowska, David Parker, and Mateusz Ujma. Verification of Markov decision processes using learning algorithms. In ATVA, pages 98–114. Springer, 2014.
- [BK08] Christel Baier and Joost-Pieter Katoen. Principles of model checking, 2008.
- [BKL⁺17] Christel Baier, Joachim Klein, Linda Leuschner, David Parker, and Sascha Wunderlich. Ensuring the reliability of your model checker: Interval iteration for markov decision processes. In CAV, pages 160–180, 2017.
- [BT00] Ronen I. Brafman and Moshe Tennenholtz. A near-optimal polynomial time algorithm for learning in certain classes of stochastic games. Artif. Intell., 121(1-2):31–47, 2000.
- [CF11] Krishnendu Chatterjee and Nathanaël Fijalkow. A reduction from parity games to simple stochastic games. In *GandALF*, pages 74–86, 2011.
- [CFK⁺13a] T. Chen, V. Forejt, M. Kwiatkowska, D. Parker, and A. Simaitis. PRISMgames: A model checker for stochastic multi-player games. In N. Piterman and S. Smolka, editors, Proc. 19th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'13), volume 7795 of LNCS, pages 185–191. Springer, 2013.
- [CFK⁺13b] Taolue Chen, Vojtech Forejt, Marta Z. Kwiatkowska, David Parker, and Aistis Simaitis. Automatic verification of competitive stochastic systems. *Formal Methods in System Design*, 43(1):61–92, 2013.
- [CH08] Krishnendu Chatterjee and Thomas A Henzinger. Value iteration. In 25 Years of Model Checking, pages 107–138. Springer, 2008.
- [CHJR10] Krishnendu Chatterjee, Thomas A. Henzinger, Barbara Jobstmann, and Arjun Radhakrishna. Gist: A solver for probabilistic games. In CAV, pages 665–669, 2010.
- [CKJ12] Radu Calinescu, Shinji Kikuchi, and Kenneth Johnson. Compositional Reverification of Probabilistic Safety Properties for Large-Scale Complex IT Systems, pages 303–329. Springer Berlin Heidelberg, Berlin, Heidelberg, 2012.
- [CKLB11] Chih-Hong Cheng, Alois Knoll, Michael Luttenberger, and Christian Buckl. GAVS+: an open platform for the research of algorithmic game solving. In ETAPS, pages 258–261, 2011.
- [CKSW13] Taolue Chen, Marta Z. Kwiatkowska, Aistis Simaitis, and Clemens Wiltsche. Synthesis for multi-objective stochastic games: An application to autonomous urban driving. In QEST, pages 322–337, 2013.

- [CMG14] Javier Cámara, Gabriel A. Moreno, and David Garlan. Stochastic game analysis and latency awareness for proactive self-adaptation. In 9th International Symposium on Software Engineering for Adaptive and Self-Managing Systems, SEAMS 2014, Proceedings, Hyderabad, India, June 2-3, 2014, pages 155–164, 2014.
- [Con92] Anne Condon. The complexity of stochastic games. Information and Computation, 96(2):203–224, 1992.
- [CY95] Costas Courcoubetis and Mihalis Yannakakis. The complexity of probabilistic verification. *Journal of the ACM*, 42(4):857–907, July 1995.
- [DJKV17a] Christian Dehnert, Sebastian Junges, Joost-Pieter Katoen, and Matthias Volk. A storm is coming: A modern probabilistic model checker. In CAV, pages 592–600, 2017.
- [DJKV17b] Christian Dehnert, Sebastian Junges, Joost-Pieter Katoen, and Matthias Volk. A storm is coming: A modern probabilistic model checker. CoRR, abs/1702.04311, 2017.
- [gam] PRISM-games case studies. prismmodelchecker.org/games/casestudies.php. Accessed: 2017-09-18.
- [HK66] A. J. Hoffman and R. M. Karp. On nonterminating stochastic games. Management Science, 12(5):359–370, 1966.
- [HM17] Serge Haddad and Benjamin Monmege. Interval iteration algorithm for mdps and imdps. *Theoretical Computer Science*, 2017.
- [KKKW18] Edon Kelmendi, Julia Krämer, Jan Křetínský, and Maximilian Weininger. Value iteration for simple stochastic games: Stopping criterion and learning algorithm. Technical Report abs/1804.04901, arXiv.org, 2018.
- [KKNP10] Mark Kattenbelt, Marta Z. Kwiatkowska, Gethin Norman, and David Parker. A game-based abstraction-refinement framework for markov decision processes. Formal Methods in System Design, 36(3):246–280, 2010.
- [KM17] Jan Kretínský and Tobias Meggendorfer. Efficient strategy iteration for mean payoff in markov decision processes. In ATVA, pages 380–399, 2017.
- [KNP11] M. Kwiatkowska, G. Norman, and D. Parker. PRISM 4.0: Verification of probabilistic real-time systems. In CAV, volume 6806 of LNCS, pages 585–591. Springer, 2011.
- [KNP12] M. Kwiatkowska, G. Norman, and D. Parker. The prism benchmark suite. 9th International Conference on Quantitative Evaluation of Systems (QEST12),pages 203204. IEEE, 2012.
- [LaV00] Steven M. LaValle. Robot motion planning: A game-theoretic foundation. Algorithmica, 26(3-4):430–465, 2000.
- [LL08] Jianwei Li and Weiyi Liu. A novel heuristic q-learning algorithm for solving stochastic games. In *IJCNN*, pages 1135–1144, 2008.
- [Mar75] Donald A Martin. Borel determinacy. Annals of Mathematics, pages 363– 371, 1975.
- [MLG05] H. Brendan Mcmahan, Maxim Likhachev, and Geoffrey J. Gordon. Bounded real-time dynamic programming: Rtdp with monotone upper bounds and performance guarantees. In *In ICML05*, pages 569–576, 2005.
- [Put14] Martin L. Puterman. Markov decision processes: Discrete stochastic dynamic programming. John Wiley & Sons, 2014.
- [SK16] María Svorenová and Marta Kwiatkowska. Quantitative verification and strategy synthesis for stochastic games. Eur. J. Control, 30:15–30, 2016.
- [TT16] Alain Tcheukam and Hamidou Tembine. One swarm per queen: A particle swarm learning for stochastic games. In *SASO*, pages 144–145, 2016.

- [Ujm15] Mateusz Ujma. On Verication and Controller Synthesis for Probabilistic Systems at Runtime. PhD thesis, Wolfson College, Oxford, 2015.
- [WT16] Min Wen and Ufuk Topcu. Probably approximately correct learning in stochastic games with temporal logic specifications. In *IJCAI*, pages 3630–3636, 2016.