

Optimistic and Topological Value Iteration for Simple Stochastic Games

Muqsit Azeem^{}, Alexandros Evangelidis^{}, Jan Křetínský^{}, Alexander Slivinskiy^{}, and Maximilian Weininger^{}

Technical University of Munich, Munich, Germany
`firstname.lastname@tum.de`

Abstract. While value iteration (VI) is a standard solution approach to simple stochastic games (SSGs), it suffered from the lack of a stopping criterion. Recently, several solutions have appeared, among them also “optimistic” VI (OVI). However, OVI is applicable only to one-player SSGs with no end components. We lift these two assumptions, making it available to *general SSGs*. Further, we utilize the idea in the context of topological VI, where we provide an efficient *precise* solution. In order to compare the new algorithms with the state of the art, we use not only the standard benchmarks, but we also design a *random generator* of SSGs, which can be biased towards various types of models, aiding in understanding the advantages of different algorithms on SSGs.

1 Introduction

Stochastic games (SGs) are a standard model for decision making in the presence of adversary and uncertainty, by combining two (opposing) non-determinisms with stochastic dynamics. Thus, they extend both Markov decision processes (MDPs), the standard model for sequential decision making and probabilistic verification, and 2-player graph games, the standard model for reactive synthesis. *Simple stochastic games* (SSGs) [12] form an important special case where the goal is to reach a given state. In technical terms, an SSG is a zero-sum two-player turn-based 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. The interest in SSGs stems from two sources. Firstly, solving an SSG is polynomial-time equivalent to solving perfect information Shapley, Everett and Gillette games [1] and further important problems can be reduced to SSGs, for instance parity games, mean-payoff games, discounted-payoff games and their stochastic extensions [6]; yet, the complexity of solving SSGs remains a long-standing open question, known to be in $\mathbf{UP} \cap \mathbf{coUP}$ [23], but with polynomial-time algorithm staying elusive. Secondly, the problem is practically relevant in verification and synthesis in stochastic environments, with many applications, e.g., [27, 9, 10, 5], surveyed in detail in [31]. Consequently, heuristics improving performance of the algorithms for solving SSGs are also practically relevant.

Algorithms used to (approximately) solve SSGs can be divided into several classes, most notably quadratic programming (QP) and dynamic programming, the latter

comprising strategy iteration (SI) and value iteration (VI). For their practical comparison, see the recent [24].

On the one hand, only when exact solutions are required, SI is mostly used. It provides a sequence of improving strategies and, accompanied by evaluation of Markov chains via systems of linear equations, can yield the precise result. On the other hand, approximate solutions (with a certain imprecision) are faster to compute and often sufficient. For this reason, VI is the technique used in practice the most, e.g., in PRISM-games [25], although not necessarily always the best. It gradually approximates (from below) the optimal probability to reach the target from each state. Interestingly, until very recently no means were known to determine the current precision, and so standard implementations terminating whenever no significant improvements occur can be arbitrarily wrong [20]. More surprisingly, this was even the case for MDPs, i.e., SSGs with a single player.

In 2014 [4,20], the first stopping criterion for MDPs was given, quantifying precision of the current approximation by providing also a sequence converging to the optimal probabilities from above. The difficulty to obtain such converging upper bound arises from cyclic dependencies of the optimal probabilities in so-called end components (ECs). For instance, an action surely self-looping on a state trivializes the equations, stating only that the probability in this state is simply equal to itself, yielding an infinity of solutions, not just the optimal one. This issue has been solved for MDPs [4,20] by “collapsing” these ECs into single states with no loops, which corresponds to identifying cyclically dependent variables into a single one.

In 2018 [18], the idea was finally extended to SSGs, giving rise to *bounded value iteration* (BVI) with the first stopping criterion for SSGs. Note that the MDP solution could not be directly used since the analog of ECs in SSGs is more complex: different states in an EC in an SSG can have different optimal probabilities and thus cannot be merged. Instead, “deflating” manipulates the values in smaller and dynamically changing “simple ECs”.

Since the first VI stopping criterion was given for MDPs, several alternatives have been proposed, most notably *sound VI* (SVI) [30] and *optimistic VI* (OVI) [21]. However, the termination proofs of both require the MDP to contain no ECs. They achieve this by collapsing ECs, which is not applicable to SSGs.

Our contribution. In this paper, we extend the idea of OVI in two ways so that we obtain algorithms for SSGs.

First algorithm The idea of OVI [21] is to run VI (converging from below) until changes are small, then to guess slightly larger values and check whether they form an upper bound. If not, the process continues. To overcome the requirement that there is no EC, we complement the procedure of [21] with the deflating of [18]. However, to ensure monotonicity of the Bellman operator, the so-called “simple” ECs must be computed differently from [18]. While the rest of the proof is analogous to [21], we try to make it simpler and more elegant by separating the core idea from the practical improvements. As a result, we obtain an OVI algorithm for SSGs.

Second algorithm We consider the classic “topological” optimization of VI [16], where the system is analysed per strongly connected component (SCC) in the bottom-up order. While such decomposition often leads to savings in runtime and memory, also when expected accumulated rewards are considered [3], the impreci-

sions from lower SCCs propagate to the upper ones, yielding the method useless whenever the system is too deep (with as few SCCs in a row as 20) even for Markov chains, see Example 1. We fix this issue by precise *and* fast computations in each SCC as follows. First, we quickly obtain an approximate solution by VI, then we optimistically guess the solution, but in contrast to OVI which guesses values, we guess optimal strategies, which turns out to require orders of magnitude fewer guesses. If the guess is not correct, a step of SI can be cheaply performed. This version of OVI can thus be also seen as a possible warm start for SI.

Comparison and model generation We compare the resulting approaches to BVI and a more recent SSG solution called “widest path” [28] (WP). While there is no clear winner, we provide insights as to which algorithm to use in different settings. As noticed already in [24], the performance of SSG algorithms is extremely sensitive to the structure of the models. Unfortunately, there are too few realistic case studies and thus a very limited number of model structures. Consequently, in order to be able to experimentally compare our algorithms in a reasonable way, we propose an approach for random SSG generation. While we prove that our approach can generate every SSG, it skews towards certain types of models. Hence we provide means for the user to skew towards model structures that they are interested in, e.g., increasing or decreasing the number of SCCs. This helps to find out which algorithms are sensitive to which model parameters, e.g., amount of SCCs. While this is only the first step towards filling this gap of random SSG (and MDP) generation, we hope to encourage more research on the topic through this effort.

Our contribution can be summarized as follows:

- We design an extension of OVI to SSGs. As a side effect, we extend OVI on MDPs, lifting the requirement of no ECs (Section 3).
- We extend the landscape by providing an efficient VI-based approach for precise solutions, using the OVI idea on strategies, rather than values (Section 4).
- We provide and evaluate a random generator of SSGs, which can be biased towards various types of models (Section 5).
- We compare the resulting methods to the state of the art (BVI, WP, SI) experimentally (Section 6).

Related work. Closest to our work, in the case of SSGs, is the work of [18] where the first stopping criterion for VI was given. It extends both normal BVI [20] and its learning-based counterpart [4] from MDPs by incorporating the so-called deflating procedure as part of their computation. Recently, another BVI variant for SSGs was proposed which introduces a global propagation of upper bounds [28]. Also, the simpler case of an SSG with one-player ECs is discussed in [33].

In general, the tools which are available for solving SGs are limited. PRISM-games [25] implements the standard VI algorithms, and it also considers other objectives apart from reachability, such as mean-payoff and ratio reward. Further, GAVS+ [11] is an algorithmic game solver with support for solving SSGs, and GIST [8] allows for the qualitative verification of SGs.

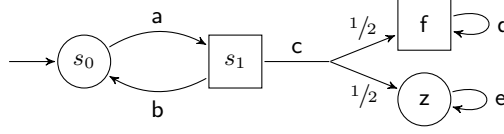


Fig. 1: An example of an SG with $S = \{s_0, s_1, f, z\}$, $S_\square = \{s_1, f\}$, $S_\circ = \{s_0, z\}$, the initial state s_0 and set of actions $A = \{a, b, c, d, e\}$; $\text{Av}(s_0) = \{a\}$ with $\delta(s_0, a)(s_1) = 1$; $\text{Av}(s_1) = \{b, c\}$ with $\delta(s_1, b)(s_0) = 1$ and $\delta(s_1, c)(f) = \delta(s_1, c)(z) = \frac{1}{2}$. For actions with only one successor, we do not depict the transition probability 1.

2 Preliminaries

2.1 Simple stochastic games

A probability distribution on a finite set X is a mapping $\delta : X \rightarrow [0, 1]$, such that $\sum_{x \in X} \delta(x) = 1$. The set of all probability distributions on X is denoted by $\text{Dist}(X)$.

Definition 1 (Stochastic game (SG), e.g., [13]). A stochastic turn-based two-player game is defined by a tuple $\mathcal{G} = \langle S, S_\square, S_\circ, s_0, A, \text{Av}, \delta \rangle$ where S is a finite set of states partitioned into a set of Minimizer (S_\circ) and Maximizer (S_\square) states, respectively. $s_0 \in S$ is the initial state. A is a finite set of actions. $\text{Av} : S \rightarrow 2^A$ assigns to every state a set of available actions. Finally, $\delta : S \times A \rightarrow \text{Dist}(S)$ is the transition function.

Note that a Markov decision process (MDP) is a special case of an SG where either $S_\circ = \emptyset$ or $S_\square = \emptyset$ and a Markov chain is a special case of an MDP where in each state there is only one available action.

Without loss of generality, we assume SGs to be non-blocking, i.e., for all $s \in S$: $\text{Av}(s) \neq \emptyset$. For convenience, we use the following notation: Given a state $s \in S$ and an action $a \in \text{Av}(s)$, the set of successor states is denoted as $\text{Post}(s, a) := \{s' \mid \delta(s, a, s') > 0\}$. For a set of states $T \subseteq S$, we use $T_\square = T \cap S_\square$ to denote all Maximizer states in T , and dually for Minimizer. Figure 1 shows an example SG.

Semantics: paths, strategies and the value. Formally, an *infinite path* ρ is defined as $\rho = s_0 a_0 s_1 a_1 \dots \in (S \times A)^\omega$, such that for every $i \in \mathbb{N}$, $a_i \in \text{Av}(s_i)$ and $s_{i+1} \in \text{Post}(s_i, a_i)$. The set of all paths in an SG \mathcal{G} is denoted as $\text{Paths}_{\mathcal{G}}$. A finite path is a prefix of an infinite path ending in a state s .

A Maximizer *strategy* is a function $\sigma : S_\square \rightarrow A$ such that $\sigma(s) \in \text{Av}(s)$ for all s ; Minimizer strategies τ are defined analogously. We restrict attention to *memoryless deterministic* strategies, because they are sufficient for the objective we consider [12]. By fixing both players' choices according to a pair of strategies (σ, τ) , we turn an SG \mathcal{G} into a Markov chain $\mathcal{G}^{(\sigma, \tau)}$ with state space S and the transition function $\delta^{\sigma, \tau}(s, s') = \delta(s, \sigma(s), s')$ for Maximizer states s and dually for Minimizer with σ replaced by τ . Given a state s , the Markov chain $\mathcal{G}^{(\sigma, \tau)}$ induces a unique probability distribution $\mathcal{P}_s^{\sigma, \tau}$ over the set of all infinite paths [2, Sec. 10.1].

Since we consider SSGs, we complement an SG with a set of goal states $F \subseteq S$ and formalize the objective of reaching F , as follows: we denote as $\Diamond F := \{\rho \mid \rho = s_0 a_0 s_1 a_1 \dots \in \text{Paths}_{\mathcal{G}} \wedge \exists i \in \mathbb{N}. s_i \in F\}$ the (measurable) set of all paths which eventually reach F . We are interested in the *value* of every state s , i.e., the probability

that s reaches a goal state if both players play optimally. Formally, for each $s \in \mathbf{S}$, its value is defined as

$$V(s) := \sup_{\sigma} \inf_{\tau} \mathcal{P}_s^{\sigma, \tau}(\Diamond F) = \inf_{\tau} \sup_{\sigma} \mathcal{P}_s^{\sigma, \tau}(\Diamond F), \quad (1)$$

where the equality follows from [12]. We use $V : \mathbf{S} \rightarrow \mathbb{R}$ to denote the function that maps every $s \in \mathbf{S}$ to its value. When comparing functions $f_1, f_2 : \mathbf{S} \rightarrow \mathbb{R}$, we use point-wise comparison, i.e., $f_1 \leq f_2$ if and only if for all $s \in \mathbf{S} : f_1(s) \leq f_2(s)$.

2.2 Value iteration and bounded value iteration

To compute the value function V for an SSG, the following partitioning of the state space is useful: firstly the goal states \mathbf{F} , secondly the set of *sink states* that do not have a path to the target $\mathbf{Z} = \{s \in \mathbf{S} \mid \nexists \rho = s_0 a_0 s_1 a_1 \dots \in \text{Paths}_G : s_0 = s \wedge \rho \in \Diamond F\}$, and finally the remaining states $\mathbf{S}^?$. For \mathbf{F} and \mathbf{Z} (which can be easily identified by graph-search algorithms), the value is trivially 1 respectively 0. Thus, the computation only has to focus on $\mathbf{S}^?$.

The well-known approach of value iteration leverages the fact that V is the least fixpoint of the *Bellman equations*, cf. [7]:

$$V(s) = \begin{cases} 1 & \text{if } s \in \mathbf{F} \\ 0 & \text{if } s \in \mathbf{Z} \\ \max_{a \in \text{Av}(s)} \left(\sum_{s' \in \mathbf{S}} \delta(s, a, s') \cdot V(s') \right) & \text{if } s \in \mathbf{S}_{\square}^? \\ \min_{a \in \text{Av}(s)} \left(\sum_{s' \in \mathbf{S}} \delta(s, a, s') \cdot V(s') \right) & \text{if } s \in \mathbf{S}_{\circ}^? \end{cases} \quad (2)$$

Now we define¹ the Bellman operator $\mathcal{B} : (\mathbf{S} \rightarrow \mathbb{R}) \rightarrow (\mathbf{S} \rightarrow \mathbb{R})$:

$$\mathcal{B}(f)(s) = \begin{cases} \max_{a \in \text{Av}(s)} \left(\sum_{s' \in \mathbf{S}} \delta(s, a, s') \cdot f(s') \right) & \text{if } s \in \mathbf{S}_{\square} \\ \min_{a \in \text{Av}(s)} \left(\sum_{s' \in \mathbf{S}} \delta(s, a, s') \cdot f(s') \right) & \text{if } s \in \mathbf{S}_{\circ} \end{cases} \quad (3)$$

Value iteration starts with the under-approximation

$$L_0(s) = \begin{cases} 1 & \text{if } s \in \mathbf{F} \\ 0 & \text{otherwise} \end{cases}$$

and repeatedly applies the Bellman operator. Since the value is the least fixpoint of the Bellman equations and $L_0 \leq V$ is lower than the value, this converges to the value in the limit [7] (formally $\lim_{i \rightarrow \infty} \mathcal{B}^i(L_0) = V$).

While this approach is often fast in practice, it has the drawback that it is not possible to know the current difference between $\mathcal{B}^i(L_0)$ and V for any given i . To address this, one can employ *bounded value iteration* (BVI, also known as interval iteration [4,20,18]) It additionally starts from an over-approximation U_0 , with $U_0(s) = 1$ for all $s \in \mathbf{S}$. However, applying the Bellman operator to this upper estimate might not converge to the value, but to some greater fixpoint instead, see [18, Section 3] for an example. The core of the problem are so called *end components*.

¹ In the definition of \mathcal{B} , we omit the technical detail that for goal states $s \in \mathbf{F}$, the value has to remain 1. Equivalently, one can assume that all goal states are absorbing, i.e., only have self looping actions.

Definition 2 (End component (EC)). A set of states T with $\emptyset \neq T \subseteq S$ is an end component if and only if there exists a set of actions $\emptyset \neq B \subseteq \bigcup_{s \in T} \text{Av}(s)$ such that:

1. for each $s \in T$, $a \in B \cap \text{Av}(s)$ we have $\text{Post}(s, a) \subseteq T$.
2. for each $s, s' \in T$ there exists a finite path $w = sa_0 \dots a_n s' \in (T \times B)^* \times T$.

An end component T is a maximal end component (MEC) if there is no other EC T' such that $T \subsetneq T'$.

Intuitively, ECs can be problematic, because the over-approximation U is higher in the EC than the value. Thus, Maximizer prefers staying in the EC and keeping the illusion of achieving the high U ; it is an illusion, because staying will never reach a target, and Maximizer actually has to use some exit of the EC. The solution proposed in [18] explicitly identifies these situations and forces all states in the EC to decrease their U by making it depend on the best exit of the EC. This operation is called *deflating*, to evoke the impression of releasing the pressure in an EC that is bloated by having too high estimates. To define deflating more formally, we need two definitions from [18]:

Definition 3 (Best exit). Given a set of states $T \subseteq S$ and a function $f : S \rightarrow \mathbb{R}$, the best exit according to f from T is defined as:

$$\text{bexit}_f(T) = \max_{\substack{s \in T, a \in \text{Av}(s) \\ \text{Post}(s, a) \not\subseteq T}} \left(\sum_{s' \in S} \delta(s, a, s') \cdot f(s, a) \right),$$

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

Definition 4 (Simple end component (SEC)). An EC T is a simple end component (SEC) if for all $s \in T$, $V(s) = \text{bexit}_V(T)$

In SSGs, states in an EC can have different values. Thus, it is necessary to find the SECs. In these simple sub-parts of the EC all states have the same value, namely that of the best exit. By setting the over-approximation to $\text{bexit}_U(T)$ for each SEC T (additionally to applying \mathcal{B}), we ensure that it converges to the value [18]. As a final complication, computing SECs is difficult, since they depend on the value V that we want to compute. The solution of [18] is to use the current under-approximation L to guess which states form a SEC and as L converges to V in the limit, eventually we guess correctly.

Thus, we can augment the Bellman operator with additional deflating and define an operator $\mathcal{B}_L^D : (S \rightarrow \mathbb{R}) \rightarrow (S \rightarrow \mathbb{R})$. Note that it depends on an L to guess the SECs. Given a function U , it proceeds as follows:

- Apply a Bellman update $\mathcal{B}(U)$.
- Guess the SECs according to L by using [18, Algorithm 2].
- For each SEC T and all states $s \in T$, set $U(s) = \min(U(s), \text{bexit}_U(T))$. The min is only to ensure monotonicity.

In summary, BVI computes two sequences: the sequence of lower bounds $L_i = \mathcal{B}^i(L)$ for $i \in \mathbb{N}$ and an additional sequence of upper bounds $U_i = (\mathcal{B}_L^D)^i(U)$. Note that for the i -th application of \mathcal{B}_L^D , it uses the current lower bound L_i . Both sequences converge to the value V in the limit [18, Theorem 2]. This allows to terminate the algorithm when the difference between the lower and upper bound is less than a pre-defined precision ε and obtain an ε -approximation of the value.

3 Optimistic Value Iteration

The idea of optimistic value iteration (OVI, [21]) is to leverage the fact that classic VI (only from below) typically converges quickly to the correct value. Indeed, the following “naive” stopping criterion results in an approximation that is ε -close in all available realistic case studies: stop when for all $s \in \mathcal{S}$ applying the Bellman update does not result in a big difference, i.e. $\text{diff}(\mathbf{L}(s), \mathcal{B}(\mathbf{L})(s)) < \varepsilon$, where we use $\text{diff}(\text{old}, \text{new}) = \text{new} - \text{old}$ to denote the absolute difference between two numbers². However, the naive stopping criterion can also terminate early when the estimate still is arbitrarily wrong [20].

OVI first performs classic VI with the naive stopping criterion, optimistically hoping that it will terminate close to the value. Additionally, it uses a *verification phase*, where it checks whether the result of VI was indeed correct. If it was, OVI terminates with the guarantee that we are ε -close to the value. Otherwise, if the result of VI cannot be verified, OVI continues VI with a higher precision ε' . By repeating this, at some point ε' is so small that when VI terminates, OVI can verify that the result is ε -precise.

Our version of OVI for SSGs is given in Algorithm 1. Lines 2-3 are the classic VI, Lines 4-9 the verification phase. Concretely, in the verification phase we first guess a candidate upper bound \mathbf{U} (Line 4), so that the difference between \mathbf{L} and \mathbf{U} is small enough that, if \mathbf{U} indeed is an upper bound, we could terminate. Formally,

for all $s \in \mathcal{S}$, $\mathbf{U}(s) = \text{diff}^+(L(s))$, where $\text{diff}_\varepsilon^+(x) = \begin{cases} 0 & \text{if } x = 0 \\ x + \varepsilon & \text{otherwise} \end{cases}$ for absolute

difference³. Then we apply the Bellman operator once (Line 6) and check whether $\mathcal{B}_L^D(\mathbf{U}) \leq \mathbf{U}$ (Line 7). If that holds, we know (by arguments from lattice theory) that $\mathbf{V} \leq \mathbf{U}$, i.e. that \mathbf{U} is a valid upper bound on the value. Thus, since \mathbf{L} and \mathbf{U} are ε -close to each other and $\mathbf{L} \leq \mathbf{V} \leq \mathbf{U}$, we return an ε -approximation of the value (Line 8). **The key difference** between the original algorithm for MDPs and the extension to SSGs is that we do not use \mathcal{B} in Line 6 any more, but the Bellman operator with additional deflating \mathcal{B}^D . On MDPs, the termination of OVI relied on the assumption that there were no ECs. This is justified, since in MDPs one can remove the ECs by “collapsing” them beforehand, cf. [4,20]. On SSGs, collapsing is not possible [18], which is why we need the new operator.

We have addressed the case that the guessed \mathbf{U} can indeed be verified as an upper bound. In the other case where we are not (yet) able to verify it, Algorithm 1 continues applying \mathcal{B}_L^D for a finite number of times (we chose $\frac{1}{\varepsilon'}$, Line 5). If for all iterations we cannot verify \mathbf{U} as an upper bound, the precision ε' for the naive stopping criterion is increased (we chose $\frac{\varepsilon'}{2}$) and we start over (Line 10).

Theorem 1. *Given an SSG \mathcal{G} and a lower bound $\mathbf{L}_0 \leq \mathbf{V}$, $\text{OVI}(\mathcal{G}, \mathbf{L}_0, \varepsilon, \varepsilon)$ terminates and returns (\mathbf{L}, \mathbf{U}) such that $\mathbf{L} \leq \mathbf{V} \leq \mathbf{U}$ and $\text{diff}(\mathbf{U}(s), \mathbf{L}(s)) \leq \varepsilon$ for all $s \in \mathcal{S}$.*

Our formulation of Algorithm 1 is simpler than [21, Algorithm 2], since we include only the key parts that are necessary for the proof of Theorem 1 (provided

² One can also use the relative difference, i.e. $\text{diff}(\text{old}, \text{new}) = \frac{\text{new} - \text{old}}{\text{new}}$.

³ $\text{diff}_\varepsilon^+(x) = x * (1 + \varepsilon)$ for relative difference.

Algorithm 1 Optimistic value iteration for SSGs.

Input: SSG \mathcal{G} , lower bound $L \leq V$, precision $\varepsilon > 0$ and naive precision $\varepsilon' > 0$
Output: (L, U) such that $L \leq V \leq U$ and $\text{diff}(U(s), L(s)) \leq \varepsilon$ for all $s \in S$

```

1: procedure OVI( $\mathcal{G}, L, \varepsilon, \varepsilon'$ )
   $\triangleright$  Classic VI with naive convergence criterion
2:   while for some state  $s \in S : \text{diff}(L(s), \mathcal{B}(L)(s)) > \varepsilon'$  do
3:      $L \leftarrow \mathcal{B}(L)$ 
   $\triangleright$  Verification phase
4:    $U \leftarrow \{s \mapsto \text{diff}_\varepsilon^+(L(s)) \mid s \in S\}$   $\triangleright$  Guess candidate upper bound
5:   for  $\frac{1}{\varepsilon'}$  times do
6:      $U' \leftarrow \mathcal{B}_L^D(U)$ 
7:     if  $U' \leq U$  then
8:       return  $(L, U)$   $\triangleright$  Found inductive upper bound
9:     For all  $s \in S^? : U(s) \leftarrow \min(U(s), U'(s))$   $\triangleright$  Ensure monotonicity
10:  return OVI( $\mathcal{G}, L, \varepsilon, \frac{\varepsilon'}{2}$ )  $\triangleright$  Try again with more precision

```

in Appendix A). Below we comment on three ways in which our algorithm can be changed, following the ideas of [21, Algorithm 2]. All these changes are not necessary for correctness or termination, but they can practically improve the algorithm.

1. We can include a check $\mathcal{B}^D(U) \geq U$. It allows to detect whether $U \leq V$, i.e. U actually is a lower bound on the value. In that case, one can immediately terminate the verification phase and use U as the new L . We include this improvement in our implementation, and it is used in almost every unsuccessful verification phase.
2. The original version continues to update the lower bound during the verification phase. This is used for an additional breaking condition if the lower bound crossed the upper bound in some state. For clarity of presentation, we chose to separate concerns and only update the upper bound in the verification phase. This improvement never made a significant difference in our experiments.
3. The original version used Gauß-Seidel VI, cf. [21, Section 3.1], for both the lower and the upper bound. Our implementation allows the user to select whether to use classic or Gauß-Seidel VI.

4 Precise topological value iteration

Topological value iteration (TVI, [16]) is a variant of VI that does not solve the whole game at once, but rather proceeds piece by piece. This can speed up convergence and help with memory issues. Concretely, it uses the insight that the strongly connected components (SCCs) of an SSG always form a directed acyclic graph. Thus, one can first solve the bottom SCCs, i.e. the last in the topological ordering, and then proceed backwards one SCC by the next, relying on the results of the already computed successor SCCs. This idea is not restricted to VI algorithms, but can also be used for other solutions methods like strategy iteration (SI) and quadratic programming [24].

The evaluation of [24] showed that this can be quite useful in some cases, but also much slower in other, possibly even running into time outs on models where

Algorithm 2 Precise topological value iteration

Input: SSG \mathcal{G}
Output: The precise value V for all states in \mathcal{G}

```

1: procedure PTVI( $\mathcal{G}$ )
2:   for every SCC  $T$  in reverse topological ordering do
3:     Select arbitrary  $\varepsilon$ 
4:      $L, U \leftarrow$  computed by some VI-algorithm with precision  $\varepsilon$ 
5:     Compute strategies  $\sigma, \tau$  which are optimal according to  $L$  and  $U$ 
6:     Precisely compute the value  $V_{\mathcal{G}^{(\sigma, \tau)}}(s)$  of  $T$  in the Markov chain  $\mathcal{G}^{(\sigma, \tau)}$ 
7:     if For all  $s \in T$  :  $\begin{cases} \sigma(s) \in \arg \max_{a \in A(s)} V_{\mathcal{G}^{(\sigma, \tau)}}(s, a) & \text{if } s \in S_{\square} \\ \tau(s) \in \arg \min_{a \in A(s)} V_{\mathcal{G}^{(\sigma, \tau)}}(s, a) & \text{if } s \in S_{\circ} \end{cases}$  then
8:       Return  $V_{\mathcal{G}^{(\sigma, \tau)}}$  as value for  $T$ .
9:     else
10:      Apply strategy iteration, using  $\sigma$  or  $\tau$  as initial strategy.
```

the normal algorithms succeed. The reason for this is a complex problem that did not occur in the proof of correctness, as it is related to machine precision: SCCs are not solved precisely, but only with ε -precision. That means that SCCs which are considered later in the computation have suboptimal information about their exits. This not only slows down convergence, but can even aggregate and lead to precision problems and non-termination when there is a chain of many SCCs, as we show in the following example.

Example 1. To exemplify TVI and show when its precision problems occur, we consider an SSG that is a chain of n SCCs, each with one state. Every state either loops or continues to the next state, both with probability 0.5. At the end of the chain, we go to the goal with 0.6 and to the sink with 0.4.

Formally, $S = S_{\square} = \{t, z, s_0, s_1, \dots, s_n\}$, where s_0 is the initial state and $t \in F$ is the only goal state. There only is one action a , so $Av(s) = A = \{a\}$ for all states $s \in S$. For every s_i with $i < n$, we have $\delta(s_i, a, s_i) = \delta(s_i, a, s_{i+1}) = 0.5$ and for s_n , we have $\delta(s_n, a, t) = 0.6$ and $\delta(s_n, a, z) = 0.4$. Both states t and z are absorbing, so they loop with probability one.

Running topological bounded VI on this SSG, we first solve the bottom SCCs, i.e. t and z , and (by graph algorithms) infer their values of 1 and 0, respectively. Then we solve the SCC $\{s_n\}$ and set both its bounds to 0.6. Next, for the SCC $\{s_{n-1}\}$ bounded VI returns an ε -precise result, as with the self-loop the precise value is only obtained in the limit. Using precision of $\varepsilon = 10^{-6}$, the resulting interval is $[0.5999994277954102, 0.6000003814697266]$. Now the imprecisions start to add up: when solving the next SCC $\{s_{n-2}\}$, we depend on the imprecise bounds for $\{s_{n-1}\}$. Thus, the progress we make in every Bellman update is smaller. This not only slows down convergence, but it also leads to the first ε -precise interval being $[0.5999994099140338, 0.6000003933906441]$. So when BVI for the SCC $\{s_{n-2}\}$ terminates, both the lower and the upper bound are less precise than in the previous SCC. In state s_{n-19} , this imprecision has aggregated such that the computation is stuck at the interval $[0.5999994000000000, 0.6000004000000001]$, where the difference is larger than ε . Even though theoretically we make progress with a Bellman update, this progress is smaller than machine precision, so practically we can neither converge nor terminate.

Note that the SSG in this example is a Markov chain, so this problem occurs not only in SSGs, but already in Markov chains and MDPs. \triangle

We address this problem by introducing the precise-topological-optimization (PTVI, see Algorithm 2). The idea of PTVI is that, after an SCC has been solved with ε -precision (Line 4), we first extract the strategies for both players from the result (Line 5) and then compute the exact value of all states in the SCC under this pair of strategies (Line 6). Finally, we use a simple local check to verify that this is indeed the optimal value (Line 7). If it is, we return the precise values that the next SCCs can safely depend on (Line 8). If it is not, then we have to continue with some precise solution method (Line 10). Since we have just extracted near-optimal strategies, it makes sense to continue with SI, see e.g., [24, Section 3.2]. For details on the selection of the strategies and the proof of Theorem 2, see Appendix B.

Theorem 2. *Algorithm 2 returns the precise solution V .*

The strength of PTVI is the simple local check that allows it to conclude that the estimates for an SCC are precise. It relies on guessing both strategies. This differs from guessing an upper bound, as OVI does; or guessing one strategy, as in SI with a warm start [24, Section 4.3]. We emphasize that even using the classical naive stopping criterion in Line 4, this local check succeeded on more than 99% of the case studies, and thus the additional steps of Line 10 are almost never necessary. Using bounded VI in Line 4, we immediately succeeded on all case studies. In contrast, the first verification phase of OVI — having the same estimates and thus “information” as PTVI has when performing the local check — succeeded only for 15% of the random case studies; in 85% of the random cases as well as several larger real case studies OVI had to perform additional verification phases.

Note that PTVI can be seen from different directions. (i) It is a practical fix of TVI [16]. (ii) It is a new way to make classical VI return a precise result, which is more efficient than running for an exponential number of steps and rounding as described in [7]. (iii) It is a warm start for SI, in the seldom case that the SI phase of the algorithm (Line 10) is necessary. (iv) Just like OVI, it optimistically iterates the lower bound and then uses guessing to verify this guess. However, unlike OVI it produces a precise result, albeit at the cost of solving a Markov chain precisely; and it uses the information available at the time of guessing more efficiently, succeeding on the first check more often OVI.

5 Random Generation of Simple Stochastic Games

In order to properly evaluate and compare our algorithms, we need a diverse set of benchmarks. However, to the best of our knowledge, there are only 12 SSG case studies modelling real world problems and 3 handcrafted models for theoretical corner cases. Since the underlying structure of a model greatly affects the runtime of algorithms [24], only scaling these few models is insufficient. Thus, we propose an algorithm for random generation of SSG case studies, which enables us to test our algorithms on a broader spectrum of models.

Moreover, as we are interested in the relation between our verification algorithms and certain features of the model structure, our implementation also allows

for skewing the probability distribution towards models that exhibit certain features. This is very useful, since it allows us to test our algorithms on models whose features were not considered before (e.g., large number of actions per state, etc.). In particular, we provide: **(i)** parameters to tune features that can be affected by parameters of single states (e.g., the size, percentage of Minimizer states, actions per state, etc.). For example, if for each state the probability of being a Maximizer or Minimizer state is equal, we get 50% Minimizer states on average. Similarly, by choosing a high probability of adding another action to a state during the generation, we obtain states with up to 90 actions and an average around 7; **(ii)** more involved guidelines to affect features which depend on the interactions of several states (e.g., the number and size of SCCs and ECs, etc.). Intuitively, to obtain an SCC or MEC of a certain size n , we have to restrict the choice of successors during the transition or action generation to ensure that there are n strongly connected states.

We provide a detailed description of our random generation algorithm in Appendix C. There, we also prove that it can generate every possible SSG with positive probability and describe and discuss the aforementioned guidelines. Additionally, we give a detailed analysis of model features for all random case studies used in the evaluation, as well as a comparison to the features of the real case studies in Appendix D.

6 Experiments

In this section we talk about the practical evaluation of our algorithms and the comparison to the state of the art. First, we describe the setup in Section 6.1. Then we give a general overview in Section 6.2 before analyzing the algorithms' performance in more detail in Sections 6.3 and 6.4.

6.1 Experimental setup

Algorithms. Our implementation is based on PRISM-games [25] and available at <https://github.com/ga67vib/Algorithms-For-Stochastic-Games>.

We compare to the following algorithms from related work: classical value iteration (VI, [7]), bounded value iteration (BVI, [18]) and the improvement of bounded value iterations based on widest paths (WP, [28]). Moreover, as a representative of a competitor yielding a precise result, we implemented a precise variant of strategy iteration (SI), which relies on linear programming for solving the opponent MDP.

The new algorithms are optimistic VI (OVI, Section 3) and the precise topological version of VI (Section 4). For the latter, we give two variants with different stopping criteria in Line 4 of the algorithm: PTVI uses the naive criterion and PTBVI the ε -guaranteed one. Finally, we consider several optimizations, but their analysis is delegated to Appendix E.1 due to space constraints. Quite surprisingly, for all optimizations, their impact can be positive or negative on different models.

Case studies. We consider case studies from three different sources: (i) all real case studies that were already used in [24], and are mainly part the PRISM benchmark suite [26]; (ii) several handcrafted corner case models: haddad-monmege (the adversarial example from [20]), BigMec and MulMec (a single big MEC or a long chain of many small MECs from [24]), as well as two new models to analyse the

behaviour of OVI and one large model with many SCCs; (iii) randomly generated models as discussed in Section 5. Note that throughout our experiments, we omitted models solved by pre-computations.

Technical details. We conducted the experiments on a server with 64 GB of RAM and a 3.60GHz Intel CPU running Manjaro Linux. We always use a precision of $\varepsilon = 10^{-6}$. The timeout was set to 15 minutes and the memory limit was 6 GB for all models except for large models ($\geq 1,000,000$ states). For the large models, the timeout was set to 30 minutes and the memory limit to 36 GB.

6.2 Overview

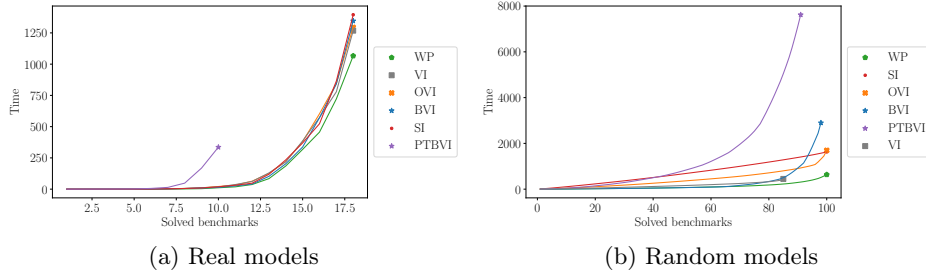


Fig. 2: Overview of the performance of the main algorithms on the *real* and *random* case studies. See Section 6.2 for a description.

Figure 2 gives an overview of the performance of the algorithms on the real and random case studies. The plots depict the number of solved benchmarks (horizontal axis) and the time it took to solve them (vertical axis). For each algorithm, the benchmarks are sorted in ascending order by verification time. A line stops when no further benchmarks could be solved. Intuitively, the further to the bottom right a line extends, the better. The algorithms shown in the legend on the right are sorted based on their performance, in descending order. Note that these plots have to be interpreted with care, as they greatly depend on the selection of benchmarks.

The precise algorithms provide harder guarantees, so we expect them to be slower. This is visible for PTBVI, which is slower and solves less benchmark than others. Still, PTBVI is optimal on certain kinds of models, as we detail in Section 6.3. Surprisingly, SI performed very well, even competing with the approximate algorithms BVI, OVI and WP. However, this comes from the model selection, particularly of the random models. Firstly, they exhibit very small transition probabilities, since we wanted the models to be hard for VI so that we can distinguish the different stopping criteria. This slows down convergence of VI, but does not affect SI. Secondly, they contain few states, so using a linear program is feasible. In Appendix E.4, we show that as model size increases, SI becomes less viable.

The algorithms giving ε -guarantees are overall quite comparable. This was also the case in the evaluation of [21], where the authors note that “*for probabilistic reachability, there is no clear winner*”. In Section 6.4, we give more details on how the performance of certain algorithms is affected by the structural features of a case study. Note that we included classical VI as a baseline, even though it gives

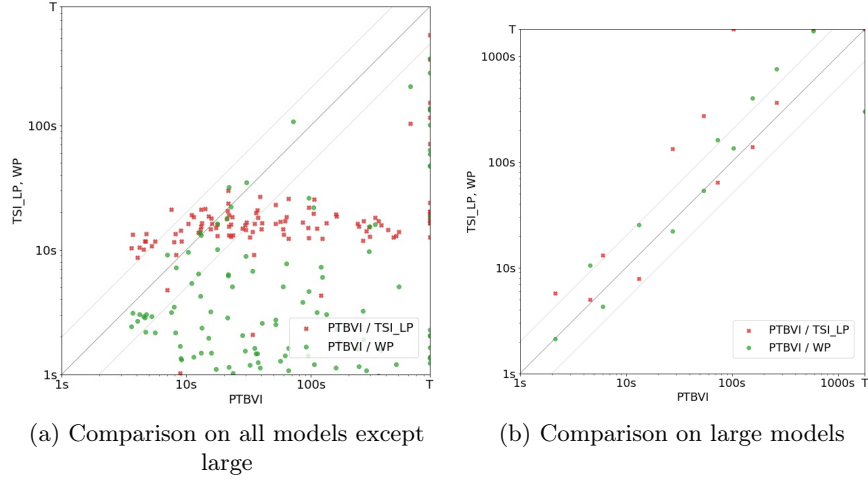


Fig. 3: PTBVI compared to SI and WP on all datasets.

no guarantees. It returned wrong results on two random models as well as the handcrafted haddad-monmege and MulMec.

Finally, it is important to note that random models of size 10,000 were already very hard for all algorithms, while some real models with more than 100,000 states could be solved quickly. This confirms the hypothesis of [24] that the graph structure of an SG (e.g., number of actions per state, depth of topological ordering, connectedness) is more important than its pure size.

6.3 Detailed analysis of precise algorithms.

PTVI and PTBVI are able to solve the chain of SCCs MulMec where normal topological VI [24] was stuck, so we achieved our original goal.

We use scatter plots to evaluate the algorithms' performance in detail. Each point in a scatter plot denotes a model. If a point is below the diagonal, the algorithm on the horizontal axis required more time to solve it than the corresponding algorithm on the vertical axis and vice versa. The two lines next to the diagonal mark the case where one algorithm was twice as fast as the other.

Figure 3 shows a scatter plot of PTBVI (which performed better than PTVI) versus the precise SI and the approximate, but very performant WP. While on smaller models PTBVI does not perform very well (Figure 3(a)), on larger models it often outperforms SI, in many cases halving the runtime or even reducing it by an order of magnitude, as shown in Figure 3(b). We conjecture that this comes from the fact that SI has to solve a linear program multiple times, while PTBVI only guesses the optimal strategies once and then solves a single Markov chain. We emphasize that PTBVI never had to resort to actually performing strategy iteration, because it guessed the correct strategies in all case studies. Moreover, PTBVI even beats the best approximate method, WP, in sufficiently large instances that contain multiple chained SCCs. In summary, PTBVI is a promising alternative to SI when needing precise solutions, especially on large models with chains of SCCs.

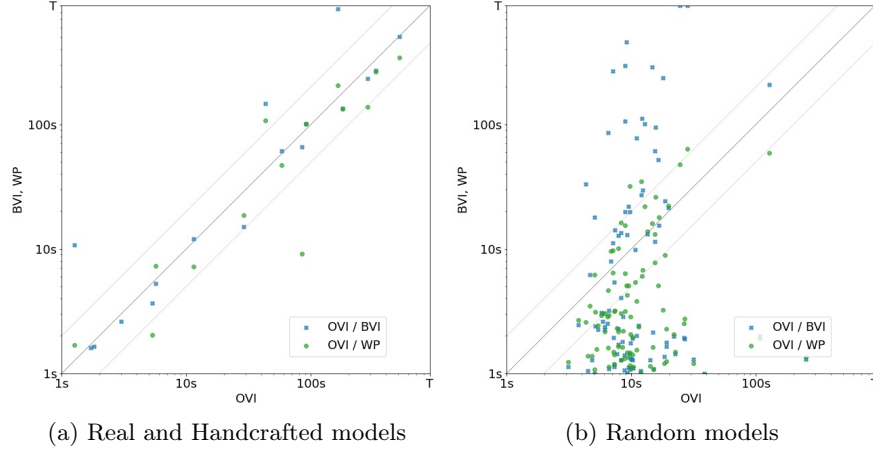


Fig. 4: OVI compared to BVI and WP on all datasets.

6.4 Detailed analysis of approximate (ε -precise) algorithms

All ε -precise algorithms perform similarly well. WP has the smallest accumulated runtime (Figure 2), no models where it is significantly worse than BVI (Appendix E.2) and only few models where it is significantly worse than OVI (Figure 4). As already observed in [28], it is particularly good when there are several or many MECs (especially on the handcrafted MulMec). Thus, it is a valid initial choice except when the models are large with a chain of big SCCs, where we concluded in Section 6.3 that PTBVI is better.

We analysed OVI in more detail to find out what features of the model affect its performance. Details validating the following statements are provided in Appendix E.3. Intuitively, OVI outperforms the other algorithms when the lower bound quickly converges, but the upper bound does not. Dually, if the lower bound converges slowly, this is problematic for OVI. Note that there are many hyper-parameters of OVI, for example the number of steps in the verification phase or the modification of the precision after a failed verification phase. We conjecture that these parameters affect the runtime and the choice can be improved; however, it is unlikely that there are parameter choices suitable for all kinds of models.

7 Conclusion

We extended optimistic VI from MDPs to SSGs. Moreover, using the “optimistic” idea, we fixed the issue of topological VI, so that it works even in the case of deeper models with more SCCs arranged in longer chains in the topological order. Besides, this fix also makes the method return the exact result. While this may be at the cost of a higher runtime, it becomes the only option when the overall model is very large, so that per-SCC analysis becomes unavoidable, and deep, so that precise values must be computed to converge at all. PTVI can be viewed as a separate algorithm or as an optimization on top of any approach from which a strategy can be extracted.

The experimental results show that the algorithms are of comparable performance, especially on real models from the standard benchmark sets. However, an in-depth analysis of the handcrafted and random models reveals that the performance of these algorithms is often sensitive to the underlying graph structure and, thus, their performance can vary accordingly. While we discuss some rules of thumb as to which algorithm is to be used for a particular benchmark, a part of the future work is to provide clearer and more algorithmic recommendations. An interesting direction here might also be to apply machine learning to recommend the most appropriate algorithm, as done for software model checkers already a few years ago, e.g., [15].

Moreover, we introduced a random generator, capable of producing various patterns even to extreme degrees. While this is very useful to find bugs and corner cases, many of the patterns need not be realistic. Consequently, we introduce a powerful set of tools to bias the generation. Nevertheless, future work shall amend this spectrum of tools with further hyper-parameters and approaches. We hope to hereby establish the platform for the community to contribute, complementary to benchmark sets [26,22].

References

1. Andersson, D., Miltersen, P.B.: The complexity of solving stochastic games on graphs. In: ISAAC. LNCS, vol. 5878, pp. 112–121. Springer (2009)
2. Baier, C., Katoen, J.P.: Principles of Model Checking. MIT Press (2008)
3. Baier, C., Klein, J., Leuschner, L., Parker, D., Wunderlich, S.: Ensuring the reliability of your model checker: Interval iteration for markov decision processes. In: CAV (1). LNCS, vol. 10426, pp. 160–180. Springer (2017)
4. Brázdil, T., Chatterjee, K., Chmelik, M., Forejt, V., Kretínský, J., Kwiatkowska, M.Z., Parker, D., Ujma, M.: Verification of markov decision processes using learning algorithms. In: ATVA. LNCS, vol. 8837, pp. 98–114. Springer (2014)
5. Cámara, J., Moreno, G.A., Garlan, D.: Stochastic game analysis and latency awareness for proactive self-adaptation. In: SEAMS 2014. pp. 155–164 (2014)
6. Chatterjee, K., Fijalkow, N.: A reduction from parity games to simple stochastic games. In: GandALF. pp. 74–86 (2011)
7. Chatterjee, K., Henzinger, T.A.: Value iteration. In: 25 Years of Model Checking. LNCS, vol. 5000, pp. 107–138. Springer (2008)
8. Chatterjee, K., Henzinger, T.A., Jobstmann, B., Radhakrishna, A.: Gist: A solver for probabilistic games. In: CAV. pp. 665–669. Springer Berlin Heidelberg (2010)
9. Chen, T., Forejt, V., Kwiatkowska, M.Z., Parker, D., Simaitis, A.: Automatic verification of competitive stochastic systems. FMSD, **43**(1), 61–92 (2013)
10. Chen, T., Kwiatkowska, M.Z., Simaitis, A., Wilsche, C.: Synthesis for multi-objective stochastic games: An application to autonomous urban driving. In: QEST. pp. 322–337 (2013)
11. Cheng, C., Knoll, A., Luttenberger, M., Buckl, C.: GAVS+: an open platform for the research of algorithmic game solving. In: TACAS. LNCS, vol. 6605, pp. 258–261. Springer (2011)
12. Condon, A.: The complexity of stochastic games. Inf. Comput. **96**(2), 203–224 (1992)
13. Condon, A.: On algorithms for simple stochastic games. In: Advances In Computational Complexity Theory. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, vol. 13, pp. 51–71. DIMACS/AMS (1993)
14. Courcoubetis, C., Yannakakis, M.: The complexity of probabilistic verification. J. ACM **42**(4), 857–907 (1995)

15. Czech, M., Hüllermeier, E., Jakobs, M., Wehrheim, H.: Predicting rankings of software verification tools. In: SWAN@ESEC/SIGSOFT FSE. pp. 23–26. ACM (2017)
16. Dai, P., Mausam, Weld, D.S., Goldsmith, J.: Topological value iteration algorithms. *J. Artif. Intell. Res.* **42**, 181–209 (2011)
17. Davey, B.A., Priestley, H.A.: Introduction to lattices and order. Cambridge university press (2002)
18. Eisentraut, J., Kelmendi, E., Křetínský, J., Weininger, M.: Value iteration for simple stochastic games: Stopping criterion and learning algorithm. *Information and Computation* (2022)
19. Grinstead, C.M., Snell, J.L.: Introduction to Probability. American Mathematical Society, second revised edition edn. (2006)
20. Haddad, S., Monmege, B.: Interval iteration algorithm for mdps and imdps. *Theor. Comput. Sci.* **735**, 111–131 (2018)
21. Hartmanns, A., Kaminski, B.L.: Optimistic value iteration. In: CAV (2). LNCS, vol. 12225, pp. 488–511. Springer (2020)
22. Hartmanns, A., Klauck, M., Parker, D., Quatmann, T., Ruijters, E.: The quantitative verification benchmark set. In: TACAS (1). LNCS, vol. 11427, pp. 344–350. Springer (2019)
23. Hoffman, A.J., Karp, R.M.: On nonterminating stochastic games. *Management Science* **12**(5), 359–370 (1966)
24. Křetínský, J., Ramneantu, E., Slivinskiy, A., Weininger, M.: Comparison of algorithms for simple stochastic games. *EPTCS* **326**, 131–148 (Sep 2020)
25. Kwiatkowska, M., Norman, G., Parker, D., Santos, G.: Prism-games 3.0: Stochastic game verification with concurrency, equilibria and time. In: CAV (2). LNCS, vol. 12225, pp. 475–487. Springer (2020)
26. Kwiatkowska, M.Z., Norman, G., Parker, D.: The PRISM benchmark suite. In: QEST 2012. pp. 203–204. IEEE Computer Society (2012)
27. LaValle, S.M.: Robot motion planning: A game-theoretic foundation. *Algorithmica* **26**(3-4), 430–465 (2000)
28. Phalakarn, K., Takisaka, T., Haas, T., Hasuo, I.: Widest paths and global propagation in bounded value iteration for stochastic games. In: CAV 2020, Los Angeles, CA, USA, July 21–24, 2020, Proceedings, Part II. LNCS, vol. 12225, pp. 349–371. Springer (2020)
29. Puterman, M.L.: Markov Decision Processes: Discrete Stochastic Dynamic Programming. Wiley Series in Probability and Statistics, Wiley (1994)
30. Quatmann, T., Katoen, J.: Sound value iteration. In: CAV (1). LNCS, vol. 10981, pp. 643–661. Springer (2018)
31. Svorenová, M., Kwiatkowska, M.: Quantitative verification and strategy synthesis for stochastic games. *Eur. J. Control* **30**, 15–30 (2016)
32. Tarjan, R.: Depth first search and linear graph algorithms. *SIAM JOURNAL ON COMPUTING* **1**(2) (1972)
33. Ujma, M.: On verification and controller synthesis for probabilistic systems at run-time. Ph.D. thesis, University of Oxford, UK (2015)

A Proof of Theorem 1

We first give an overview of the proof and delegate the technical details to the proofs of the lemmata below.

Correctness: Since L is computed by a classic VI variant, we know that $L \leq V$ by e.g., [7]. We have guessed U such that $\text{diff}(U(s), L(s)) \leq \varepsilon$ for all $s \in S$. It remains to show that upon termination $V \leq U$, i.e. that U is a correct upper bound.

The algorithm terminates if and only if it has found an upper bound with $U' = \mathcal{B}_L^D(U) \leq U$ (see Lines 7 and 8). To show that this implies our goal $U \geq V$, we apply standard arguments from lattice theory [17, Theorem 8.20]⁴. This is the point of Lemma 2. For it, we need that \mathcal{B}^D is a monotonic operator, which is shown in Lemma 1.

Termination: We prove this in Lemma 3 by a case distinction on the relation between U and V , similar to [21], but aggregating some cases and including an argument about the additional deflating.

Lemma 1. \mathcal{B}_L^D is monotonic, i.e. $f_1 \leq f_2$ implies $\mathcal{B}_L^D(f_1) \leq \mathcal{B}_L^D(f_2)$.

Before we prove this, we mention this is not true when using \mathcal{B}^D as BVI from [18] does. There, the L used for guessing the SECs is always updated and changed in every iteration. This can theoretically lead to non-monotonic behaviour. Thus, it is important that we fix L (and the implied set of SEC-candidates) for the whole verification phase. For Algorithm 1, this is obviously true as we do not update L during the verification phase. However, we must not continue iterating the lower bound during verifications as in [21] or depend on the upper bound when guessing the SECs.

Proof. $\mathcal{B}_L^D(f)$ modifies the given function two times: first by applying the Bellman operator \mathcal{B} , which is monotonic, as addition, multiplication, taking maximum or minimum and their combinations are monotonic. Secondly, $\mathcal{B}_L^D(f)$ applies deflating. Since L is fixed, the set of states affected by this is constant, namely the union of all SEC-candidates according to L . For every state in this set, f is updated to $\min(f(s), \text{bexit}_f(T))$. The min ensures that the estimate cannot increase, and the computation of the best exit again contains only max, multiplication and summation. Thus, \mathcal{B}_L^D is a monotonic operator.

Lemma 2. $\mathcal{B}_L^D(U) \leq U$ implies that $V \leq U$.

Proof. We restate [17, Theorem 8.20]: Given a complete lattice (P, \leq) and a monotonic self-map f on this lattice, $f(x) \leq x$ implies $\text{lfp}(f) \leq x$, where $\text{lfp}(f)$ denotes the least fixpoint of f .

We consider the complete lattice $(\{U \mid U : S \rightarrow \mathbb{R}\}, \leq)$ of functions mapping S to probabilities with the ordering according to pointwise comparison. The bottom element of the lattice is the vector that maps every state to 0, the top element maps every state to 1.

\mathcal{B}_L^D is a self-map on this lattice, as it takes a function of values for every state and returns another such function. Moreover, it is monotonic by Lemma 1. So, using [17,

⁴ This is what was called *Park induction* in [21]. Since we were not able to access the work cited in that paper, we use a variant of the same claim from another textbook.

Theorem 8.20] and instantiating f with \mathcal{B}^D and x with U , we get $\mathcal{B}^D(U) \leq U$ implies $\text{lfp}(\mathcal{B}^D) \leq U$. Thus, to prove our goal, it remains to show that $\text{lfp}(\mathcal{B}^D) = V$.

Deflating is sound and monotonic by [18, Lemma 3], i.e. it can never decrease a function U with $U \geq V$ below V and it can never increase a function. So deflating any $x \leq V$ just returns x and thus for $x \leq V$, \mathcal{B}^D does the same thing as \mathcal{B} . From this and using that V is the least fixpoint of the Bellman operator \mathcal{B} , see e.g., [7], we conclude

$$V = \text{lfp}(\mathcal{B}) = \text{lfp}(\mathcal{B}_L^D).$$

Lemma 3. *Given an SG \mathcal{G} and a lower bound $L_0 \leq V$, $\text{OVI}(\mathcal{G}, L_0, \varepsilon, \varepsilon)$ terminates.*

Proof. We prove this by a case distinction on the relation between U and V , similar to [21]. Our proof mainly differs in that we aggregate Cases 2 and 3 in their proof into case 1 of our proof and include an argument about the additional deflating.

- $U(s) \geq V(s)$ for all $s \in S$: If U truly is an upper bound on the value, we want to terminate. Note that this is also the only case in which we can terminate, as otherwise we recursively call OVI again with more precision for VI and more iterations in the verification phase. Thus, since the precision of VI increases, we can assume that L has converged close enough to V such that all SECs are detected correctly. Note that it is possible that we terminate even if the SECs are not yet guessed correctly. However, to prove termination, this assumption is necessary.

By [18, Theorem 2], the value function V is the unique fixpoint of the operator \mathcal{B}_L^D , since Bellman updates and deflating together ensure that U converges to the value.

As in [21], we made the Bellman update monotonic by using Line 9, to avoid a particular corner case where there is an alternating increase and decrease, see [21, Section 4.2].

The remaining argument is the same as in Cases 2 and 3 of the original termination proof. Intuitively, given enough steps (and the number of steps in the verification phase increases, so there eventually will be enough steps), all states see a decrease in value. The corner case of a state already having the correct value and hence not decreasing is handled in the same pragmatic way: we just abort the verification phase and try again with another upper bound that is necessarily higher than before.

- Otherwise, so if there exists a state $s \in S$ with $U(s) < V(s)$: In such a case, the verification phase will certainly be aborted, as we terminate if and only if $\mathcal{B}_L^D(U) \leq U$, which cannot happen when some $U(s) < V(s) = \text{lfp}(\mathcal{B}_L^D)(s)$. Thus, we need to abort the verification phase with this U , which happens after a finite number of iterations (see Line 5). Note that the additional breaking conditions mentioned at the end of Section 3 help to detect this case faster.

It remains to show that we cannot stay in this case forever, but that eventually we guess a $U \geq V$. This happens, because L converges to V , see e.g., [7]. Since the increase of $\text{diff}_\varepsilon^+$ for guessing the upper bound is constant⁵, we will eventually guess an upper bound that is greater than the value and our algorithm will eventually detect that and terminate.

⁵ Or, in the case of relative error, can be lower bounded by $L(s) * \varepsilon$

B Proofs for Section 4

B.1 Details on PTVI (Algorithm 2)

Computing optimal strategies: How to compute the strategies is a heuristic that should get the best from the given information of the VI-algorithm. We phrased it such that the computation depends on both a lower bound L and an upper bound U in order to show how to use it with both BVI and OVI. When using naive VI, we set $U := L$. We use U to guess the Minimizer strategy and L to guess the Maximizer strategy, as this is using the most conservative estimate that we have available.

From an estimate function, we can derive a strategy by picking actions that are optimal according to this estimate, similar to the Bellman equations 2.2. Indeed, for Minimizer this is sufficient, so we set

$$\tau \leftarrow \{(s, a) \mid s \in S_{\bigcirc}, a \text{ picked randomly from } \arg \min_{a' \in A(s)} U(s, a')\}$$

For Maximizer this is not sufficient, because Maximizer might have actions with optimal value that however do not make progress towards the target. Imagine for example a Maximizer state with an action a that loops and an action b leading to the goal. While picking a for a finite number of times does not decrease the value, as the state can still play b and reach the goal surely, committing to playing a infinitely often reduces the value to 0. Thus, when deriving a Maximizer strategy, we cannot just pick some optimal action.

There are several ways to deal with this problem: firstly, to obtain a memoryless deterministic strategy, one can use graph algorithms to compute a distance measure between states in S_{\circ} and goal states; then the strategy picks actions that not only are optimal according to the estimate, but additionally reduce the distance to the goal. This is similar to the construction used in [1, Theorem 2] for deriving strategies of Maximizer in a reachability game. However, we want to avoid the additional computation time for these graph algorithms. Thus, we use a second approach: strengthening the notion of strategy to allow for randomization. A randomized strategy is map $S \rightarrow \text{Dist}(A)$ that for every state s returns a probability distribution over available actions $A(s)$. We write $(s, a, p) \in \sigma$ to say that $\sigma(s)(a) = p$; a strategy is well-defined when all triples with $p > 0$ are enumerated. Thus, we select the Maximizer strategy as

$$\sigma \leftarrow \{(s, a, p) \mid s \in S_{\square}, a \in \arg \max_{a' \in A(s)} L(s, a'), p = \frac{1}{|\arg \max_{a' \in A(s)} L(s, a')|}\}$$

This is correct, because every action that is optimal (according to the estimate) has positive probability to be selected. The “staying” actions do not decrease the value, and almost surely we eventually pick an action that makes progress towards the goal. In the words of [1, Lemma 5]: the strategy is safe and stopping. Note that this does not drastically decrease the transition probabilities in the resulting Markov chain, as there typically are less than 3 actions per state. Hence it poses no problem for the Markov chain solving.

Markov chain solving: To solve the induced Markov Chains, we use the standard equation approach as described in [19, Chapter 11]. First, the MC is represented in a transition matrix. Next, one has to find the states $s \in S^?$ which cannot reach any target in the MC and remove their rows and columns from the matrix. To do so, it is usually necessary to perform an all-pairs-shortest-path algorithm like Floyd-Warshall. However, since this requires $O(|S|^3)$ operations, and we already have L from the VI process, we can obtain said states instead directly by checking whether $L(s) = 0$. Removing states that cannot reach any target from the transition matrix guarantees non-singularity, and thus invertability. Lastly, the part of the matrix corresponding to $S^?$ needs to be inverted and multiplied with the part of the transition matrix corresponding to F , yielding the reachability probability for every state $s \in S^?$ in the Markov chain.

In our implementation, we use the JAMA library⁶ to perform matrix operations like multiplications or inversions.

Strategy iteration in the case of non-optimal strategies: The local check we perform in Line 7 gives guarantees if and only if it succeeds. So if it fails, we know nothing about the strategies and have to perform normal SI. This means we can only fix one of them, as solving the resulting MDP might return a different strategy for the other player. However, it is quite likely that both strategies are good. Note that for SI we need memoryless deterministic strategies, so using the randomized σ as we defined above is not possible. Thus it makes sense to resort to starting with τ .

B.2 Proof of Theorem 2

Theorem 2. *Algorithm 2 returns the precise solution V .*

Proof. The argument for correctness of solving the SG in topological order is the same as in [24, Section 4.4].

In the case that the condition in Line 7 evaluates to false, there is nothing to prove, because in this case the algorithm falls back to using strategy iteration which is known to return the precise solution, see e.g. [13]. So we only have to prove that, if the check in Line 7 evaluates to true, this implies that σ and τ are optimal strategies and $V_{\mathcal{G}}^{(\sigma, \tau)}$ is the value function of the SG.

In the following, for any SG (or MDP or Markov chain) \mathcal{G} , we use $V_{\mathcal{G}}$ to denote its value function. Further, $\mathcal{G}^{(\sigma, \cdot)}$ and $\mathcal{G}^{(\cdot, \tau)}$ are the MDP with Maximizer strategy σ respectively Minimizer strategy τ fixed, and $\mathcal{G}^{(\sigma, \tau)}$ is the Markov chain with both strategies fixed.

Intuitively, our local check amounts to performing SI on the MDPs $\mathcal{G}^{(\sigma, \cdot)}$ and $\mathcal{G}^{(\cdot, \tau)}$. For the proof, first consider the MDP $\mathcal{G}^{(\sigma, \cdot)}$. SI in this MDP starts by fixing an arbitrary Minimizer strategy and we choose τ . Then it computes the value of the resulting Markov chain $\mathcal{G}^{(\sigma, \tau)}$ and checks whether Minimizer wants to switch the strategy in any state, i.e. whether there exists an $s \in S_{\bigcirc}$ such that $\tau(s) \notin \arg \min_{a \in A(s)} V_{\mathcal{G}^{(\sigma, \tau)}}(s, a)$. If Minimizer wants to switch, τ is not optimal and the check evaluates to false. Otherwise, τ is an optimal strategy in the MDP $\mathcal{G}^{(\sigma, \cdot)}$, cf. [29]. Thus, we have $V_{\mathcal{G}^{(\sigma, \cdot)}} = V_{\mathcal{G}^{(\sigma, \tau)}}$.

⁶ <https://math.nist.gov/javanumerics/jama/>

Dually, the check for the Maximizer strategy verifies that $V_{\mathcal{G}^{(\cdot, \tau)}} = V_{\mathcal{G}^{(\sigma, \tau)}}$ by proving that σ is an optimal strategy in the MDP $\mathcal{G}^{(\cdot, \tau)}$. Note here that, since we use a randomized σ , we have to check that every action to which σ assigns positive probability is optimal, i.e. $\{a \mid \exists p > 0 : (s, a, p) \in \sigma\} \subseteq \arg \max_{a \in A(s)} V_{\mathcal{G}^{(\sigma, \tau)}}(s, a)$.

Finally, we have to relate $V_{\mathcal{G}^{(\sigma, \cdot)}}$ and $V_{\mathcal{G}^{(\cdot, \tau)}}$ to the value on the actual game $V_{\mathcal{G}}$. Observe that by fixing an arbitrary Maximizer strategy, we can only decrease the value, as it might be suboptimal, but we can never increase the value. Formally, $V_{\mathcal{G}^{(\sigma, \cdot)}} \leq V_{\mathcal{G}}$. Dually, $V_{\mathcal{G}} \leq V_{\mathcal{G}^{(\cdot, \tau)}}$. Putting everything together, we have

$$V_{\mathcal{G}^{(\sigma, \tau)}} = V_{\mathcal{G}^{(\sigma, \cdot)}} \leq V_{\mathcal{G}} \leq V_{\mathcal{G}^{(\cdot, \tau)}} = V_{\mathcal{G}^{(\sigma, \tau)}}.$$

This proves that indeed both strategies are optimal in \mathcal{G} and $V_{\mathcal{G}^{(\sigma, \tau)}}$ is the correct value function.

C Additional details for randomly generated models

C.1 Our algorithm for random model generation

We use Algorithm 3 to create any random stochastic game that is connected from the initial state. After initialization, the algorithm has two phases: the forward and the backward procedure. During initialization, we generate a random number n and create a set of states $S := [n]$. Also, we assign states at random to either Maximizer or Minimizer. In the forward procedure, we iterate over every state $s \in S$ and make sure that a previous state s' is connected to it by providing an action with positive transition probability to s from s' . This guarantees that the initial state can reach every state in the stochastic game. The backward procedure then adds arbitrary actions to arbitrary states to enable generating every possible SG.

To generate the actions of a state, we use Algorithm 4. It receives a state-action pair (s, a) where $\sum_{s' \in S} \delta(s, a, s') < 1$. It then increases the transition probability of a randomly selected state $s' \in S$ where $\delta(s, a, s') = 0$. This is repeated until $\sum_{s' \in S} \delta(s, a, s') \geq 1$ or there is no state s' that is not yet reached by (s, a) . In case $\sum_{s' \in S} \delta(s, a, s') < 1$ holds but the state-action pair is reaching every state in S , we increase the most recently increased transition probability such that the resulting distribution is a probability distribution. If we reach $\sum_{s' \in S} \delta(s, a, s') > 1$, we reduce the transition probability we increased most recently. After applying Algorithm 4, (s, a) is a valid transition distribution, i.e. its probabilities sum up to 1.

Lemma 4. *Algorithm 3 creates formally correct stochastic games.*

Proof. S is finite since its size is determined by a random number $n \in \mathbb{N}$ (Line 1). Line 2 ensures that we have a partition of S into S_{\square} and S_{\circ} . Next, Lines 3-5 provide an initial state and a target. Thus, we only need to argue that Av is truly a mapping of $S \rightarrow 2^A$ and that the transition function yields a probability distribution. When we introduce state-action pairs (Lines 9-14 and Lines 18-21), we introduce a new action that we add to A . Also, we add the state-action pair to Av (Lines 12 and 19). Thus, any Av is function of $S \rightarrow 2^A$.

Algorithm 3 Generating random models connected from initial state**Output:** Stochastic game \mathcal{G} where the initial state is connected to any $s \in \mathbf{S}$

```

1: Create  $\mathbf{S}$  with a random  $n \in \mathbb{N}$ 
2: Partition  $\mathbf{S}$  uniformly at random into  $\mathbf{S}_{\square}$  and  $\mathbf{S}_{\circ}$ 
3: Enumerate  $s \in \mathbf{S}$  in any random order from 0 to  $n-1$ 
4: Set  $s_0$  to the state with index 0
5: Set  $\mathbf{F} = \{s_{n-1}\}$ 
6: for  $s = 1 \rightarrow n - 1$  do ▷ Forward Procedure
7:   if  $s$  does have an incoming transition then Continue (Skip iteration)
8:   else
9:     Pick any state  $s'$  with index smaller than  $s$ 
10:    Create an action  $a$  that starts at  $s'$ 
11:    Assign to  $(s', a)$  a positive probability of reaching  $s$ 
12:    Create a valid probability distribution for  $(s', a)$  by applying FillAction( $s', a$ )
13:    Add  $a$  to  $\mathbf{Av}(s')$ 
14:    Add  $a$  to  $\mathbf{A}$ 
15: for  $s = n - 1 \rightarrow 0$  do ▷ Backward Procedure
16:   Pick a random number  $m \in \mathbb{N}$  ▷ Add as many actions as possible
17:   if  $|\mathbf{Av}(s)| = 0$  then  $m \leftarrow \max\{m, 1\}$  ▷ Every state must have at least one action
18:   for  $i = 1 \rightarrow m$  do
19:      $(s, a_i) = \text{FillAction}(s, a_i)$ 
20:     Add  $a_i$  to  $\mathbf{Av}(s)$ 
21:     Add  $a$  to  $\mathbf{A}$ 

```

Algorithm 4 FillAction(s, a)**Input:** outgoing state s , action a **Output:** action a that has a valid underlying transition probability distribution

```

1: repeat
2:   Pick a random state  $s'$  where  $\delta(s, a, s') = 0$ 
3:   Increase  $\delta(s, a, s')$  by a random number  $\in (0, 1]$ 
4: until either  $\sum_{s' \in \mathbf{S}} \delta(s, a, s') \geq 1$  or  $\forall s' \in \mathbf{S} : \delta(s, a, s') > 0$ 
5: if  $\sum_{s' \in \mathbf{S}} \delta(s, a, s') > 1$  then
6:   Decrease the most recently modified  $\delta(s, a, s')$  so  $\sum_{s' \in \mathbf{S}} \delta(s, a, s') = 1$ 
7: else if  $\sum_{s' \in \mathbf{S}} \delta(s, a, s') < 1$  then ▷ Loop terminated because  $\forall s' \in \mathbf{S} : \delta(s, a, s') > 0$ 
8:   Increase the most recently modified  $\delta(s, a, s')$  so  $\sum_{s' \in \mathbf{S}} \delta(s, a, s') = 1$ 
return  $(s, a)$ 

```

Next we need to prove that for every $s \in \mathbf{S}$ and every action $a \in \mathbf{Av}(s)$ the transition function $\delta(s, a)$ yields a probability distribution. In other words we need to validate that (i) for every state $s' \in \mathbf{S} : \delta(s, a, s') \in [0, 1]$ and (ii) $\sum_{s' \in \mathbf{S}} \delta(s, a, s') = 1$ are true.

To prove (i), note that whenever we introduce an action a to the set of enabled actions $\mathbf{Av}(s)$ of a state $s \in \mathbf{S}$, we have $\delta(s, a, s') = 0$ for all $s' \in \mathbf{S}$. We increase $\delta(s, a, s')$ in Line 11 of Algorithm 3 and Lines 3 and 8 of Algorithm 4. We increase transition probabilities by numbers in $[0, 1]$. Due to the condition in Line 4 of Algorithm 4, $\delta(s, a, s')$ can only be increased a second time in Line 8 of Algorithm 4. However, per state-action pair (s, a) with $a \in \mathbf{Av}(s)$ Line 8 of Algorithm 4 may be executed only once and we increase only by a $\Delta > 0$ such that $\delta(s, a, s') + \Delta + \sum_{s'' \in \mathbf{S} \setminus \{s'\}} \delta(s, a, s'') = 1$. Since every state $s'' \in \mathbf{S} \setminus \{s'\}$ was increased only once, $\delta(s, a, s') + \Delta \leq 1$.

To prove (ii), note that every pair (s, a) where $a \in \text{Av}(s)$ is given to Algorithm 4. Once the loop (Lines 1-4) terminates, it either holds that (a) for every state $s' \in \mathcal{S} : \delta(s, a, s') > 0$ or that (b) $\sum_{s' \in \mathcal{S}} \delta(s, a, s') \geq 1$. If $\sum_{s' \in \mathcal{S}} \delta(s, a, s') = 1$ the algorithm terminates correctly. Hence, we consider the other two cases: In case (a) Line 8 increases the most recently added triple (s, a, s') by a $\Delta > 0$ such that $\delta(s, a, s') + \Delta + \sum_{s'' \in \mathcal{S} \setminus \{s'\}} \delta(s, a, s'') = 1$. In case (b) $\sum_{s' \in \mathcal{S}} \delta(s, a, s') > 1$. Due to the exit condition of the loop (Line 4), without the most recently increased transition (s, a, s') the sum of the probabilities must be below 1, i.e. $\sum_{s'' \in \mathcal{S} \setminus \{s'\}} \delta(s, a, s'') < 1$. Thus, there is a $\Delta \in (0, 1), \Delta < \delta(s, a, s'')$ such that $\delta(s, a, s') - \Delta + \sum_{s'' \in \mathcal{S} \setminus \{s'\}} \delta(s, a, s'') = 1$. We decrease (s, a, s') by this Δ . In conclusion, every pair (s, a) that is provided to Algorithm 4 yields a valid transition distribution where $\sum_{s' \in \mathcal{S}} \delta(s, a, s') = 1$. Thus, Algorithm 3 generates formally correct stochastic games.

To argue about the SGs that the algorithm produces, we introduce some notation. Let $\mathcal{G}_{\text{Algo}}$ be the set of SGs that Algorithm 3 can produce and $\mathcal{G}_{\text{reach}}$ be the set of all SGs where every state is reachable from the initial state. Note that every SG that is not an element of $\mathcal{G}_{\text{reach}}$ can be transformed into an SG with the same value and optimal strategies by removing all unreachable states. Hence, we only care about producing SGs in $\mathcal{G}_{\text{reach}}$.

Lemma 5. *Algorithm 3 creates all SGs where every state is reachable from the initial state, i.e. $\mathcal{G}_{\text{Algo}} = \mathcal{G}_{\text{reach}}$.*

Proof. We first show that $\mathcal{G}_{\text{Algo}} \subseteq \mathcal{G}_{\text{reach}}$:

For this statement to hold, any $\mathcal{G} \in \mathcal{G}_{\text{Algo}}$ must be connected from the initial state. Proof by induction over the indices i of the states along their enumeration assigned during Algorithm 3:

Basis: $i = 0$: s_0 is the initial state. The initial state can reach itself within 0 steps.

Hypothesis: Let i be arbitrary but fixed with $i \leq n - 1$, where $n = |\mathcal{S}|$. For every $j \leq i$ it holds that s_0 can reach s_j .

Inductive Step: $i \leftarrow i + 1$: Due to the forward procedure it holds that

$$\exists s_j \in \mathcal{S}, j < i, a \in \text{Av}(s_j) : \delta(s_j, a, s_i) > 0$$

However, according to our hypothesis s_0 is connected to s_j and thus also to s_i .

Now we show that $\mathcal{G}_{\text{reach}} \subseteq \mathcal{G}_{\text{Algo}}$:

Pick an arbitrary but fixed stochastic game $\mathcal{G} \in \mathcal{G}_{\text{reach}}$. Next, we show that there is a run of our algorithm that will return a stochastic game $\mathcal{G}' \in \mathcal{G}_{\text{Algo}}$ where \mathcal{G}' is an automorphism to \mathcal{G} . Thus, \mathcal{G}' and \mathcal{G} are the same except for the state enumeration and the labels of the actions.

For this, we need several statements to hold at once:

1. The number of states in \mathcal{G} and \mathcal{G}' is equal.
2. The partition of \mathcal{S} to \mathcal{S}_{\square} and \mathcal{S}_{\bigcirc} is the same for \mathcal{G} and \mathcal{G}' .
3. \mathcal{G} and \mathcal{G}' have the same initial states and targets.
4. All state-action pairs in \mathcal{G} and \mathcal{G}' yield the same probability distributions in δ .

5. Every state in \mathcal{G} and \mathcal{G}' has the same actions.

We decide randomly on the number of states in Line 1 of Algorithm 3 and partition S into S_{\square} and S_{\circ} at random in Line 2 of Algorithm 3. Thus, for every $\mathcal{G} \in \mathcal{G}_{\text{reach}}$ there exists $\mathcal{G}' \in \mathcal{G}_{\text{Algo}}$ that have the same number of states to which there is an enumeration such that they are partitioned equally in both stochastic games. Since the states can be arranged in any order, we pick the initial state of \mathcal{G}' such that is the same as in \mathcal{G} . All targets of \mathcal{G} can be mapped to the singular target of \mathcal{G}' , since they only have self-loops and behave identically to f of \mathcal{G}' . Thus, there exists a run where Statements 1, 2, and 3 hold.

When using Algorithm 4 to create a probability distribution for a state-action pair (s, a) , we increase transition probabilities until they sum up to 1. Thus, any summation $\sum_{s' \in S} \delta(s, a, s') = 1$ is possible. In consequence, an action may lead into arbitrary states, have an arbitrary number of positive transition probabilities between 1 and $|S|$, and may have arbitrary probability distributions on the transitions as long as they sum up to 1. So out of all runs where Statements 1, 2, and 3 hold, there also must be at least one run where statement 4 holds too.

To show that Statement 5 holds in addition, note that each $\mathcal{G} \in \mathcal{G}_{\text{reach}}$ has a minimal set of state-action tuples such that the initial state is connected to every state. Taking this set, we can perform a breadth-first search from the initial state to provide an enumeration of the states. If we iterate over the states along this enumeration, we can reproduce each of the actions in the minimal set during the forward process. Due to the enumeration, to each state s except for the initial state, there is a state with a smaller index s' such that s' has an action a with a positive transition probability of reaching s . Since every other transition of (s', a) can lead into arbitrary states and the probability distribution of (s', a) can be arbitrary, we can recreate the minimal set of state-action tuples in \mathcal{G}' .

The remaining state-action pairs of \mathcal{G} can be added to \mathcal{G}' during the backward process, where every state may add arbitrarily many actions with arbitrary transition distributions.

Note that although $\mathcal{G}_{\text{Algo}} = \mathcal{G}_{\text{reach}}$, in general Algorithm 3 does not sample $\mathcal{G}_{\text{reach}}$ uniformly at random. Due to the forward procedure, states with smaller indices tend to have more actions than the states with higher ones. Additionally, creating the transition distributions as described in Algorithm 4 favors state-action pairs to have few transitions. If we pick the transition probabilities between $(0, 1]$ uniformly at random, around 83,33% of all actions have two or three transitions with positive probability and none with one transition.

C.2 Parameters and guidelines for model construction

We implemented a constrained version of Algorithm 3 from Appendix C.1 to randomly generate models. The real-world implementation has to be constrained due to the natural restriction that a computer cannot generate arbitrarily large stochastic games and arbitrarily small transitions due to finite memory. Additional constraints on the real-world implementation are the pseudo-randomness while taking decisions, as well as floating point machine precision. Moreover, we want to give the users some control over the properties of the resulting models like the number of states or the partitioning into Minimizer and Maximizer states. If all model

properties vary significantly, it is very hard to deduce why an algorithm performs differently on two models. Thus, we provide several parameters which can be set to restrict the randomization. Some examples of the parameters that one can control are the number of states, number of models, smallest probability that is allowed to occur, number of transitions of an action.

Limits and additional guideline options

Although the parameters we expose for random generation can directly influence various structural properties of the resulting models, there are other structural properties that are hard to influence. For example, there is no direct way to affect the size and number of SCCs of a model, or to guarantee that every state in a model has a certain number of actions when using Algorithm 3. We provide some guidelines that can be followed to influence such properties.

1. **RandomTree guideline.** We refer to the guideline that controls how many actions a state has as *RandomTree* because it creates a *tree-like* graph structure where the initial state is the root. Every node of the tree has k actions and at most k children, where k is a parameter. Every action has an assigned child to which it has a positive transition probability. The rest of the probability distribution of the action is assigned at random. An inner node of the tree may have less than k children if adding k children would exceed the requested number of states n for the model. Also, leaves are not required to have k actions. Their actions are only introduced during the backward process and are there to enable the generation of end components.
2. **RandomSCC guideline** This is a guideline that controls the SCCs of a model. The procedure requires a minimal and maximal size boundary $[a, b]$ for every SCC and the total number of states in the stochastic game n . First, we create subgames of a randomly chosen size in $[a, b]$. The subgames are created by the algorithm described in Appendix C.1. We then use Tarjan’s algorithm for strongly connected components [32] to identify the SCCs of the created subgame. Next, we unify all the SCCs of the subgame by using the topological enumeration Tarjan’s algorithm provides. We circularly connect the SCCs along the enumeration, making the whole subgame an SCC. Next, we make sure that the subgame is connected to the rest of the stochastic game by making sure a previously created subgame has an action leading into this subgame. We repeat this procedure until we have at least n states in the stochastic game, resulting in a stochastic game \mathcal{G} that is connected from the root where the user has an easy way of controlling the number and size of the SCC in \mathcal{G} .

D Details on model analysis

In this section, we discuss the feature distribution of the real and randomly generated models. We claim that such a feature analysis of benchmark sets is important in order to find out biases, as well as judge why certain sets of benchmarks are hard. We first give a few noteworthy conclusions of our analysis.

- Often large parts of the real case studies are solved by pre-computations. Hence, instead of the size of the state space $|S|$ we should consider the size of the unknown part $|S^?|$.

- In the real case studies we have, there are usually only few actions and successors per state, few MECs and few non-trivial SCCs. Thus, to assess general algorithm performance, we should generate models with these features, be it by hand, randomly or by finding realistic problems with this structure.
- In general, the models generated by our algorithm without using the guidelines are hard because they form very connected models with many actions and transitions as well as low transition probabilities. Apart from that, the models generated do not exhibit large size, many MECs or long chains of SCCs.

We provide two figures, each analysing 16 features of a set of case studies. Figure 5 shows the feature distribution of the real case studies. We refer to this set as REAL. Figure 6 shows the feature distribution of a randomly generated set of models with many actions. We refer to this set as RANDOM. For our analysis we use box plots, where for every feature distribution, the plot shows the median (orange line), average (green triangle), 25 and 75 percentile (box), 10 and 90 percentile (whiskers) and outliers (circles). In every figure, the box plots are grouped by and coloured according to the following categories:

- Green outlines are for features related to states.
- Blue outlines are for features related to actions.
- Cyan outlines are for features related to transitions.
- Red outlines are for features related to MECs.
- Orange outlines are for features related to SCCs.

For every feature we point out differences and similarities between the sets. We also motivate why certain features are interesting and our choices of parameters for the set RANDOM.

- NumStates: Number of states in a model, given in log scale. For REAL, this has a spread between several dozen up to a few million states. For RANDOM, we always chose 10,000, so that the complexity of the model could not come from its size, but the size was still non-trivial. To analyze the complexity induced by the sheer size of a model, we used the additional handcrafted scalable model described in Appendix E.4.
- Sinks% and Unknown%: These features describe the percentage of states that are sinks (Z) or unknown states ($S^?$). Note that all randomly generated models have exactly one goal state. Hence, we have that Sinks% + Unknown% sums up to (a little less than) 1. For REAL, the percentage of unknown states is almost always less than 40% and less than 20% in half the cases. This means that most of the model is solved by pre-computations, being either sinks or goal states. We highlight one implication of this: an experimental analysis relating the number of states to the performance of algorithms is flawed, since large parts of the state space are quickly solved by pre-computations shared by all algorithms. Thus, such an analysis should rather consider the number of unknown states. In RANDOM, the percentage of unknown states is around 60%, i.e., typically more of the model requires an actual computation of the solution algorithm.
- MinStates%: Percentage of states belonging to Minimizer. Note that the percentage of Maximizer states is 1 minus this. For all sets, the average is around 50%. For REAL, the variance is higher.

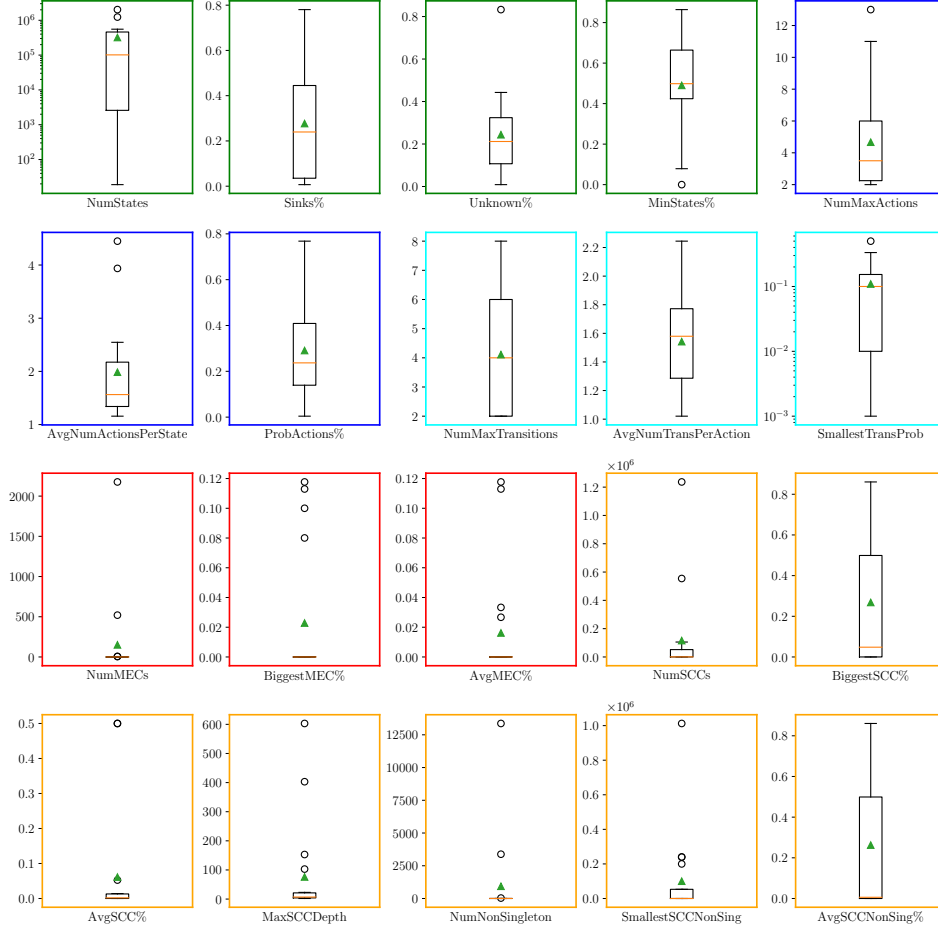


Fig. 5: REAL: Box plots for analysing the features of the *real case studies*. A description of how to read box plots is provided in Section 5.

- NumMaxActions, AvgNumActionsPerState, ProbActions%, NumMaxTransitions, AvgNumTransPerAction: These features intuitively describe the “breadth” or “branching” of the model. They are, in order: the maximum number of actions per state occurring in the model, the average number of actions per state, the percentage of probabilistic actions (with more than one successor state), the maximum number of transitions occurring for a state-action pair, and the average number of transitions occurring for a state-action pair.

For REAL, we typically have slightly less than 2 actions per state, usually no more than 5 and never more than 13. Only a third of the actions are probabilistic. Typically, we have between one and two successors, seldom more than 4 and never more than 8.

Since we wanted to explore graph structures that are not present in REAL, for the generation of RANDOM we allowed our algorithm to create models with more actions and transitions and higher branching. We typically have

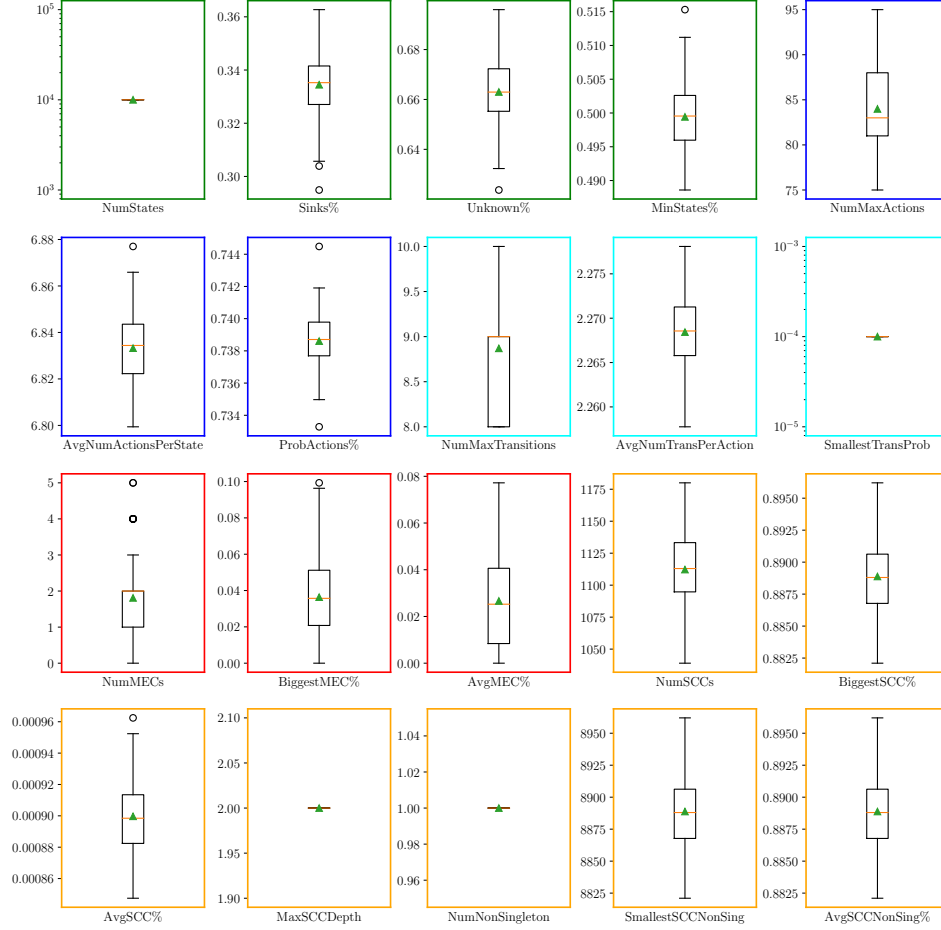


Fig. 6: RANDOM: Box plots for analysing the features of the set of *randomly generated models*. A description of how to read box plots is provided in Section 5.

around 7 actions per state, but also up to 95. Three quarters of the actions are probabilistic. We have an average of slightly more than 2 transitions per action and a maximum of 10. The resulting models are more connected and thus often harder to solve for the algorithms.

- SmallestTransProb: The smallest transition probability occurring in the model. In REAL, this goes from 0.5 (the largest possible non-trivial probability) to 10^{-3} . In RANDOM, we always set this to 10^{-4} , making the models generally hard for value iteration, so that non-trivial solution times occur.
- NumMECs, BiggestMEC%, AvgMEC%, NumSCCs, BiggestSCC%, AvgSCC%: These features are related to the important graph theoretic concepts of MEC and SCC. We give their number (Num), the size of the biggest occurring MEC/SCC in percentage of the state space (Biggest%) as well as the average size of MECs/SCCs. Note that we only count MECs in the unknown part of the state space.

Most of the models in REAL have very few and very small MECs, with a few exceptions going up to around 2000 MECs. Since we often only have one MEC, the biggest and average MEC size box plot look very similar. In contrast, in RANDOM we usually have around 2 non-trivial MECs, but they also are small. To analyze the impact of many or big MECs, we used the handcrafted model MulMec. However, we can also offer guidelines to create random models with certain numbers and sizes of MECs, analogous to the RandomSCC guideline. For SCCs, note that the scale of the plot for the number of SCCs is multiplied by one million. In REAL we usually have around 100,000 SCCs, with an outlier having more than a million. The biggest SCC is often large, comprising a third of the model on average and going up to 50% in three quarters of the cases. Since there are many transient states (which form a singleton SCC of size 1), the average SCC size is small. Since this does not give a lot of information, we also give the average size of non-singleton SCCs, see below. In RANDOM, the variance is smaller. There are around 1,000 SCCs, and the biggest SCC makes up around 90% of the model, making it almost completely strongly connected. This is because our random generation picks successor states randomly, and chances are good that this induces cycles throughout most of the model. This is why we offer the RandomSCC guideline, which allows to create models that less connected.

- MaxSCCDepth: This is important for topological algorithms, as it is the depth of the DAG forming the graph. Note that Example 1 only needed a chain of 20 SCCs to show the problems of topological value iteration. In REAL, the average depth is around 100 and it can go up to 600. In RANDOM, this depth is low, because long chains of SCCs are unlikely to occur with the random generation. We analyze these chains using the handcrafted model from Appendix E.4. Alternatively, we could use the RandomSCC guideline.
- NumNonSingleton, SmallestSCCNonSing, AvgSCCNonSing%: Since the many singleton (i.e., trivial one state) SCCs make the features relating to all SCCs hard to interpret, these features analyze only the non-trivial SCCs. In REAL, in three quarters of the models there are very few non-trivial SCCs. From the SmallestSCCNonSing (again with the axis multiplied with a million) and AvgSCCNonSing% plot (which looks very similar to that for biggest overall SCC) we can see that in half the models even non-trivial SCCs are small. However, there are models where even the smallest SCC is big, indicating that there is one big SCC with chains of transient states around it. This is also the structure of the models from RANDOM, where the singleton smallest and average size plot are the same as the overall biggest SCC plot.

E Additional details on the experimental evaluation

E.1 Details on the optimizations

In this appendix, we describe the results of our evaluation of the different optimizations. In principle, we can enhance every VI algorithm with any combination of them (with the exception that combining T and PT is the same as just having PT).

- G: The Gauß-Seidel variant of VI does not perform the Bellman update on all states at once, but rather proceeds state by state. This allows to immediately use the new estimates of the states that were updated before and can thus speed up the computation.
- D: Deflating is a costly operation, since computing the SECs always requires a MEC decomposition. While this is possible in polynomial time [14], it is slow in practice. Thus, delaying deflating and only applying it only every n steps can speed up the computation. However, for high n , the algorithm can also waste time waiting for the next deflation to happen. Deflating every 100 steps was observed to be a good compromise in [18]. However, for OVI, we omitted this optimization, as there the whole point of the verification phase is to decrease the upper bound, and there is not point in waiting (for e.g., the lower bound to converge).
- T: The topological variant of VI [16] first computes a partitioning of the state space into SCCs. Instead of solving the whole SG at once, topological VI proceeds SCC by SCC, starting from those at the bottom of the topological ordering. This allows for using less memory at once and can speed up the computation.
- PT: The precise variant of topological VI which we introduced in Section 4. For a description of how to read scatter plots, see Section 6.3.

For assessing the impact of the G, D and T optimizations on the algorithms' performance, we compare the optimized versions of BVI, OVI, and SI to their unoptimized ones, as shown in the scatter plots (Figures 7, 8 and 9). PT is analyzed in Section 6.3.

G optimization for BVI and OVI : Figure 7 indicates that using the Gauß-Seidel optimization may reduce both the time required to solve a model by 1-4 times in most cases. The Gauß-Seidel optimization is slower in some cases because the values are computed sequentially to enable the use of the already computed results. Also, the unoptimized version uses vector operations instead, which turn out to be faster sometimes. Furthermore, it is possible that there are more iterations required to solve. This is because OVI and BVI may find different ECs to deflate depending on whether Gauß-Seidel is used or not. In some cases, the unoptimized version is able to find more favourable sets of ECs and thus requires less iterations. Note that the order in which the states are updated by Gauß-Seidel VI is random. However, changing the order of computation of the states to be a reverse topological enumeration of the states did not yield improvements in our experiments.

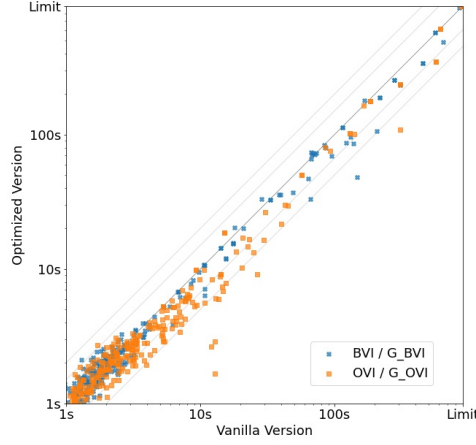


Fig. 7: Scatter plot with logarithmic scale, comparing verification times of vanilla (unoptimized) BVI, OVI and with their Gauß-Seidel variant. Below the diagonal means with G is better, above means vanilla is better.

D optimization for BVI :

Figure 8 clearly indicates that although DBVI may sometimes solve models faster, for the majority of our models it could not compete with BVI.

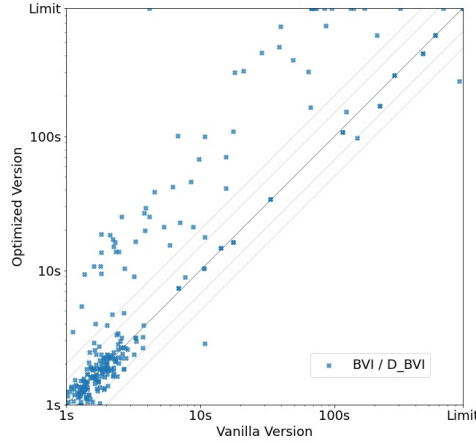


Fig. 8: Scatter plot with logarithmic scale, comparing verification times of vanilla (unoptimized) BVI, OVI and SI_{LP} with their deflating variant. Below the diagonal means with D is better, above means vanilla is better.

T optimization for BVI, OVI and SI_{LP} :

As the scatter plot in Figure 9 shows, on both types of models, the topological addition to SI neither increases nor decreases its performance considerably. However, most models have very few SCCs, so the topological optimization does not contribute a lot. The data point where SI is significantly faster is on the real case study "dice", where every state is an SCC on its own. Obviously, this is the best-case scenario for topological algorithms.

Adding the topological optimization to BVI often makes it significantly worse, mainly because of the precision issues addressed, when introducing PTBVI in Section 4. For OVI, the picture is slightly more mixed, but using the topological variant often helps.

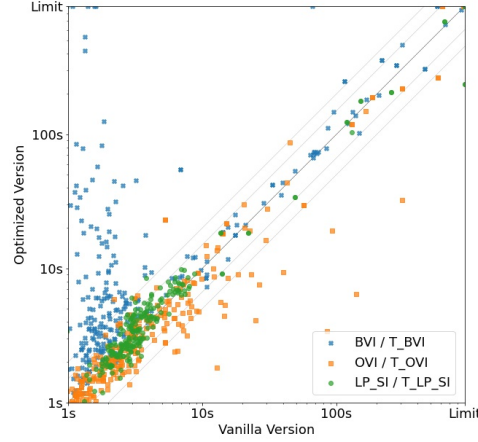


Fig. 9: Scatter plot with logarithmic scale, comparing verification times of vanilla (unoptimized) BVI, OVI and SI with their topological variant. Below the diagonal means with T is better, above means vanilla is better.

E.2 Additional scatter plot for BVI compared to WP and OVI

Figure 10 shows the scatter plots comparing BVI with WP and OVI, complementing those in Figure 4. These two figures give all relations between the three algorithms.

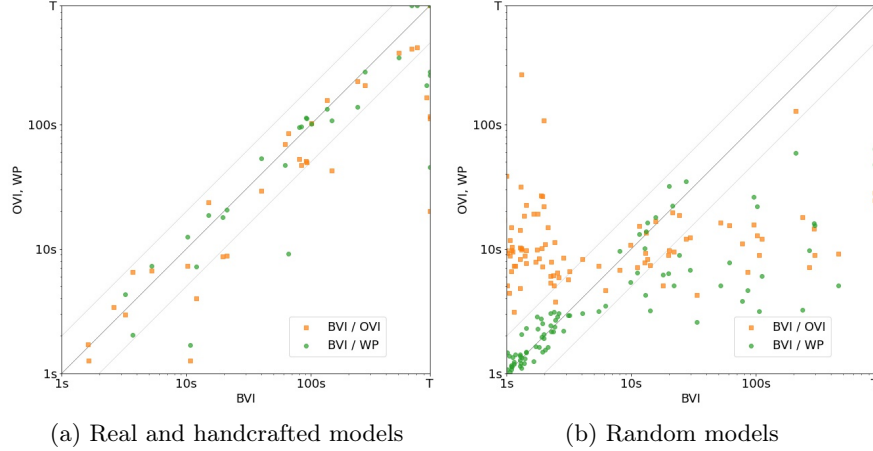


Fig. 10: BVI compared to WP and OVI.

E.3 Additional details and plots for OVI

Handcrafted models for OVI. As an extreme example where the lower bound converges faster, consider a chain of Maximizer states, where every state has two

actions: One that immediately yields a value of 0.5 and one that continues to the next state in the chain with low chance (0.01) and self-loops with a high chance (0.99). These kind of self loops always slow down VI algorithms, see also the haddad-monmege model [20]. The last state in the chain has an action with value 0.49. The lower bound of all states can immediately be set to 0.5, as all states have the first action that guarantees this value. However, to ensure that continuing in the chain is certainly the worse option, i.e., to have a convergent upper bound, BVI has to wait for the information to propagate over the self-loops. In contrast, OVI quickly knows the correct lower bound and then can verify it. Concretely, BVI cannot solve the model within 2 minutes as soon as the chain has more than 200 states, while OVI can deal with more than 10,000 states, see Figure 11.

As noted in Section 6.4, the dual situation where the lower bound converges slowly, is problematic for OVI. For example, if we remove the action with the value of 0.5 from the game above, we get a Markov chain where every state has a low chance of progressing towards the target⁷. On this type of chain, OVI is consistently around 4 times slower than BVI, as it is more time expensive to wait for the lower bound to converge than to converge from both sides.

Number of verification phases Theoretically, in the worst case the number of iterations for the first verification phase is $\frac{1}{\varepsilon'}$ (which is 1,000,000 in our case). If the verification fails, ε' is halved and number of iterations in the next verification phase doubles. Thus, it is possible that the verification phase is aborted, although the game is almost solved. OVI will iterate unnecessary extra steps until it reaches its new precision that verifies that OVI may indeed terminate. For large models, it is more likely that a model is so complex that it requires multiple verification phases. Thus, it is also more probable that the precision is halved at a point that will lead to unnecessary iterations. In addition to that, iterations are more costly for larger models since the state space is large.

⁷ Note that this is different from the completely adversarial example of [20], where every state does not only have self-loop, but even goes back to the initial state.

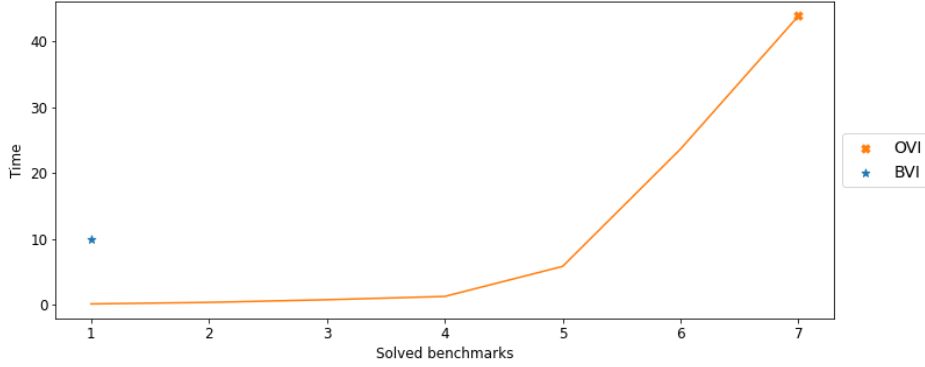


Fig. 11: Measuring number of handcrafted OVI-model solved against aggregated runtime. The legend sorts the algorithms by their aggregated runtimes in descending order. The model is always the same, but with 100, 500, 900, 1000, 5000, 9000 and 10000 states. OVI is fast on this model, while BVI without optimizations cannot solve the model as soon as it has over 200 states. Topological optimizations handle the model very well, since every state is an SCC.

Choice of parameters Two important parameters are the number of steps in a verification phase and the modification of the precision after a verification phase.

If a verification phase cannot verify the upper bound, we should abort it as soon as possible. However, it might be necessary to iterate for a long number of steps before the decrease of the upper bound propagates to all states, since it might be “hidden” behind some state switching its action. Thus, we may also not choose it too small.

After a failed verification phase, we have to increase the lower bound enough that trying to guess the next upper bound has good chances of being correct. However, if the precision becomes too small, the VI-phase might have to run for a very long time before another verification phase is attempted. On the other hand, keeping the precision too similar results in lots of aborted verification phases.

Finding a good trade-off, possibly even by dynamically changing these parameters during a run of the algorithm, is a task we leave for future work.

E.4 Analysis of large models

To analyze how the algorithms scale on a large model with many SCCs, we handcrafted a model, called `simple_n_m_SCC`. It contains n states and m SCCs. Every SCC forms a tree. The inner nodes have deterministic actions, leading to the next level. The leaves have probabilistic actions, leading to the root of the current or of the next tree, making the tree strongly connected.

Table 1 shows the number of states and SCCs (visible in the name of the model) and the resulting verification times. We compared PTBVI with the fastest approximate algorithm, WP, and two variants of topological strategy iteration: one using linear programming for solving the opponent MDP (TSI_{LP}) and one using strategy iteration for solving the opponent MDP (TSI_{SI}).

We see that PTBVI scales well on these models, while both variants of SI struggle time out on models with more than 2 million states.

Model name	PTBVI	TSI _{LP}	TSI _{SI}	WP
simple_50000_1_SCC	2.154	5.773	7.079	2.135
simple_50000_5_SCC	4.564	4.999		10.548
simple_100000_1_SCC	6.049	13.189	14.694	4.33
simple_100000_5_SCC	13.179	7.889	12.927	25.34
simple_500000_1_SCC	27.422	133.114		22.127
simple_500000_5_SCC	73.477	64.124		160.747
simple_1000000_1_SCC	53.707	274.134		53.319
simple_1000000_5_SCC	156.327	138.902		401.318
simple_2000000_1_SCC	103.149			134.549
simple_2000000_5_SCC	263.625	366.236		753.834
simple_3000000_1_SCC	165.242			227.212
simple_3000000_5_SCC	389.154			1244.764
simple_4000000_1_SCC				190.158
simple_4000000_5_SCC	470.012			
simple_5000000_1_SCC				301.517
simple_5000000_5_SCC	589.712			1723.37
simple_10000000_1_SCC	581.999			804.652
simple_10000000_5_SCC	1433.221			

Table 1: Performance comparison of PTBVI against TSI_{LP}, TSI_{SI}, WP