

Stopping Criteria for Value Iteration on Concurrent Stochastic Reachability and Safety Games

Marta Grobelna^{*✉}, Jan Křetínský^{†✉}, Maximilian Weininger^{‡§✉}

^{*}Technical University of Munich
Munich, Germany
marta.grobelna@tum.de

[†]Masaryk University
Brno, Czech Republic
jan.kretinsky@fi.muni.cz

[‡]Ruhr-University Bochum
Bochum, Germany
maximilian.weininger@rub.de

[§]Institute of Science and
Technology Austria
Klosterneuburg, Austria

Abstract—We consider two-player zero-sum concurrent stochastic games (CSGs) played on graphs with reachability and safety objectives. These include degenerate classes such as Markov decision processes or turn-based stochastic games, which can be solved by linear or quadratic programming; however, in practice, value iteration (VI) outperforms the other approaches and is the most implemented method. Similarly, for CSGs, this practical performance makes VI an attractive alternative to the standard theoretical solution via the existential theory of reals.

VI starts with an under-approximation of the sought values for each state and iteratively updates them, traditionally terminating once two consecutive approximations are ϵ -close. However, this stopping criterion lacks guarantees on the precision of the approximation, which is the goal of this work. We provide *bounded* (a.k.a. interval) VI for CSGs: it complements standard VI with a converging sequence of over-approximations and terminates once the over- and under-approximations are ϵ -close.

Index Terms—Formal methods, foundations of probabilistic systems and games, verification, model checking

I. INTRODUCTION

Concurrent stochastic games (CSGs, e.g., [14]): We consider two-player zero-sum games played on a graph. Every vertex represents a *state*. Edges are directed, originating from one state and leading to one or several other states. An edge is associated with a *probability distribution* over the successor states. A *play* proceeds through the graph as follows: starting from an initial state, both players simultaneously and independently choose an *action*, determining the edge to follow. Then, the next state is sampled according to the probability distribution, and the process is repeated in the successor.

We focus on infinite-horizon reachability and safety objectives [14]. The goal of reachability is to maximize the probability of reaching a given goal state. In contrast, the safety objective aims to maximize the probability of staying within a given set of states. The two objectives are dual, as instead of maximizing the probability of reaching a set of target states, one can minimize the probability of staying within the set of non-target states. Thus, we refer to both types collectively as CSGs. Popular subclasses of CSGs include turn-based stochastic games (TSGs), where the players make

decisions in turns, or Markov decision processes (MDPs), which involve only one player.

The value problem: In CSGs, memoryless (a.k.a. stationary) strategies suffice for both players, meaning they yield the same supremum probability as history-dependent strategies. However, unlike in TSGs, the strategies for CSGs require *randomization*, meaning players choose distributions over actions rather than single actions. Additionally, while the safety objective player, Player \mathcal{S} , can attain optimal strategies [39], the reachability objective player, Player \mathcal{R} , only possesses ϵ -optimal strategies for a given $\epsilon > 0$ [21]. As a result, the problem of deciding whether the supremum probability (a.k.a. *the value*) is at least p for $p \in [0, 1]$, is thus more subtle than for the mentioned subclasses. While for MDPs the value problem is in P, and for TSGs it is known to be in $\text{NP} \cap \text{co-NP}$, for CSGs it can be elegantly encoded into the *existential theory of reals* (ETR), which is only known to be decidable in PSPACE (although not known to be complete for it) [20]. Unfortunately, algorithms for ETR are practically even worse than the more general, doubly exponential methods for the first-order theory of reals [40]. “Finding a practical algorithm remains a very interesting open problem” [25].

Practical approximation: We focus on algorithms *approximating* the value with a predefined precision $\epsilon > 0$. Both for MDPs and TSGs, dynamic programming techniques such as *value iteration* (VI) or *strategy iteration* (SI) are practically more efficient than mathematical programming (linear or quadratic, respectively) [27], [30]. Thus, VI algorithms are prevalently used and implemented in popular tools such as PRISM-GAMES [33], motivating the focus on VI here.

Problem and our contribution: In VI, the lowest possible value is initially assigned to each state and then iteratively improved, computing an under-approximation of the value, converging to it in the limit. The algorithm (in practical implementations) terminates once two consecutive approximations are ϵ -close. However, the result can then be arbitrarily imprecise [23]. In this work, we introduce *bounded value iteration for CSGs*, following its previous success for MDPs [1], [5] or TSGs [18]. Its main idea is to enhance standard VI by introducing an over-approximation of the values computed in parallel with the under-approximation. Once the upper and lower bounds are ϵ -close, VI terminates, ensuring that the true value is at most ϵ away from the obtained approximation. Since the naïve formulation of an upper bound does not converge to the value

This research was funded in part by the German Research Foundation (DFG) project 427755713 GOPro, the MUNI Award in Science and Humanities (MUNI/I/1757/2021) of the Grant Agency of Masaryk University, the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 101034413, and the ERC Starting Grant DEUCE (101077178).

in general, previous approaches, notably including [7], have attempted to fix this but have failed. In this paper, we finally provide a valid solution.

Technical challenge: The fundamental technical difficulties arise from the following. The non-convergence of upper-bound approximations is primarily due to cyclic components, so-called *end components* (ECs), see [14], [18], [29]. Notably, non-convergence, for this reason, is an issue present already in MDPs and TSGs; see [18], [29]. Solutions have been developed over the past decade for these two subclasses. Indeed, for MDPs with reachability objective, these end components can effectively be removed from the graph without altering the value [6], [24]. Other objectives were considered in [1], [3]. TSGs with reachability objectives already require more careful analysis, decomposing the end components into sub-parts, so-called *simple ECs* [18]. A comprehensive framework for various quantitative objectives was proposed in [29]. Unfortunately, the idea of simple ECs is not easily extendable from TSGs to CSGs due to the absence of optimal strategies.

Summary of our contribution: We provide a stopping criterion for VI on CSGs, solving an open problem with erroneous solution attempts in the literature (see the related work in Subsec. I-A below). To this end, we unravel the recursive hierarchical structure of end components in CSGs (see Rem. 31) and adapt the bounded VI algorithm.

A. Related Work

Available approaches: The PSPACE-algorithm introduced in [20] for deciding whether the value of a given game is at least p , for $p \in [0, 1]$, allows for a trivial stopping criterion by iteratively executing this algorithm for a suitable sequence of $(p_i)_{i \in \mathbb{N}}$ (intuitively, we choose p_i such that alternatingly, the value of the game is above and below, while the distance between two consecutive p_i 's monotonically decreases). However, this criterion is impractical since it uses the existential theory of reals [20]. The best known complexity upper bound comes from [22] and states that the problem of approximating the value of a CSG is in TFNP[NP], i.e. total function from NP with an oracle for NP. However, the proposed algorithm is not practical, as it relies on guessing a floating point representation of the value and optimal strategies for both players. A recursive bisection algorithm was introduced in [26], which is also impractical as its time complexity is best-case doubly exponential in the number of states. For the algorithms commonly employed for in the non-concurrent case, VI and SI, [25] provide doubly exponential lower and upper bounds on the number of iterations that VI requires in the worst-case for computing an ε -approximation. Their counter-example uses a CSG where all states have value 1. Thus, the worst-case complexity of our approach is the same since an additional over-approximation does not speed up convergence in this example. Nonetheless, these results are worst-case bounds, i.e. they hold a priori for all games; earlier termination is possible, but necessarily requires a stopping criterion, which has so far been elusive. Finally, in [38], an algorithm is provided that, unlike all other known algorithms, only has a single-

exponential dependency on the number of states. A practical comparison of our value iteration and [38] is an interesting future step, as better worst-case complexity of an algorithm need not translate to better practical performance on typical instances; for example, in MDPs, worst-case exponential VI and SI typically outperform the polynomial approach of linear programming [27].

Previous attempts at stopping criterion: A stopping criterion for SI and VI on CSGs was first presented in [7], but later found to contain an irreparable mistake [10]. Specifically, the algorithm returned over-approximations smaller than the actual values in certain situations, as detailed in [10]. We analyze the counter-example from [10] in App. D-A Later, [19] proposed a stopping criterion for VI, which also contains a fundamental flaw: it fails to converge for CSGs with ECs. We analyze the counter-example to this approach in App. D-B. *Our work thus delivers the first stopping criterion in this context.*

Further directions of related work: Variants of CSGs have appeared very early, under the names of Everett, Gillette, or Shapley game. See [26] for an explanation of all game types, their relations, and algorithms to solve them. These games also consider discounted payoff or limit-average payoff, generalizing the reachability and safety CSGs we consider here. A generalization of CSGs to ω -regular objectives has been considered in [8], [15]. An insightful characterization of optimal strategies in concurrent games with various objectives can be found in [4].

II. PRELIMINARIES

A. Concurrent Stochastic Games

Probability Distributions: For a countable set X , a function $\mu: X \rightarrow [0, 1]$ is called a *distribution* over X if $\sum_{x \in X} \mu(x) = 1$. The *support* of μ is $\text{Supp}(\mu) := \{x \mid \mu(x) > 0\}$. The set of all distributions over X is denoted by $\text{Dist}(X)$.

Concurrent Stochastic Games: A *concurrent stochastic game* (CSG) [16] is a tuple $G := (S, \mathcal{A}, \Gamma_{\mathcal{R}}, \Gamma_{\mathcal{S}}, \delta, s_0, T)$, where S is a finite set of *states*, $\mathcal{A} := A \times B$ is a finite set of *actions* with $A := \{a_1, \dots, a_l\}$ and $B := \{b_1, \dots, b_m\}$ the sets of actions available for player \mathcal{R} and \mathcal{S} , respectively, $\Gamma_{\mathcal{R}}: S \rightarrow (2^A \setminus \emptyset)$ and $\Gamma_{\mathcal{S}}: S \rightarrow (2^B \setminus \emptyset)$ are two *enabled actions* assignments and $\delta: S \times A \times B \rightarrow \text{Dist}(S)$ is a *transition function*, where $\delta(s, a, b)(s')$ gives the *probability of a transition* from state s to state s' when player \mathcal{R} chooses action $a \in \Gamma_{\mathcal{R}}(s)$ and player \mathcal{S} action $b \in \Gamma_{\mathcal{S}}(s)$, $s_0 \in S$ is an *initial state*, and $T \subseteq S$ is a set of *target states*. A **CSG** is *turn-based* if for every state s only one player has a meaningful choice, i.e. either $\Gamma_{\mathcal{R}}(s)$ or $\Gamma_{\mathcal{S}}(s)$ is a singleton; we call such game a *turn-based stochastic game* (TSG).

Example 1 (CSGs). Consider the **CSG** Hide-Run-or-Slip depicted in Fig. 1 (an adaption of the Hide-or-Run game in [16], [21], [31]). Circles represent states and black dots depict a probabilistic transition with uniform distribution. Each edge is labeled with a pair of actions, the left for player \mathcal{R} and the right for player \mathcal{S} ; \square is a placeholder for an arbitrary action. We have $S := \{s_{\text{hide}}, s_{\text{home}}, s_{\text{wet}}\}$, with $\Gamma_{\mathcal{R}}(s_{\text{hide}}) := \{\text{hide}, \text{run}\}$

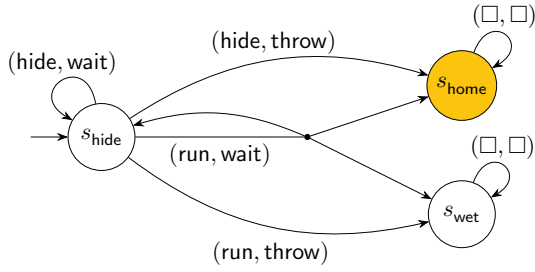


Fig. 1: Example CSG called Hide-Run-or-Slip.

and $\Gamma_{\mathcal{S}}(s_0) := \{\text{wait, throw}\}$, s_0 is the initial state denoted by the arrow with no predecessor state, and $T := \{s_{\text{home}}\}$. The game has the following intuitive interpretation: Player \mathcal{R} wants to get home without getting wet. Player \mathcal{S} has a single snowball and can make player \mathcal{R} wet by throwing it at player \mathcal{R} . If player \mathcal{R} runs and player \mathcal{S} throws the ball, player \mathcal{R} gets wet. If player \mathcal{R} runs but player \mathcal{S} waits, with a probability of $\frac{1}{3}$ the player reaches home or the player slips instead and with a probability of $\frac{1}{3}$ does not move at all or with a probability of $\frac{1}{3}$ falls on the ground and gets wet. \triangle

Plays: A play π of a CSG G is an infinite sequence of states $s_0 s_1 s_2 \dots$, such that for all $i \in \mathbb{N}$ there are actions $a \in \Gamma_{\mathcal{R}}(s_i)$ and $b \in \Gamma_{\mathcal{S}}(s_i)$ with $\delta(s_i, a, b)(s_{i+1}) > 0$. $\text{Play}(G)$ is the set of all plays and $\text{Play}_s(G)$ the set of all plays $s_0 s_1 s_2 \dots$ with $s_0 = s$.

Strategies: A strategy for player \mathcal{R} (or \mathcal{S}) is a function $\rho: S \rightarrow \text{Dist}(A)$ (or $\sigma: S \rightarrow \text{Dist}(B)$) that assigns a distribution over actions available to player \mathcal{R} (or \mathcal{S}) to each state, i.e., for all $s \in S$, $\text{Supp}(\rho(s)) \subseteq \Gamma_{\mathcal{R}}(s)$ (or $\text{Supp}(\sigma(s)) \subseteq \Gamma_{\mathcal{S}}(s)$).¹ We call a Player \mathcal{R} (or Player \mathcal{S}) strategy ρ *pure* if all distributions it returns are Dirac distributions, i.e., at each $s \in S$ we have a unique action $a \in \Gamma_{\mathcal{R}}(s)$ (or $b \in \Gamma_{\mathcal{S}}(s)$) such that $\rho(s)(a) = 1$ (or $\sigma(s)(b) = 1$). Otherwise, the strategy is *mixed*. For player \mathcal{R} (or \mathcal{S}) we denote the set of strategies by \mathcal{R} (or \mathcal{S}) and a single strategy by ρ (or σ).

Markov Decision Processes: Given a CSG G , if we fix a strategy $\rho \in \mathcal{R}$ of player \mathcal{R} , the game becomes a \mathcal{S} -Markov Decision Process (MDP, [41]) G_ρ with the transition function

$$\delta_\rho(s, b)(s') := \sum_{a \in \Gamma_{\mathcal{R}}(s)} \delta(s, a, b)(s') \cdot \rho(s)(a),$$

for all $s \in S$ and $a \in \Gamma_{\mathcal{R}}(s)$. The MDP induced by a fixed strategy $\sigma \in \mathcal{S}$ is defined analogously.

Markov Chains: Similarly, by fixing a pair of strategies $(\rho, \sigma) \in \mathcal{R} \times \mathcal{S}$, we obtain a Markov chain $G_{\rho, \sigma}$ with the same state space S , the initial state s_0 , and the transition probabilities P given by

$$\delta_{\rho, \sigma}(s)(s') := \sum_{(a, b) \in A} \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b).$$

¹Since memoryless strategies are sufficient for the objectives considered in this paper, we do not introduce general history-dependent strategies to avoid clutter. We refer to [16] for more details.

Thus, a pair of strategies (ρ, σ) induces a unique probability measure $\mathbb{P}_{s_0}^{\rho, \sigma}$ over plays in the Markov chain as usual, see [2, Chap. 10.1], where the set of paths starting in s_0 has measure 1.

Objectives: We partition S into T , denoting the set of states player \mathcal{R} wants to reach, and $F := S \setminus T$, denoting the set of states player \mathcal{S} wants to confine the game in. We denote the *reachability* objective by $\Diamond T := \{s_0 s_1 s_2 \dots \mid \exists i \in \mathbb{N} : s_i \in T\}$ and the *safety* objective by $\Box F := \{s_0 s_1 s_2 \dots \mid \forall i \in \mathbb{N} : s_i \in F\}$. The *value* of the objective $\Diamond T$, i.e. $V_{\mathcal{R}}(s)$, and the objective $\Box F$, i.e. $V_{\mathcal{S}}(s)$, at state s are given by

$$V_{\mathcal{R}}(s) := \sup_{\rho \in \mathcal{R}} \inf_{\sigma \in \mathcal{S}} \mathbb{P}_s^{\rho, \sigma}(\Diamond T)$$

$$V_{\mathcal{S}}(s) := \sup_{\sigma \in \mathcal{S}} \inf_{\rho \in \mathcal{R}} \mathbb{P}_s^{\rho, \sigma}(\Box F).$$

By the determinacy of CSGs and the duality of these objectives [21], [34], it holds that $V_{\mathcal{R}}(s) + V_{\mathcal{S}}(s) = 1$. Consequently, the task of approximating $V_{\mathcal{R}}(s)$ with a given precision is equivalent to approximating $V_{\mathcal{S}}(s)$. Further, the objective of minimizing the reachability for T is equivalent to the objective of maximizing safety for F for the *same* player. Consequently, in the following we only focus on maximizing reachability as both minimization and the safety objectives can be reduced to it.

Example 2 (Optimal Strategies Need Not Exist). In CSGs, an optimal strategy for player \mathcal{R} might not exist [31], meaning that at some states, the value is attainable only in the limit. Consider our running example from Fig. 1, Hide-Run-or-Slip. Assume for the moment that there is no chance of slipping, i.e. upon playing run and wait, the target state is reached. To win, Player \mathcal{R} has to run eventually. However, Player \mathcal{S} can utilize a strategy that throws with positive probability at all points in time. Thus, Player \mathcal{R} cannot win almost surely.

However, Player \mathcal{R} has the possibility of *limit-sure winning* in [16]: By running with vanishingly low probability ε in every round, the probability of winning is $1 - \varepsilon$. This is because Player \mathcal{S} has the highest probability ε of hitting Player \mathcal{R} by throwing in the first round; throwing in a later round n only has hitting probability ε^n . For any $\varepsilon > 0$, this strategy of Player \mathcal{R} achieves $1 - \varepsilon$. The value, being the supremum over all strategies, is 1.

This notion of obtaining a value only in the limit is not restricted to sure winning: By adding the chance of slipping, the value of the game becomes 0.5. However, by the same argument as above, Player \mathcal{R} cannot win with probability 0.5, but only with a probability $0.5 - \varepsilon$ for all $\varepsilon > 0$. \triangle

For $s \in S$, $a \in \Gamma_{\mathcal{R}}(s)$ and $b \in \Gamma_{\mathcal{S}}(s)$, the *set of potential successors* of s is denoted by $\text{Post}(s, a, b) := \text{Supp}(\delta(s, a, b))$. We lift the notation to strategies $\rho \in \mathcal{R}$ and $\sigma \in \mathcal{S}$ by

$$\text{Post}(s, \rho, \sigma) = \bigcup_{a \in \text{Supp}(\rho(s))} \bigcup_{b \in \text{Supp}(\sigma(s))} \text{Post}(s, a, b).$$

We denote by $W_{\mathcal{S}} := \{s \in S \mid V_{\mathcal{R}}(s) = 0\}$ the *sure winning region* of player \mathcal{S} . It can be computed in at most $|S|$ steps by iteration $W_{\mathcal{S}}^0 := (S \setminus T)$ and $W_{\mathcal{S}}^{k+1} := \{s \in S \setminus T \mid$

$\exists b \in \Gamma_{\mathcal{S}}(s) : \forall a \in \Gamma_{\mathcal{R}}(s) : \text{Post}(s, a, b) \subseteq W_{\mathcal{S}}^k$ for all $k \in \mathbb{N}$ [15]. Consequently, we can assume without loss of generality that T and $W_{\mathcal{S}}$ are both singletons and absorbing.

Example 3 (The Sets T , F , and $W_{\mathcal{S}}$). In Fig. 1, player \mathcal{R} wants to reach $T = \{s_1\}$, while player \mathcal{S} aims to stay in $F = \{s_0, s_2\}$. Since s_2 is absorbing, $W_{\mathcal{S}} = \{s_2\}$. \triangle

Matrix Games: At each state of a CSG the players \mathcal{R} and \mathcal{S} play a two-player zero-sum matrix game [35], [43]. In general, a matrix game is a tuple $Z := (N, A, u)$ [43] where, $N := \{1, \dots, n\}$ is a finite set of players, $A := \{\alpha_1, \dots, \alpha_m\}$ is a finite set of actions available to each player, and $u : A \rightarrow \mathbb{Q}$ is a utility function. In CSGs, the matrix game played at a specific state s can be represented by a matrix $Z_{V_{\mathcal{R}}}(s) \in \mathbb{Q}^{l \times m}$, where $\Gamma_{\mathcal{R}}(s) = \{a_1, \dots, a_l\}$ and $\Gamma_{\mathcal{S}}(s) = \{b_1, \dots, b_m\}$. The entries of the matrix correspond to the utility, i.e., the value attainable upon choosing a pair of actions $(a_i, b_j) \in A$. Thus, the i -th row and the j -th column is given by $Z_{V_{\mathcal{R}}}(s)(i, j) := \sum_{s' \in S} \delta(s, a_i, b_j)(s') \cdot V_{\mathcal{R}}(s')$.

Example 4 (Matrix Game). Consider the CSG in Fig. 1. The matrix game played at state s_{hide} is given by the following matrix.

$$Z_{V_{\mathcal{R}}}(s_{\text{hide}}) = \begin{pmatrix} \text{throw} & & \text{wait} \\ 0 & \frac{1}{3} \cdot V_{\mathcal{R}}(s_{\text{hide}}) + \frac{1}{3} & \\ 1 & V_{\mathcal{R}}(s_{\text{hide}}) & \end{pmatrix} \begin{matrix} \text{run} \\ \text{hide} \end{matrix} \quad (1)$$

Player \mathcal{R} is the so called *row player* while Player \mathcal{S} is the *column player*. \triangle

In a matrix game, a player's *strategy* is a distribution over the available actions at a specific state. To distinguish between strategies of a CSG and strategies of a matrix game, we refer to strategies of a matrix game as *local strategies* and strategies of a CSG as *global strategies*. The set of all local strategies at a state s is denoted by $\mathcal{R}(s)$ or $\mathcal{S}(s)$ for player \mathcal{R} or \mathcal{S} , respectively. The existence of optimal (local) strategies in a matrix game for both players is guaranteed by Nash's Theorem [36], [37]. The payoff that is attainable with an optimal local strategy is called *value* that we denote by $V(Z_{V_{\mathcal{R}}})$ for a matrix game $Z_{V_{\mathcal{R}}}$. It can be calculated using linear programming (e.g., [28], see App. A-A).

End Components: A non-empty set of states $C \subseteq S$ is called an *end component* (EC) if (i) there exists a pair of strategies $(\rho, \sigma) \in \mathcal{R} \times \mathcal{S}$ such that for each $s \in C$ it holds that $\text{Post}(s, \rho, \sigma) \subseteq C$; and (ii) for every pair of states $s, s' \in C$ there is a play $s_0 s_1 \dots$ such that $s_0 = s$ and $s_n = s'$ for some n , and for all $0 \leq i < n$, it holds $s_i \in C$ and $s_{i+1} \in \text{Post}(s_i, \rho, \sigma)$.

Intuitively, an EC is a set of states where a play can stay forever under some pair of strategies. In other words, the players can cooperate to keep the play inside the EC (this is the usual way to lift the definition of [14] from MDP to games). Thus, we can compute ECs in a CSG by computing ECs in the corresponding MDP with both players unified, i.e., every pair of actions is interpreted as an action in the MDP. Efficient algorithms for this exist [2], [12], [46]. An EC C is

called *inclusion maximal* (short maximal) if there exists no EC C' such that $C \subsetneq C'$.

B. Value Iteration

Value iteration (VI, e.g. [11]) assigns an initial value estimate to each state and then iteratively updates it. In classical VI, which approximates the reachability value from below, the initial estimates are 1 for states in T and below the actual value otherwise, e.g. 0. Each iteration backpropagates the estimate by maximizing the expectation of the value player \mathcal{R} can ensure with respect to the previous estimate.

Formally, we capture estimates as valuations, where a *valuation* $v : S \rightarrow [0, 1]$ is a function mapping each state s to a real number representing the (approximate or true) value of the state. For two valuations v, v' , we write $v \leq v'$ if $v(s) \leq v'(s)$ for every $s \in S$.

To compute the expected value at a state s , the matrix game $Z_v(s)$ has to be solved, meaning its value, $V(Z_v(s))$, has to be estimated. This computation is, especially in the turn-based setting, also referred to as *Bellman update*. Formally,

$$V(Z_v(s)) := \mathfrak{B}(v)(s) := \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(v)(s, \rho, \sigma),$$

where $\mathfrak{B}(v)(s, \rho, \sigma) :=$

$$\sum_{(a,b) \in A} \sum_{s' \in S} \rho(a) \cdot \sigma(b) \cdot \delta(s, a, b)(s') \cdot v(s').$$

Convergent Under-approximation: We recall VI from below as in [9]: starting from the initial valuation L^0 , we perform the Bellman update on every state to obtain a new valuation. We denote by L^k the valuation obtained in the k -th iteration. Formally:

$$L^0(s) := \begin{cases} 1, & \text{if } s \in T; \\ 0, & \text{else,} \end{cases} \quad L^{k+1}(s) := \mathfrak{B}(L^k)(s). \quad (2)$$

Since T and $W_{\mathcal{S}}$ are absorbing, for all $k \in \mathbb{N}$ we have $L^k(s) = 1$ for all $s \in T$, and $L^k(s) = 0$ for all $s \in W_{\mathcal{S}}$. The updated valuation, i.e. $L^{k+1}(s)$, is computed by solving the corresponding matrix game.

Theorem 5 (VI converges from below [17, Thm. 1]). *VI from below converges to the value, i.e. $\lim_{k \rightarrow \infty} L^k = V_{\mathcal{R}}$.*

Bounded Value Iteration: While Thm. 5 proves that VI from below converges in the limit, this limit may not be reached in finitely many steps and can be irrational [17]. Thus, we do not know when to stop the algorithm to guarantee certain precision of the approximation, as there is no practical bound how close any valuation L^k is to the actual value. We merely have the worst-case bound of [25]: Running for a number of iterations that is doubly-exponential in the number of states allows to conclude that the lower bound is ε -close to the value.

To obtain a practical stopping criterion, we use the approach of Bounded Value Iteration (BVI), shown in Alg. 1. In addition to the lower bound L , it maintains an upper bound on the value

Algorithm 1 Bounded value iteration procedure for CSGs.

```

1: Algorithm BVI(CSG  $G$ , threshold  $\varepsilon > 0$ )
2:    $W_{\mathcal{S}} \leftarrow \{s \in S \mid V_{\mathcal{R}}(s) = 0\}$   $\triangleright$  Winning region for  $\mathcal{S}$ 
3:    $L^0, U^0$  initialized by Eq. (1) and (2), respectively
4:    $\text{MEC} \leftarrow \text{FIND\_MECs}(G)$   $\triangleright$  Find all MECs in the game
5:    $k \leftarrow 0$ 
6:   repeat
7:     for  $s \in S$  do  $\triangleright$  Standard Bellman update of both bound
8:        $L^{k+1}(s) \leftarrow \mathfrak{B}(L^k)(s)$ 
9:        $U^{k+1}(s) \leftarrow \mathfrak{B}(U^k)(s)$ 
10:    for  $C \in \text{MEC}$  do
11:       $U^{k+1} \leftarrow \text{DEFLATE}(G, U^{k+1}, C)$ 
12:     $k \leftarrow k + 1$ 
13:  until  $U^{k+1} - L^{k+1} \leq \varepsilon$ 

```

U that is meant to converge to the value from above. Naïvely, this upper bound is defined as follows:

$$U^0(s) := \begin{cases} 0, & \text{if } s \in W_{\mathcal{S}}; \\ 1, & \text{else,} \end{cases} \quad U^{k+1}(s) := \mathfrak{B}(U^k)(s). \quad (3)$$

Given a precision $\varepsilon > 0$, the algorithm terminates once the under- and the over-approximations are ε -close, i.e., when both approximations are at most ε -away from the actual value. However, applying Bellman updates does not suffice for the over-approximation to converge in the presence of ECs, as the following example shows:

Example 6 (Non-convergent Over-approximations). Consider the CSG *Hide-Run-or-Slip* in Fig. 1. To compute $U^{k+1}(s_{\text{hide}})$ according to Eq. (3), in each iteration we solve the matrix game $Z_{U^k}(s_{\text{hide}})$ is given by Eq. (1) where the unknown $V_{\mathcal{R}}$ is replaced by U^k , i.e.:

$$Z_{U^k}(s_{\text{hide}}) = \begin{pmatrix} \text{throw} & \text{wait} \\ 0 & \frac{1}{3} \cdot U^k(s_{\text{hide}}) + \frac{1}{3} \\ 1 & U^k(s_{\text{hide}}) \end{pmatrix} \begin{matrix} \text{run} \\ \text{hide} \end{matrix} \quad (4)$$

Table I shows the updates of the lower and upper bounds, $L^k(s_{\text{hide}})$ and $U^k(s_{\text{hide}})$, respectively. While the lower bound converges to 0.5, the upper bound stays at 1. This is because for player \mathcal{R} , action *hide* always “promises” a valuation of 1 in the next step, as the lower row of the matrix yields 1 for all player \mathcal{S} strategies. \triangle

III. CONVERGENT OVER-APPROXIMATION: OVERVIEW

Here, we describe the structure of our solution. Ex. 6 shows that the naïve definition of the over-approximation need not

TABLE I: BVI for the game *Hide-Run-or-Slip* (Fig. 1), where the over-approximations do not converge.

k	0	1	2	...	∞
$L^k(s_{\text{hide}})$	0.0	0.25	0.36	...	0.5
$U^k(s_{\text{hide}})$	1.0	1.0	1.0	...	1.0

converge to the true value. In particular, in the presence of ECs, Bellman updates do not have a unique fixpoint. Thus, our goal is to define a function **DEFLATE** (usage highlighted in Alg. 1, definition in Alg. 3) that, intuitively, decreases the “bloated” upper bounds inside each EC to a realistic value substantiated by a value promised *outside* of this EC. Formally, we ensure that Alg. 1 produces a monotonically decreasing sequence of valuations (i) over-approximating the reachability value and (ii) converging to it in the limit. This idea has been successfully applied for TSGs [18], [29].

Remark 7 (Inflating for Safety). Since over-approximations need not converge for $V_{\mathcal{R}}$, dually under-approximations need not converge for $V_{\mathcal{S}}$ (as is the case in TSGs, see [29]). Thus, to directly solve a safety game, one needs an inflating operation dual to deflating. As described when introducing the objectives, we take the conceptually easier route of reducing everything to maximizing reachability objectives.

We proceed in two steps. First, in Sec. IV we prove that indeed ECs are the source of non-convergence, in particular what we call *Bloated End Components*. Intuitively, these are ECs where at each state both players prefer local strategies which all successor states belong to the EC. Second, in Sec. V, we define the **DEFLATE** algorithm, which essentially ensures that we focus on player \mathcal{R} strategies that do not make the game stuck in a EC but rather progress towards the target. To this end, we lift the notion of *best exit* from TSGs [18, Definition 3] to CSGs.

IV. THE CORE OF THE PROBLEM: CHARACTERIZING BLOATED END COMPONENTS

A locally optimal strategy of Player \mathcal{R} does not coincide with an optimal global strategy of Player \mathcal{R} because the latter must eventually leave ECs, while the former is under the illusion that staying is optimal. Thus, in this section, we want to find properties that local strategies (of both players) must fulfill in order to leave an EC in a way that is globally optimal. Thus, since this section mainly concerns *local* strategies, we use the word strategy to speak about local strategies and explicitly make clear when we talk about global ones.

Remark 8. Throughout the technical sections and the appendix, we always fix a CSG $G := (S, \mathcal{A}, \Gamma_{\mathcal{R}}, \Gamma_{\mathcal{S}}, \delta, s_0, T)$.

A. Convergence without ECs

As a first step, we prove that ECs are the only source of non-convergence, and without them, the naïve BVI using only Bellman updates converges.

Theorem 9 (Convergence without ECs — Proof in App. C-A). *Let G be a CSG where all ECs are trivial, i.e. for every EC C we have $C \subseteq W_{\mathcal{S}} \cup T$. Then, the over-approximation using only Eq. (3) converges, i.e. $\lim_{k \rightarrow \infty} U^k = V_{\mathcal{R}}$.*

Proof sketch. This proof is an extension of the proof of [18, Theorem 1] for turn-based games to the concurrent setting. The underlying idea is the same, and can be briefly summarized as follows: We assume towards a contradiction that $\lim_{k \rightarrow \infty} U^k =:$

TABLE II: **BVI** for the **CSG** in Ex. 10, where the over-approximations converge.

k	0	1	2	3	\dots	∞
$L^k(s_{\text{hide}})$	0.0	$\frac{1}{3}$	$\frac{4}{9}$	0.4815	\dots	0.5
$U^k(s_{\text{hide}})$	1.0	$\frac{2}{3}$	$\frac{5}{9}$	0.5185	\dots	0.5

$U^* \neq V_{\mathcal{R}}$, and find a set \mathcal{X} that maximizes the difference between upper bound and value. We show that every pair of strategies leaving the set \mathcal{X} decreases the difference $U^* - V_{\mathcal{R}}$. However, $V_{\mathcal{R}}$ and U^* are fixpoints of the Bellman update, from [17, Theorem 1] and Lem. 48, respectively. Consequently, optimal strategies need to remain in the set. However, in the absence of **ECs**, optimal strategies have to leave the set, which yields a contradiction and proves that $U^* = V_{\mathcal{R}}$.

The key difference to the proof of [18, Theorem 1] is that we cannot argue about actions anymore, but have to consider mixed strategies. This significantly complicates notation. Additionally, and more importantly, the former proof crucially relied on the fact that for a state of Player \mathcal{R} , we know that its valuation is at least as large as that of any action, and dually for a state of Player \mathcal{S} , its valuation is at most as large as that of any action. In the concurrent setting, this is not true. The optimal strategies need not be maximizing nor minimizing the valuation and, moreover, they can be maximizing for one valuation and minimizing for another. Thus, we found a more general, and in fact simpler, way of proving that “no state in \mathcal{X} can depend on the outside” [18, Statement 5] and deriving the contradiction. \square

Interestingly, not all **ECs** cause non-convergence of Eq. (3) as the following example illustrates.

*Example 10 (Unproblematic **EC**).* We modify the **CSG** Hide-Run-or-Slip (Fig. 1) such that the matrix game played at s_{hide} is $Z'_{U^k}(s_{\text{hide}})$ below.

$$Z'_{U^k}(s_{\text{hide}}) = \begin{pmatrix} \text{throw} & \text{wait} \\ \begin{pmatrix} 1 & \frac{1}{3} \cdot U^k(s_{\text{hide}}) + \frac{1}{3} \\ 0 & U^k(s_{\text{hide}}) \end{pmatrix} & \begin{pmatrix} \text{run} \\ \text{hide} \end{pmatrix} \end{pmatrix}$$

The difference is that $Z'_{U^k}(s_{\text{hide}})(\text{run}, \text{throw}) = 1$ and $Z'_{U^k}(s_{\text{hide}})(\text{hide}, \text{throw}) = 0$, switching the values as compared to the original **CSG** (see Eq. (4)). Here, both bounds converge to 0.5 despite the presence of the **EC** $\{s_{\text{hide}}\}$, as shown in Table II. \triangle

B. Towards Characterizing Bloated End Components

Intuition: In Ex. 10, the best strategy of Player \mathcal{R} leaves the **EC** almost surely against all counter-strategies of Player \mathcal{S} , and hence **BVI** converges. In contrast, in Ex. 6, the best strategy of Player \mathcal{R} is one where Player \mathcal{S} has a counter-strategy that forces the play to stay inside the **EC**; this causes non-convergence. Generalizing these ideas, we see that a problem occurs if Player \mathcal{R} has a strategy that is locally optimal but non-leaving, i.e. Player \mathcal{S} has a counter-strategy that keeps the play inside an **EC**.

Outline: We formalize these ideas in the following definitions: First, Def. 12 formalizes optimal (local) strategies using *weakly dominant strategies* in matrix games, extending the standard definition (e.g. [35]) to sets of strategies. This extension is not straightforward, and there are several technical intricacies that we comment on. Next, Def. 15 captures leaving and staying strategies. We differentiate strategies that are leaving (irrespective of the opponent’s strategy), staying (irrespective of the opponent’s strategy), and non-leaving (where there exists an opponent’s strategy that leads to staying) with respect to a given set of states. Based on this, we formally describe hazardous strategies in Def. 16, which are (locally) optimal strategies of Player \mathcal{R} that are non-leaving; additionally, to be problematic for convergence, they are better than all leaving strategies. Using these, we can precisely characterize the Bloated End Components (**BECs**, Def. 18) that cause non-convergence.

*Remark 11 (Additional Challenges Compared to **TSGs**).* The core problem is the same as in **TSGs**: Player \mathcal{R} is under the illusion that staying inside an **EC** yields a better valuation than leaving. However, in **TSGs**, the definitions of optimality and leaving are straightforward, since every state belongs to a single player and pure strategies are optimal; the definitions of hazardous and trapping strategies are not even necessary. In contrast, the definitions in **CSGs** are technically involved, as we have to take into account the interaction of the players and the possibility of optimal mixed strategies. In particular, [19] defined a straightforward extension of leaving based only on actions, not strategies. This is incorrect, as we demonstrate in App. D-B.

Definition 12 (Dominating Sets of Strategies). Let v be a valuation, $s \in S$ a state, $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}' \subseteq \mathcal{R}(s)$ and $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}' \subseteq \mathcal{S}(s)$ sets of local strategies. We now define two notions of domination for sets of strategies, namely *weak domination* and being *not worse*. Both of these depend on the player.

Definition for Player \mathcal{R} : We write $\mathcal{R}_2 \prec_{v, \mathcal{S}'} \mathcal{R}_1$ to denote that \mathcal{R}_1 *weakly dominates* \mathcal{R}_2 under the set of counter-strategies \mathcal{S}' with respect to v . Formally, $\exists \rho_1 \in \mathcal{R}_1. \forall \rho_2 \in \mathcal{R}_2$:

- (i) $\inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_2, \sigma) \leq \inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_1, \sigma)$, and
- (ii) $\exists \sigma' \in \mathcal{S}'$ such that $\mathfrak{B}(v)(s, \rho_2, \sigma') < \mathfrak{B}(v)(s, \rho_1, \sigma')$.

If only Condition (i) is satisfied, we write $\mathcal{R}_2 \preceq_{v, \mathcal{S}'} \mathcal{R}_1$ to denote that the set \mathcal{R}_1 is *not worse* than \mathcal{R}_2 under \mathcal{S}' with respect to v .

Definition for Player \mathcal{S} : Dually, we write $\mathcal{S}_2 \prec_{v, \mathcal{R}'} \mathcal{S}_1$ to denote that set \mathcal{S}_1 *weakly dominates* \mathcal{S}_2 under \mathcal{R}' with respect to v . Formally, $\exists \sigma_1 \in \mathcal{S}_1. \forall \sigma_2 \in \mathcal{S}_2$:

- (i) $\sup_{\rho \in \mathcal{R}'} \mathfrak{B}(v)(s, \rho, \sigma_2) \geq \sup_{\rho \in \mathcal{R}'} \mathfrak{B}(v)(s, \rho, \sigma_1)$, and
- (ii) $\exists \rho' \in \mathcal{R}'$ such that $\mathfrak{B}(v)(s, \rho', \sigma_2) > \mathfrak{B}(v)(s, \rho', \sigma_1)$.

If only Condition (i) is satisfied, we write $\mathcal{S}_2 \preceq_{v, \mathcal{R}'} \mathcal{S}_1$ to denote that the set \mathcal{S}_1 is *not worse* than \mathcal{S}_2 under \mathcal{R}' with respect to v .

Example 13 (Dominating Sets of Strategies). Consider the matrix game defined in Eq. (4) and the valuation $U^k(s_{\text{hide}}) =$

$U^k(s_{\text{home}}) = 1$ and $U^k(s_{\text{wet}}) = 0$. Here, for Player \mathcal{R} , the pure strategy $\{\text{hide} \mapsto 1\}$ dominates the pure strategy $\{\text{run} \mapsto 1\}$:

$$\{\text{run} \mapsto 1\} \prec_{U^k, \mathcal{S}(s_{\text{hide}})} \{\text{hide} \mapsto 1\}.$$

This is because when Player \mathcal{S} throws the ball, hiding yields 1 while running yields 0. Note that this is in fact independent of the valuation, so it also holds for $V_{\mathcal{R}}$.

Let $\text{RunPositive} := \{(\text{run} \mapsto \varepsilon, \text{hide} \mapsto 1 - \varepsilon) \mid \varepsilon > 0\}$ be the set of all strategies that put positive probability on running. We have

$$\text{RunPositive} \prec_{U^k, \mathcal{S}(s_{\text{hide}})} \{\text{hide} \mapsto 1\}.$$

Again, this is true even when using $V_{\mathcal{R}}$ as valuation. Note that we have this weak dominance even though the supremum over the set yields the optimum valuation, namely $\sup_{\rho \in \text{RunPositive}} \inf_{\sigma \in \mathcal{S}(s_{\text{hide}})} \mathfrak{B}(U^k)(s, \rho, \sigma) = 1$. This exemplifies the strictness of our notion of dominance. It is crucial that our notion of domination can distinguish these sets of strategies: The set RunPositive contains all strategies that leave the EC. However, none of them is optimal (even though the supremum over all of them is), which is exactly the reason why VI chooses the staying strategy $\{\text{hide} \mapsto 1\}$ for updating the valuation, and thus is stuck. \triangle

We remark on several technicalities of Def. 12:

- “Weak” dominance: The term “weak” might be misleading. We choose to use the word for consistency with [35, Def. 4.12]. There, *weak* domination concerns Condition (ii), only requiring that there *exists* a counter-strategy where the inequality is strict; strict domination requires Condition (ii) *for all* counter-strategies. One might be tempted to use weak domination to denote what we call “not worse”, i.e. only require that there is a strategy in the first set that has optimal valuation at least as good as all in the other set; or to think it only refers to the numerical comparators, e.g. \geq and $>$ (as is sometimes the case when only comparing single strategies).
- Set-related challenges: The definition is challenging since we cannot speak about actions, but have to consider sets of — possible mixed — strategies. The exact quantification of the strategies is relevant. Further, it depends not only on the two sets we are comparing, but also on the counter-strategies of the opponent. Thus, we provide the definition explicitly for both players, to avoid confusion that could arise from just saying that they are analogous.
- All-quantification instead of optima: The definition uses all-quantification instead of optima. Concretely, weak dominance for Player \mathcal{R} uses $\forall \rho_2 \in \mathcal{R}_2$ instead of writing $\sup_{\rho_2 \in \mathcal{R}_2}$. The latter definition cannot sufficiently distinguish sets of strategies, since the supremum of a set need not be contained in it, as we exemplified in Ex. 13. This fact is extremely important, as in the proof of Thm. 21, we pick the maximum from a set of strategies, and the existence of this maximum is guaranteed only because of the correct definition of dominance.
- Locally optimal strategies are not worse than any other: Formally, this claim is that for all locally optimal strategy

$\rho \in \mathcal{R}(s)$ with respect to v , we have $\mathcal{R}(s) \preceq_{v, \mathcal{S}(s)} \{\rho \mapsto 1\}$ (and dually for Player \mathcal{S}). This is immediate from Def. 12, since a locally optimal strategy maximizes $\inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho, \sigma)$, and thus satisfies Condition (i) when compared to all other strategies. We will use this fact throughout the paper.

- Notation: When the valuation is clear from the context, we omit it for the sake of readability. Further, if we say that a strategy ρ_1 weakly dominates another strategy ρ_2 with respect to a counter-strategy σ , then we mean that $\{\rho_2\} \prec_{\{\sigma\}} \{\rho_1\}$.

We prove a lemma about the relation of weak domination and not being worse that is useful and instructive. The proof in App. C-B works by straightforward unfolding of definitions and rewriting.

Lemma 14 (Negating Weak Domination — Proof in App. C-B). *Let v be a valuation, $s \in S$ a state, $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}' \subseteq \mathcal{R}(s)$ and $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}' \subseteq \mathcal{S}(s)$ sets of local strategies.*

If for some sets of strategies we do not have $\mathcal{R}_2 \prec_{v, \mathcal{S}'} \mathcal{R}_1$, then we have $\mathcal{R}_1 \preceq_{v, \mathcal{S}'} \mathcal{R}_2$. Analogously, not $\mathcal{S}_2 \prec_{v, \mathcal{R}'} \mathcal{S}_1$ implies $\mathcal{S}_1 \preceq_{v, \mathcal{R}'} \mathcal{S}_2$.

To complete our intuitive understanding of the definition of domination, we point out a connection to the standard definition of weak domination: “A rational player does not use a dominated strategy.” [35, Asm. 4.13] If a strategy is not dominated, by Lem. 14 it is not worse than any other strategy. This is exactly what we argued above: Locally optimal strategies are not worse than any other. We often use this fact throughout the paper.

Next, we formally define *leaving and staying strategies*. Given a set of states, a leaving strategy ensures that the set of successor states contains states outside the given set of states for all given counter-strategies of the opponent player. A strategy is staying if all successor states belong to the given set of states for all given counter-strategies of the opponent player. Note that a strategy can be neither leaving nor staying, if the set is exited for some, but not all counter-strategies of the opponent.

Definition 15 (Leaving and Staying Strategies). Consider a set of states $\mathcal{X} \subseteq S$ and a state $s \in \mathcal{X}$. Let $\mathcal{R}' \subseteq \mathcal{R}(s)$ and $\mathcal{S}' \subseteq \mathcal{S}(s)$ be sets of strategies of Player \mathcal{R} and \mathcal{S} , respectively. The set of (local) *leaving strategies* for Player \mathcal{R} with respect to \mathcal{S}' , is given by

$$\mathcal{R}_L(\mathcal{S}', \mathcal{X}, s) := \{\rho \in \mathcal{R}(s) \mid \forall \sigma \in \mathcal{S}'.(s, \rho, \sigma) \text{ leaves } \mathcal{X}\},$$

and for Player \mathcal{S} with respect to \mathcal{R}' by

$$\mathcal{S}_L(\mathcal{R}', \mathcal{X}, s) := \{\sigma \in \mathcal{S}(s) \mid \forall \rho \in \mathcal{R}'.(s, \rho, \sigma) \text{ leaves } \mathcal{X}\}.$$

A strategy that is not leaving is called *non-leaving*. The set of all non-leaving Player \mathcal{R} strategies is denoted by $\mathcal{R}_{\bar{L}}(\mathcal{S}', \mathcal{X}, s)$ (or $\mathcal{S}_{\bar{L}}$ for Player \mathcal{S}).

In contrast, the set of *staying* strategies at a state $s \in \mathcal{X}$ for Player \mathcal{R} with respect to \mathcal{S}' , is given by

$$\mathcal{R}_S(\mathcal{S}', \mathcal{X}, s) := \{\rho \in \mathcal{R}(s) \mid \forall \sigma \in \mathcal{S}'.(s, \rho, \sigma) \text{ staysIn } \mathcal{X}\},$$

and for Player \mathcal{S} with respect to \mathcal{R}' by

$$\mathcal{S}_S(\mathcal{R}', \mathcal{X}, s) := \{\sigma \in \mathcal{S}(s) \mid \forall \rho \in \mathcal{R}'.(s, \rho, \sigma) \text{ staysIn } \mathcal{X}\}.$$

Notation: If we consider leaving (or staying) strategies with respect to all counter-strategies, then we omit the set of counter-strategies, i.e. instead of $\mathcal{R}_L(\mathcal{S}(s), \mathcal{X}, s)$ (or $\mathcal{R}_S(\mathcal{S}(s), \mathcal{X}, s)$) we write $\mathcal{R}_L(\mathcal{X}, s)$ (or $\mathcal{R}_S(\mathcal{X}, s)$). We use the same shorthand notion for leaving (staying) strategies of Player \mathcal{S} .

We often speak about a leaving/staying pair of local strategies, so we provide the following shorthand notations: For a tuple $(s, \rho, \sigma) \in \mathcal{S} \times \mathcal{R} \times \mathcal{S}$, we say that (s, ρ, σ) *leaves* \mathcal{X} if and only if $\text{Post}(s, \rho, \sigma) \cap (\mathcal{S} \setminus \mathcal{X}) \neq \emptyset$. Analogously, we say that (s, ρ, σ) *staysIn* \mathcal{X} if and only if $\text{Post}(s, \rho, \sigma) \cap (\mathcal{S} \setminus \mathcal{X}) = \emptyset$ (or, equivalently, $\text{Post}(s, \rho, \sigma) \subseteq \mathcal{X}$).

Intuition of Hazardous Strategies: Using the definitions of dominance and leaving or staying, we can now classify strategies of Player \mathcal{R} that can lead to non-convergence. Intuitively, a hazardous strategy is one that Player \mathcal{R} chooses, even though it can be staying for some counter-strategies. Thus, such a strategy (i) is non-leaving (i.e. there exist counter-strategies that make it staying), and (ii) it is not worse than any other strategy so that Player \mathcal{R} may choose it. Moreover, to be problematic for convergence, (iii) the strategy weakly dominates all leaving strategies, i.e. leaving strategies are not chosen for the update.

Definition 16 (Hazardous Strategy). Let $\mathcal{X} \subseteq \mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ be as set of states, \mathbf{v} a valuation, and $s \in \mathcal{X}$. A strategy $\rho \in \mathcal{R}(s)$ is called hazardous with respect to \mathbf{v} if it satisfies:

- (i) $\rho \in \mathcal{R}_L(\mathcal{X}, s)$,
- (ii) $\mathcal{R}(s) \setminus \{\rho\} \preceq_{\mathcal{S}(s)} \{\rho\}$, and
- (iii) $\mathcal{R}_L(\mathcal{X}, s) \prec_{\mathcal{S}(s)} \{\rho\}$.

$\text{Hazard}_{\mathbf{v}}(\mathcal{X}, s)$ denotes the set of all hazardous strategies at state s with respect to a set of states \mathcal{X} and a valuation \mathbf{v} .

We mention a corner case: In a state where Player \mathcal{R} possesses no leaving strategies, all optimal strategies are hazardous (note in particular that Condition (iii) is trivially satisfied, since the dominated set of strategies $\mathcal{R}_L(\mathcal{X}, s)$ is empty).

Example 17 (Hazardous strategy). Consider again the matrix game defined in Eq. (4) and the initial valuation $U^0(s_{\text{hide}}) = U^0(s_{\text{home}}) = 1$ and $U^0(s_{\text{wet}}) = 0$. The strategy $\rho' := \{\text{hide} \mapsto 1\}$ is hazardous because: (i) It is non-leaving. (ii) It is an optimal strategy, i.e. it is not worse than any other strategy. (iii) It weakly dominates the set of all leaving strategies, see Ex. 13. \triangle

Definition 18 (Bloated End Component (BEC)). An $\text{EC } \mathcal{X} \subseteq \mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ is called *bloated end component* (BEC) with respect to a valuation \mathbf{v} if for all $s \in \mathcal{X}$ it holds that $\text{Hazard}_{\mathbf{v}}(\mathcal{X}, s) \neq \emptyset$.

Example 19 (Bloated End Component). Consider the **CSG** Hide-Run-or-Slip from Fig. 1. As discussed in Ex. 17, there exists a hazardous strategy in state s_{hide} . Moreover, $\{s_{\text{hide}}\}$ is an **EC**, since under the pair of strategies that plays hide and wait, the play stays in it. Consequently, $\{s_{\text{hide}}\}$ is a **BEC** and therefore VI does not converge in this state, see Ex. 6. \triangle

We provide a lemma that captures the intuition of what it means to (not) be a **BEC**, and that is also useful in several proofs:

Lemma 20 (Negating Bloated — Proof in App. C-B). *If an $\text{EC } \mathcal{X} \subseteq \mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ is not bloated for a valuation \mathbf{v} , then there exists a state $s \in \mathcal{X}$ that has a locally optimal strategy that is leaving, formally $\exists \rho \in \mathcal{R}_L(\mathcal{X}, s). \mathcal{R}(s) \preceq_{\mathbf{v}, \mathcal{S}(s)} \{\rho\}$.*

C. Convergence in the Absence of BECs

Now we can prove that **BECs** indeed are the reason that VI does not converge for over-approximations.

Theorem 21 (Non-convergence implies **BECs** — Proof in App. C-B). *Let $U^* := \lim_{k \rightarrow \infty} U^k$ be the limit of the naive upper bound iteration (Eq. (3)) on the **CSG** G . If VI from above does not converge to the value in the limit, i.e. $U^* > \mathbf{V}_{\mathcal{R}}$, then the **CSG** G contains a **BEC** in $\mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ with respect to U^* .*

Proof sketch. This proof builds on the proof of Thm. 9. There, we constructed a set \mathcal{X} maximizing the difference between U^* and $\mathbf{V}_{\mathcal{R}}$ and showed that if there is a pair of optimal strategies leaving \mathcal{X} , then we can derive a contradiction: The upper bound decreases, which contradicts the fact that it is a fixpoint. In the context of the other proof, that allowed us to show that without **ECs**, VI converges, because without **ECs** it is impossible to have a set of states where all optimal strategies stay in that set.

In the presence of **ECs**, states can indeed have a positive difference between U^* and $\mathbf{V}_{\mathcal{R}}$, see e.g. Ex. 6. Our goal is to prove that at least one of these **ECs** is bloated. Thus, we assume for contradiction that no **EC** is bloated under U^* . Then, by Lem. 20, there is an optimal leaving strategy for Player \mathcal{R} . Using that, we can repeat the argument from Thm. 9, showing that in this case U^* would decrease. Again, this is a contradiction because it is a fixpoint of applying Bellman updates (Lem. 48). Thus, the initial assumption that no **EC** is bloated is false, and we can conclude that there exists a **BEC**. \square

Remark 22 (Relation to [18]). Def. 18 of **BEC** is more general than the definition of **BEC** for **TSGs** in [18, Definition 4]. The differences are that in [18], an **EC** is only called bloated if it is bloated with respect to $\mathbf{V}_{\mathcal{R}}$, whereas we extended that definition to speak about a concrete valuation, similar to [29, Def. 3]. Further, the definition for **TSGs** speaks about the best exit value, which is the optimum among the available actions; in contrast, in **CSGs** the definition of best exit is technically involved and dependent on hazardous strategies, see Def. 27. Thus, our definition of **BEC** does not analyze the exit values, but instead uses a fundamental analysis of the strategies.

Key Contribution: The key novelty of this section is the correct definition of **BEC** that captures when VI from above does not converge. We highlight that this definition contains many technical intricacies: Lifting the notion of an exiting action [18, Section 2.2] from **TSGs** to **CSGs** requires considering sets of local strategies that leave against all opponent-strategies (Def. 15), and considering the additional complication that strategies can be neither leaving nor staying. Further, the exact definition of dominance is very important, as the supremum over all leaving strategies can be a staying strategy, see Ex. 13 and the related discussion in the item “All-quantification instead of optima” after the example.

V. RESOLVING BLOATED END COMPONENTS

A. Solution in TSGs

We have identified **BECs** as the cause of non-convergence. Our method for fixing the over-approximation is again based on the ideas for **TSGs**. We explain the intuition of their solution: Firstly, staying actions yield the valuation that is bloated; thus, we need an additional update of the over-approximation that depends only on leaving actions. The valuation to which we reduce the over-approximation is the *best exit* from the **EC**, which in **TSGs** simply is the leaving action attaining the highest value over all states of Player \mathcal{R} [18, Definition 3]. Secondly, not all states in an **EC** need to have the same value, since Player \mathcal{S} can prevent Player \mathcal{R} from reaching the state that attains the best exit valuation. Hence, an **EC** is decomposed into parts that share the same value, called *simple ECs* (SECs) in [18, Definition 5]. Repeatedly finding these SECs and *deflating* their valuation by setting it to the best exit from the SEC suffices for convergence in **TSGs**.

When generalizing these ideas to **CSGs**, we encounter the following problems: Firstly, the definition of best exit is more involved, since staying and leaving depends not only on actions, but on strategies. Additionally, the supremum over all leaving strategies can be a non-leaving strategy, see Ex. 13. (This is also the reason why globally optimal strategies need not exist in **CSGs**, see Ex. 2). This was the fundamental reason why the solution proposed in [19] did not work, as it was based on actions. Secondly, we need to decompose the **EC** into parts. For this, we use a recursive approach, removing states that have been successfully deflated and checking whether there are further problematic states in the remainder of the **EC**.

Outline: In Sec. V-B, we develop a strategy-based definition of best exit (Def. 28), which relies on identifying the *trapping strategies* (Def. 24) that Player \mathcal{S} uses to keep the play inside the **BEC**; and the *deflating strategies* (Def. 25), the best response of Player \mathcal{R} , namely the leaving strategies that should be played with arbitrarily small probability ε . In Sec. V-C, we provide the **FIND_MBEC** algorithm that finds all maximal **BECs** that are present in a given **MEC**. A maximal **BEC** is a set of states $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ that is a **BEC** and there exists no another set of states $\mathcal{X}' \subseteq S \setminus (T \cup W_{\mathcal{S}})$ that is a **BEC** and $\mathcal{X} \subsetneq \mathcal{X}'$. Finally, Sec. V-D provides the full deflating procedure for **CSGs** in Alg. 3. In particular, it uses

a recursive call to decompose a given **MEC** and deflate all relevant parts of it.

B. Defining the Best Exit

The key problem of a **BEC** is that in all states, all leaving strategies of Player \mathcal{R} are dominated by hazardous strategies. Player \mathcal{S} can play a trapping strategy and thereby make the Bellman update self-dependent. If the current valuation is too high, this prevents convergence. Intuitively, the “pressure” inside the **BEC** is too high, and we want to “deflate” it, by adjusting it to the pressure, i.e. valuation, outside of the **BEC**. Since non-trivial **BECs** neither contain target states nor belong to the winning region of Player \mathcal{S} , there has to exist a state in a **BEC** where the supremum over leaving and staying strategies is equal. To “equalize the pressure” between the **BEC** and the rest of the states, we need to estimate the pressure outside the **BEC**. To do so, we estimate the valuation attainable upon leaving the **BEC** at every state of the **BEC**, called exit value of the state. The best exit value is the maximum of all exit values. Reducing the upper bounds of the states inside the **BEC** to the best exit value, in case the best exit is smaller than the current valuation, “decreases the pressure”. However, since the valuations of the exiting strategies can depend on the valuations of the states that belong to the **BEC**, this procedure may still only converge in the limit, already in **TSGs**. We provide an illustrative example:

Example 23 (Deflating BECs). Consider again the **CSG** **Hide-Run-or-Slip** (Fig. 1). Under the initial valuation (see Eq. (3)) the matrix game played at s_{hide} is given by

$$Z_{U^0}(s_{\text{hide}}) = \begin{pmatrix} \text{throw} & \text{wait} \\ 0 & \frac{2}{3} \\ 1 & 1 \end{pmatrix} \begin{matrix} \text{run} \\ \text{hide} \end{matrix}$$

Due to the hazardous strategy $\{\text{hide} \mapsto 1\}$, the Bellman update cannot improve the initial upper bound of s_{hide} , but remains at 1. However, by the same argument as in Ex. 2, Player \mathcal{R} can use the strategy of running with a probability $\varepsilon > 0$ that is arbitrarily small. This yields a value arbitrarily close to $\frac{2}{3}$. Consequently, we can deflate, i.e. decrease the upper bound of s_{hide} to $\frac{2}{3}$ which is the valuation attainable upon leaving the **BEC** at s_{hide} . After the Bellman update, the matrix game played at s_{hide} is given by

$$Z_{U^1}(s_{\text{hide}}) = \begin{pmatrix} \text{throw} & \text{wait} \\ 0 & \frac{5}{9} \\ 1 & \frac{2}{3} \end{pmatrix} \begin{matrix} \text{run} \\ \text{hide} \end{matrix}$$

The strategy $\{\text{hide} \mapsto 1\}$ is still hazardous, but we can deflate the upper bound of s_{hide} to $\frac{5}{9}$. By repeating these steps, the upper bound converges to $\frac{1}{2}$ in the limit, and only in the limit, similar to the lower bound in Table I. \triangle

How can we find the best exit value in general?: In Ex. 23, we used an argument about playing a leaving action with vanishingly small probability in order to figure out which entry in the matrix we choose for deflating. We provide an

alternative intuition: Player \mathcal{R} plays a hazardous strategy most of the time. The best response of Player \mathcal{S} is to play a weakly dominant strategy that stays in the **EC**, trapping the play. Thus, we can ignore the other strategies of Player \mathcal{S} and consider only the columns of the matrix that correspond to *trapping strategies*. Now, we need to select a leaving strategy of Player \mathcal{R} which then is played with vanishingly low probability. Thus, we restrict the matrix further to use only rows corresponding to leaving strategies. In the example, we end up with the top right entry. We now formalize how we can construct this sub-matrix, called the exiting sub-game.

Definition 24 (Trapping Strategy). Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ be a set of states, v a valuation, and $s \in \mathcal{X}$. A strategy $\sigma \in \mathcal{S}(s)$ is called *trapping strategy* if two conditions are satisfied:

- (i) $\sigma \in \arg \min_{\sigma' \in \mathcal{S}(s)} \max_{\rho \in \mathcal{R}(s)} \mathfrak{B}(v)(s, \rho, \sigma')$, and
- (ii) $\forall \rho \in \text{Hazard}_v(\mathcal{X}, s) : (s, \rho, \sigma) \text{ stays in } \mathcal{X}$.

$\text{Trap}_v(\mathcal{X}, s)$ denotes the set of all trapping strategies at state s with respect to a set of states \mathcal{X} and a valuation v .

Definition 25 (Deflating Strategies). Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ be a set of states. A Player \mathcal{R} strategy $\rho \in \mathcal{R}(s)$ is called *deflating* if two conditions are satisfied:

- (i) $\exists \sigma \in \text{Trap}_v(\mathcal{X}, s)$ such that (s, ρ, σ) leaves \mathcal{X} , and
- (ii) $\text{Supp}(\rho) \cap \bigcup_{\rho' \in \text{Hazard}_v(\mathcal{X}, s)} \text{Supp}(\rho') = \emptyset$.

$\text{Defl}_v(\mathcal{X}, s)$ denotes the set of all deflating strategies at state s with respect to a set of states \mathcal{X} and a valuation v .

Definition 26 (Exiting sub-game). Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ be a set of states, $s \in \mathcal{X}$ a state, and v a valuation. Further, let $Z_v(s)$ be the matrix game played at state $s \in \mathcal{X}$. If $\text{Trap}_v(\mathcal{X}, s) \neq \emptyset$ then, the *exiting sub-game* played at state s , denoted by $Z_v^{\text{exit}}(s)$, is the matrix game where Player \mathcal{R} has the actions in $\bigcup_{\rho \in \text{Defl}_v(\mathcal{X}, s)} \text{Supp}(\rho)$ and Player \mathcal{S} has the actions in $\bigcup_{\sigma \in \text{Trap}_v(\mathcal{X}, s)} \text{Supp}(\sigma)$. The value of the exiting sub-game is given by $V(Z_v^{\text{exit}}(s)) := \max(0, \sup_{\rho \in \text{Defl}_v(\mathcal{X}, s)} \inf_{\sigma \in \text{Trap}_v(\mathcal{X}, s)} \mathfrak{B}(v)(s, \rho, \sigma))$.

We explain how this exiting sub-game and its value are well-defined: The set of deflating strategies can be empty, namely in a state which has no leaving strategies. In this case, the value of the exiting sub-game is the supremum over an empty set, i.e. the smallest possible value, commonly minus infinity and 0 in our case. We highlight this possibility by explicitly taking the maximum of 0 and the value of the exiting subgame when we compute it. If the set of deflating strategies is non-empty, then the set of trapping strategies necessarily is non-empty, too, since the Item (i) of Def. 25 requires existence of a trapping strategy.

Definition 27 (Exit value). Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ be a set of states, $s \in \mathcal{X}$ a state, and v an over-approximation. Then,

the *exit value* from \mathcal{X} attainable at state s is given by

$$\text{exitVal}_v(\mathcal{X}, s) := \begin{cases} \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(v)(s, \rho, \sigma), & \text{if } \text{Hazard}_v(\mathcal{X}, s) = \emptyset \\ \vee \text{Trap}_v(\mathcal{X}, s) = \emptyset; \\ V(Z_v^{\text{exit}}(s)), & \text{else.} \end{cases}$$

For **BECs** that consists of more than one state, the exit values attainable at different states within the **BEC** may differ. Consequently, to ensure that deflating does not reduce the value of any of the states in the **BEC** below its actual value we estimate the exit values at each state of the **BEC**, and finally select the maximal exit value, i.e. the best exit value, for deflating.

Definition 28 (Best Exit). Let $\mathcal{X} \subseteq (S \setminus (T \cup W_{\mathcal{S}}))$ be a set of states and v a valuation. The *best exit value* with respect to a set \mathcal{X} and a valuation v is given by

$$\text{bestExitVal}_v(\mathcal{X}) := \max_{s \in \mathcal{X}} \text{exitVal}_v(\mathcal{X}, s).$$

The *best exit*, denoted by $\text{bestExit}_v(\mathcal{X})$, is a state obtaining $\text{bestExitVal}_v(\mathcal{X})$, and the set of all best exits is denoted by $\text{bestExits}_v(\mathcal{X})$.

Lemma 29 (bestExitVal is sound). Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{S}})$ be an **EC**, and $U \in [0, 1]^{|S|}$ be a valuation with $U \geq V_{\mathcal{R}}$. Then, for all states $s \in \mathcal{X}$, we have $\text{bestExitVal}_U(\mathcal{X}) \geq V_{\mathcal{R}}(s)$.

Proof. This proof relies on the technical Lem. 49 stating that

- (i) $\mathcal{X}' := \{s \in \mathcal{X} \mid V_{\mathcal{R}}(s) \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)\} \neq \emptyset$, and
- (ii) $\max_{s \in \mathcal{X}'} V_{\mathcal{R}}(s) \geq \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$.

Let $\mathcal{X}' \subseteq \mathcal{X}$ be a set of states satisfying Condition (i). Further, choose $e \in \arg \max_{t \in \mathcal{X}'} V_{\mathcal{R}}(t)$ as one of the exits from \mathcal{X}' . By Item (ii) of Lem. 49, we have that for all $s \in \mathcal{X}$: $V_{\mathcal{R}}(s) \leq V_{\mathcal{R}}(e)$. It remains to show $\text{bestExitVal}_U(\mathcal{X}) \geq V_{\mathcal{R}}(e)$. We conclude as follows:

$$\begin{aligned} V_{\mathcal{R}}(e) &\leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, e) && (\text{Since } e \in \mathcal{X}') \\ &= \text{bestExitVal}_{V_{\mathcal{R}}}(\mathcal{X}) && (\text{By the choice of } e) \\ &\leq \text{exitVal}_U(\mathcal{X}, e) && (\text{By Lem. 50}) \\ &\leq \text{bestExitVal}_U(\mathcal{X}). && (\text{By Def. 28}) \end{aligned}$$

We want to highlight that the statement of Lem. 50 is non-trivial and the proof is technically involved. \square

C. Finding Maximal BECs

Since a **BEC** might contain other **BECs**, we want to find *maximal BECs*. A **BEC** \mathcal{X} is maximal if there exists no **BEC** \mathcal{X}' such that $\mathcal{X} \subsetneq \mathcal{X}'$. The existence of maximal **BECs** is proven in App. C-C.

Given a **CSG** G , a **MEC** C and the current upper bound estimate U , Alg. 2 finds all maximal **BECs** within C as follows: The set B contains all states for which hazardous strategies exist with respect to the set C . In case B is non-empty, it might consist of multiple disjoint **MECs**. Therefore, the algorithm calls the `FIND_MECS` procedure on B and returns all **MECs**

Algorithm 2 Algorithm for finding maximal BECs.

```

1: Algorithm FIND_MBECs(CSG G, MEC C, upper bound
   estimate U)
2:    $B := \{s \in C \mid \text{Hazard}_U(C, s) \neq \emptyset\}$ 
3:   if  $B \neq \emptyset$  then ▷ C contains BECs
4:     return FIND_MECS(G, B)
5:   else
6:     return  $\emptyset$  ▷ There exists no BEC in C

```

Algorithm 3 Algorithm for deflating BECs.

```

1: Algorithm DEFLATE(CSG G, upper bound estimate U,
   MEC C)
2:   for  $\mathcal{X} \in \text{FIND\_MBECs}(G, C, U)$  do
3:      $u \leftarrow \text{bestExitVal}_U(\mathcal{X})$ 
4:     for  $s \in \mathcal{X}$  do
5:        $U(s) \leftarrow \min(U(s), u)$ 
6:     for  $\mathcal{E} \in \text{FIND\_MECS}(\mathcal{X} \setminus \text{bestExits}_U(\mathcal{X}))$  do
7:        $U \leftarrow \text{DEFLATE}(G, U, \mathcal{E})$ 
▷ Recursively deflate sub-BECs
8:   return U

```

that are in B . Otherwise, B is empty, which means that C does not contain any BECs, and the empty set is returned.

Lemma 30 (FIND_MBECs is correct— Proof in App. C-C). *For a CSG, a MEC C and a valid upper bound U , it holds that $\mathcal{X} \in \text{FIND_MBECs}(G, C, U)$ if and only if \mathcal{X} is a BEC in C and there exists no $T \subseteq C$ that is a BEC and $\mathcal{X} \subsetneq T$.*

Proof sketch. To prove the “ \Rightarrow ” direction, we assume towards a contradiction that \mathcal{X} is a maximal EC but not a BEC. Then, all Player \mathcal{R} strategies must violate at least one of the conditions posed by Def. 16. We consider each condition separately and derive a contradiction to the assumption that $\mathcal{X} \in \text{FIND_MBECs}(G, C, U)$.

To prove the opposite direction, i.e., “ \Leftarrow ”, we assume towards a contradiction that \mathcal{X} is a maximal BEC but $\mathcal{X} \notin \text{FIND_MBECs}(G, C, U)$. Here, we again make a case distinction. In the first case, some state $s \in \mathcal{X}$ was removed by the algorithm because at least one player has an optimal strategy that can leave \mathcal{X} . However, this is a contradiction to the assumption that \mathcal{X} is a BEC. In the second case, we assume towards a contradiction that $\mathcal{X} \notin \text{FIND_MECS}(G, B)$ (where B is defined as in FIND_MBECs). This is then a contradiction to the assumption that \mathcal{X} is a maximal EC. \square

D. Convergent Bounded Value Iteration for CSGs and the Recursive Structure of ECs

Finally, Alg. 3 describes our main goal: the deflating procedure for BECs, to be plugged into Alg. 1. The algorithm takes a CSG G , the current upper bound estimate U , and a MEC C as input. First, the algorithm searches for all maximal BECs that might be contained in the MEC. The current upper bound estimate is returned if no BEC exists in the MEC. Otherwise, at least one BEC exists and must be deflated. If multiple BECs

are found, each maximal BEC is deflated separately (see the for-loop in Line 2).

To deflate a maximal BEC \mathcal{X} , first, the *best* exit value $\text{bestExitVal}_U(\mathcal{X})$ is estimated. Next, each state of the BEC is considered and if the best exit value is smaller than the current upper bound estimate at that state, then it is reduced to the best exit value (as nothing better can be reached).

However, notice that within a BEC \mathcal{X} there might exist another BEC \mathcal{X}' . From \mathcal{X}' , Player \mathcal{R} might not be able to get to the best exit of \mathcal{X} (the globally best one in \mathcal{X}) but only to a worse one, locally optimal for \mathcal{X}' . Hence, for such a sub-BEC \mathcal{X}' , more aggressive deflating to $\text{bestExitVal}_U(\mathcal{X}')$ is due. In other words, after deflating to $\text{bestExitVal}_U(\mathcal{X})$ we have handled all states where \mathcal{R} can ensure reaching this best exit, and we can ignore these states for the moment; we can also ignore this best exit and have a fresh look at which states are *now* in a BEC and can reach the *second best* exit, i.e., the best exit in this remainder. Subsequently, we continue with the third best option etc.

Consequently, on Lines 6-7, DEFLATE is called *recursively* on all MECs that are contained in $\mathcal{X} \setminus \text{bestExits}_U(\mathcal{X})$, i.e. after removing all best exits. The procedure is repeated independently for each maximal BEC contained in C . A full example showing how DEFLATE works on a more complex BEC is included in App. B-B.

Remark 31 (Structure of ECs). This elucidates the hierarchical structure of ECs and their corresponding best exits. The recursive call of DEFLATE exposes the partial order over the ECs, their sub-ECs, and “internally” transient states (those not within sub-ECs after removing the best exit since they are bound to the just removed exit or another sub-EC that is a BEC with a lower value). This hierarchy captures (i) the ordering of exits by their values and (ii) the “independence” of exits with possibly different values when visiting one from another cannot be ensured.

Our goal for the remainder of this work is to show that Alg. 1 with deflation is correct and converges, i.e. that complementing Bellman updates \mathfrak{B} with deflation results in a sequence of upper bounds U that converges to the value from above. For the sake of simplicity, we denote by $\mathfrak{D} : [0, 1]^{|S|} \rightarrow [0, 1]^{|S|}$ the operator that performs DEFLATE on a given valuation for all MECs in the CSG (which is reasonable since Alg. 1 performs DEFLATE on all MECs (in an arbitrary but fixed ordering)).

Remark 32 (Valid upper bounds). In the following, whenever we quantify an upper bound (a.k.a. over-approximation), we require it to be *valid*; meaning that it was obtained by iteratively applying deflating and the Bellman update on the initial over-approximation U^0 from Eq. (3), i.e. $U = (\mathfrak{D} \circ \mathfrak{B})^k(U^0)$ for some $k \in \mathbb{N}$. The reason for this restriction is that for convergence, we require \mathfrak{D} to be order-preserving. While \mathfrak{D} is order-preserving on valid upper bounds (see Lem. 36), in general it is not monotonic for arbitrary $U \in \mathbb{R}^{|S|}$. We illustrate this in App. B-C.

Definition 33 (Valid upper bound). An upper bound $U \in [0, 1]^{|S|}$ is called *valid* if there exists $k \in \mathbb{N}$ such that $U = (\mathfrak{D} \circ \mathfrak{B})^k(U^0)$, or if $U = V_{\mathcal{R}}$.

We proceed as follows: After proving fundamental properties of both operators \mathfrak{B} and \mathfrak{D} in Lem. 34, we use these properties to show that valid upper bounds are indeed upper bounds, i.e. they are always greater or equal than the value in Lem. 35. With correctness established, Lem. 36 shows that on valid upper bounds, $(\mathfrak{D} \circ \mathfrak{B})$ is order-preserving, which is a necessary ingredient for convergence. Finally, Thm. 37 concludes by proving soundness and completeness of the full algorithm.

Lemma 34 (Properties of \mathfrak{D} and \mathfrak{B}). *Let $v \in [0, 1]^{|S|}$ be a valuation. If $v \geq V_{\mathcal{R}}$, then:*

- (i) $\mathfrak{B}(v) \leq v$ and $\mathfrak{D}(v) \leq v$.
- (ii) $\mathfrak{B}(v) \geq V_{\mathcal{R}}$ and $\mathfrak{D}(v) \geq V_{\mathcal{R}}$.

Proof. For the Bellman operator, both properties follow from the fact that the value $V_{\mathcal{R}}$ is the least fixpoint of the Bellman operator, see [17, Thm. 1]. Thus, given a valuation greater than the value, it cannot increase, and it cannot become smaller than the value.

For \mathfrak{D} , observe that deflation only updates the valuation in Line 5 of Alg. 3 when setting the valuation of a state to $\min(U(s), \text{bestExitVal}_U(\mathcal{X}))$ for some EC \mathcal{X} . Item (i) holds because of taking the minimum with the current valuation, so it can only decrease. For Item (ii), we show in Lem. 29 that for every EC \mathcal{X} and state $s \in \mathcal{X}$, we have that $\text{bestExitVal}_U(\mathcal{X}) \geq V_{\mathcal{R}}(s)$. This proves our goal, as the only update of DEFLATE keeps the valuation greater than $V_{\mathcal{R}}$ in every state. The proof of Lem. 29 is technically involved, having to unfold many definitions in order to show the following intuitive fact: No state can have a larger value than that of some exit from the EC it is contained in. \square

Lemma 35 (Soundness of valid upper bounds). *For all $k \in \mathbb{N}$ it holds that $(\mathfrak{D} \circ \mathfrak{B})^k(U^0) \geq V_{\mathcal{R}}$.*

Proof. We proceed by induction over k .

Base case: $k = 0$, thus $(\mathfrak{D} \circ \mathfrak{B})^0(U^0) = U^0 \geq V_{\mathcal{R}}$.

Induction hypothesis: For all $k \geq 0$, we assume that

$$(\mathfrak{D} \circ \mathfrak{B})^k(U^0) \geq V_{\mathcal{R}}.$$

Induction step: To show: $(\mathfrak{D} \circ \mathfrak{B})^{k+1}(U^0) \geq V_{\mathcal{R}}$. We know by induction hypothesis that $(\mathfrak{D} \circ \mathfrak{B})^k(U^0) \geq V_{\mathcal{R}}$. Applying Item (ii) of Lem. 34 for \mathfrak{B} , we obtain $\mathfrak{B}((\mathfrak{D} \circ \mathfrak{B})^k(U^0)) \geq V_{\mathcal{R}}$. From this and Item (ii) for \mathfrak{D} , we get $\mathfrak{D}(\mathfrak{B}((\mathfrak{D} \circ \mathfrak{B})^k(U^0))) \geq V_{\mathcal{R}}$, proving our goal. \square

Lemma 36 ($(\mathfrak{D} \circ \mathfrak{B})$ is order-preserving on valid upper bounds). *Let v_1, v_2 be valid upper bounds with $v_1 \geq v_2$. It holds that $(\mathfrak{D} \circ \mathfrak{B})(v_1) \geq (\mathfrak{D} \circ \mathfrak{B})(v_2)$.*

Proof. We know by Lem. 35 that all valid upper bounds are greater or equal to the value (including the value itself). Thus, if $v_2 = V_{\mathcal{R}}$, the statement holds, and if $v_1 = V_{\mathcal{R}}$, then

$v_2 = V_{\mathcal{R}}$ as well, since $v_1 \geq v_2$. It remains to consider the case that both valuations come from repeated application of the deflation and Bellman operators, i.e. $v_1 = (\mathfrak{D} \circ \mathfrak{B})^i(U^0)$ and $v_2 = (\mathfrak{D} \circ \mathfrak{B})^j(U^0)$. We assume $v_1 \neq v_2$, since otherwise the claim trivially holds.

By Item (i) of Lem. 34, every application of the operators can only decrease the resulting value; the item remains applicable, since the resulting valuations are always greater than or equal to the value. Consequently, $i < j$, as $v_1 > v_2$. Using this and applying Item (i) of Lem. 34 again, we conclude by stating

$$\begin{aligned} (\mathfrak{D} \circ \mathfrak{B})(v_1) &= (\mathfrak{D} \circ \mathfrak{B})^{i+1}(U^0) \geq \\ &(\mathfrak{D} \circ \mathfrak{B})^{j+1}(U^0) = (\mathfrak{D} \circ \mathfrak{B})(v_2). \end{aligned}$$

\square

Theorem 37 (Soundness and completeness - Proof in App. C-D). *For CSGs Alg. 1, using Alg. 3 as DEFLATE, produces monotonic sequences L under- and U over-approximating $V_{\mathcal{R}}$, and terminates for every $\varepsilon > 0$.*

Proof sketch. Soundness and convergence of lower bounds is classical [17, Thm. 1], and our algorithm does not modify the computation of under-approximations. The soundness of the upper bounds is immediate from Lem. 35, since all upper bounds computed by the algorithm are valid, and thus greater or equal than the value. Proving the convergence of the upper bounds is the main challenge. First, in Lem. 52 we provide the technical proof that the upper bound in Alg. 1 indeed converges to a fixpoint, using that the operators are order-preserving (Lem. 36) and arguments from lattice theory. Then, we use the same idea we have utilized in the proofs of Thms. 9 and 21: We assume for contradiction that there exists a state where the difference between the fixpoint of the upper bound in Alg. 1 and the true value is strictly greater than zero. The states with the largest such difference contain a BEC (by Thm. 21), that eventually will be found and deflated; since deflation depends on the outside of the BEC, this decreases the upper bound. This causes a contradiction to the fact that the upper bound has converged to a fixpoint. Consequently, there can be no state with a positive difference, and the upper bounds converge, too. \square

VI. CONCLUSION AND FUTURE WORK

We have introduced a convergent over-approximation for concurrent stochastic games with reachability and safety objectives, thus giving value iteration the first sound stopping criterion and turning it into an anytime algorithm. Since the games are concurrent and (ε) -optimal strategies may need to be randomized, we could not use the technique of simple end components of [18]. Instead, we identify bloated end components where the play can get stuck forever and recursively deflate the over-approximations of these states to the best possible value attainable upon leaving the end component. We leave an efficient implementation for future work as an extension of the standard model checker PRISM-GAMES [32].

REFERENCES

- [1] Pranav Ashok, Krishnendu Chatterjee, Przemysław Daga, Jan Křetínský, and Tobias Meggendorfer. Value Iteration for Long-Run Average Reward in Markov Decision Processes. In Rupak Majumdar and Viktor Kunčák, editors, *Computer Aided Verification*. Springer International Publishing, 2017.
- [2] Christel Baier and Joost-Pieter Katoen. *Principles of Model Checking*. MIT Press, April 2008.
- [3] Christel Baier, Joachim Klein, Linda Leuschner, David Parker, and Sascha Wunderlich. Ensuring the Reliability of Your Model Checker: Interval Iteration for Markov Decision Processes. In Rupak Majumdar and Viktor Kunčák, editors, *Computer Aided Verification*. Springer International Publishing, 2017.
- [4] Benjamin Bordaïs, Patricia Bouyer, and Stéphane Le Roux. Subgame Optimal Strategies in Finite Concurrent Games with Prefix-Independent Objectives. In Orna Kupferman and Paweł Sobocinski, editors, *Foundations of Software Science and Computation Structures*, pages 541–560. Springer Nature Switzerland, 2023.
- [5] Tomáš Brázdil, Krishnendu Chatterjee, Martin Chmelík, Vojtěch Forejt, Jan Křetínský, Marta Kwiatkowska, David Parker, and Mateusz Ujma. Verification of Markov decision processes using learning algorithms. In *International Symposium on Automated Technology for Verification and Analysis*, pages 98–114. Springer, 2014.
- [6] Tomáš Brázdil, Krishnendu Chatterjee, Martin Chmelík, Vojtěch Forejt, Jan Křetínský, Marta Kwiatkowska, David Parker, and Mateusz Ujma. Verification of Markov Decision Processes Using Learning Algorithms. In Franck Cassez and Jean-François Raskin, editors, *Automated Technology for Verification and Analysis*. Springer International Publishing, 2014.
- [7] Krishnendu Chatterjee, Luca de Alfaro, and Thomas A. Henzinger. Termination criteria for solving concurrent safety and reachability games. In *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms*, pages 197–206. SIAM, 2009.
- [8] Krishnendu Chatterjee, Luca de Alfaro, and Thomas A. Henzinger. Qualitative concurrent parity games. *ACM Transactions on Computational Logic*, July 2011.
- [9] Krishnendu Chatterjee, Luca de Alfaro, and Thomas A. Henzinger. Strategy improvement for concurrent reachability and safety games. *arXiv preprint arXiv:1201.2834*, 2012.
- [10] Krishnendu Chatterjee, Luca de Alfaro, and Thomas A. Henzinger. Strategy improvement for concurrent reachability and turn-based stochastic safety games. *Journal of computer and system sciences*, 79(5):640–657, 2013.
- [11] Krishnendu Chatterjee and Thomas A. Henzinger. Value Iteration. In Orna Grumberg and Helmut Veith, editors, *25 Years of Model Checking: History, Achievements, Perspectives*. Springer, 2008.
- [12] Costas Courcoubetis and Mihalis Yannakakis. The complexity of probabilistic verification. *Journal of the ACM*, 42, July 1995.
- [13] B. A. Davey and H. A. Priestley. *Introduction to Lattices and Order*. Cambridge University Press, Cambridge, 2 edition, 2002.
- [14] Luca De Alfaro. How to specify and verify the long-run average behaviour of probabilistic systems. In *Proceedings. Thirteenth Annual IEEE Symposium on Logic in Computer Science (Cat. No. 98CB36226)*, pages 454–465. IEEE, 1998.
- [15] Luca de Alfaro and Thomas A. Henzinger. Concurrent omega-regular games. In *Logic in Computer Science, 2000. Proceedings. 15th Annual IEEE Symposium on*, pages 141–154. IEEE, 2000.
- [16] Luca De Alfaro, Thomas A. Henzinger, and Orna Kupferman. Concurrent reachability games. *Theoretical Computer Science*, 2007.
- [17] Luca de Alfaro and Rupak Majumdar. Quantitative solution of omega-regular games. *Journal of Computer and System Sciences*, 68(2):374 – 397, 2004. Special Issue on STOC 2001.
- [18] Julia Eisentraut, Edon Kelmendi, Jan Křetínský, and Maximilian Weininger. Value iteration for simple stochastic games: Stopping criterion and learning algorithm. *Information and Computation*, 2022.
- [19] Julia Eisentraut, Jan Křetínský, and Alexej Rotar. Stopping Criteria for Value and Strategy Iteration on Concurrent Stochastic Reachability Games, 2019.
- [20] Kousha Etessami and Mihalis Yannakakis. Recursive Concurrent Stochastic Games. *Logical Methods in Computer Science*, Volume 4, Issue 4, November 2008.
- [21] H. Everett. *RECURSIVE GAMES*, pages 47–78. Princeton University Press, 1957.
- [22] Søren Kristoffer Stil Frederiksen and Peter Bro Miltersen. Approximating the Value of a Concurrent Reachability Game in the Polynomial Time Hierarchy. Springer, 2013.
- [23] Serge Haddad and Benjamin Monmege. Interval iteration algorithm for mdps and imdps. *Theoretical Computer Science*, 735:111–131, 2018.
- [24] Serge Haddad and Benjamin Monmege. Interval Iteration Algorithm for MDPs and IMDPs. *Theoretical Computer Science*, 2018.
- [25] Kristoffer Arnsfelt Hansen, Rasmus Ibsen-Jensen, and Peter Bro Miltersen. The Complexity of Solving Reachability Games Using Value and Strategy Iteration. *Theory of Computing Systems*, 55, 2014.
- [26] Kristoffer Arnsfelt Hansen, Michal Koucký, Niels Lauritzen, Peter Bro Miltersen, and Elias P. Tsigaridas. Exact algorithms for solving stochastic games: Extended abstract. In *Proceedings of the Forty-Third Annual ACM Symposium on Theory of Computing*, 2011.
- [27] Arnd Hartmanns, Sebastian Junges, Tim Quatmann, and Maximilian Weininger. A Practitioner’s Guide to MDP Model Checking Algorithms. In Sriram Sankaranarayanan and Natasha Sharygina, editors, *Tools and Algorithms for the Construction and Analysis of Systems*. Springer Nature Switzerland, 2023.
- [28] Frederick S. Hillier and Gerald J. Lieberman. *Introduction to Operations Research*. McGraw-Hill Higher Education, 2010.
- [29] Jan Křetínský, Tobias Meggendorfer, and Maximilian Weininger. Stopping Criteria for Value Iteration on Stochastic Games with Quantitative Objectives. In *2023 38th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, 2023.
- [30] Jan Křetínský, Emanuel Ramneantu, Alexander Slivinskiy, and Maximilian Weininger. Comparison of algorithms for simple stochastic games. *Information and Computation*, 2022.
- [31] P. R. Kumar and T. H. Shiao. Existence of Value and Randomized Strategies in Zero-Sum Discrete-Time Stochastic Dynamic Games. *SIAM Journal on Control and Optimization*, 19, 1981.
- [32] Marta Kwiatkowska, Gethin Norman, David Parker, and Gabriel Santos. Automated verification of concurrent stochastic games. In *International Conference on Quantitative Evaluation of Systems*, pages 223–239. Springer, 2018.
- [33] Marta Kwiatkowska, Gethin Norman, David Parker, and Gabriel Santos. PRISM-games 3.0: Stochastic Game Verification with Concurrency, Equilibria and Time. In Shuvendu K. Lahiri and Chao Wang, editors, *Computer Aided Verification*. Springer International Publishing, 2020.
- [34] Donald A. Martin. The determinacy of blackwell games. 1998.
- [35] Michael Maschler, Shmuel Zamir, and Eilon Solan. *Game Theory*. Cambridge University Press, June 2020.
- [36] John Nash. Non-Cooperative Games. 1951.
- [37] John F. Nash. Equilibrium points in n-person games. 1950.
- [38] Miquel Oliu-Barton. New algorithms for solving zero-sum stochastic games. *Mathematics of Operations Research*, 46(1):255–267, 2021.
- [39] T. Parthasarathy. Discounted, positive, and noncooperative stochastic games. *International Journal of Game Theory*, 2(1):25–37, Dec 1973.
- [40] Grant Olney Passmore and Paul B. Jackson. Combined decision techniques for the existential theory of the reals. In *Calculemus/MKM*, volume 5625 of *Lecture Notes in Computer Science*, pages 122–137. Springer, 2009.
- [41] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, 2009.
- [42] T. E. S. Raghavan and J. A. Filar. Algorithms for stochastic games — A survey. *Zeitschrift für Operations Research*, 35(6):437–472, 1991.
- [43] G. H. R. Santos. *Automatic Verification and Strategy Synthesis for Zero-Sum and Equilibria Properties of Concurrent Stochastic Games*. PhD thesis, University of Oxford, 2020.
- [44] Dana Scott. Outline of a Mathematical Theory of Computation. *Kiberneticheskij Sbornik. Novaya Seriya*, 14, January 1977.
- [45] Stephen Simons. Minimax Theorems and Their Proofs. In Ding-Zhu Du and Panos M. Pardalos, editors, *Minimax and Applications*. Springer US, 1995.
- [46] Anton Wijs, Joost-Pieter Katoen, and Dragan Bošnački. Efficient GPU algorithms for parallel decomposition of graphs into strongly connected and maximal end components. *Formal Methods in System Design*, 48(3), June 2016.

APPENDIX A
FURTHER DEFINITIONS AND CONCEPTS

A. Solving Matrix Games

Given a valuation v for all states in a two-player CSG, we can update the valuation at state s by solving a Linear Program (LP). Let the matrix game played at s be given by a matrix $Z \in \mathbb{Q}^{l \times m}$, where l (resp. m) is the number of actions available to player \mathcal{R} (resp. \mathcal{S}) at the state s . Then the LP that yields the value is the following [43]: Maximize $v(s)$ subject to the constraints:

$$\begin{aligned} v(s) &\leq x_1 \cdot z_{1j} + \cdots + x_l \cdot z_{lj} \text{ for } 1 \leq j \leq m \\ x_i &\geq 0 \text{ for } 1 \leq i \leq l \\ 1 &= x_1 + \cdots + x_l \end{aligned}$$

where $z_{ij} = Z(s)(i, j) = \sum_{s' \in S} \delta(s, a_i, b_j)(s') \cdot v(s')$, and x_i is the probability that player \mathcal{R} will take action i . Thus, by solving the LP, we not only obtain the value but also the optimal local strategy for player \mathcal{R} .

B. Domination of Strategies

Definition 38 (Domination [35, Def. 4.12]). Given a valuation v , a state $s \in S$ and sets of strategies $\mathcal{R}(s)$ and $\mathcal{S}(s)$ available at the state s for player \mathcal{R} and \mathcal{S} , respectively.

- A strategy $\rho \in \mathcal{R}(s)$ is *weakly dominated* if there exists another strategy $\rho' \in \mathcal{R}(s)$ satisfying the following two conditions:
 - $\inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(v)(s, \rho, \sigma) \leq \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(v)(s, \rho', \sigma)$, and
 - $\exists \sigma' \in \mathcal{S}(s)$ such that $\mathfrak{B}(v)(s, \rho, \sigma') < \mathfrak{B}(v)(s, \rho', \sigma')$.
- Dually, a strategy $\sigma \in \mathcal{S}(s)$ is *weakly dominated* if there exists another strategy $\sigma' \in \mathcal{S}(s)$ satisfying the following two conditions:
 - $\sup_{\rho \in \mathcal{R}(s)} \mathfrak{B}(v)(s, \rho, \sigma) \geq \sup_{\rho \in \mathcal{R}(s)} \mathfrak{B}(v)(s, \rho, \sigma')$, and
 - $\exists \rho' \in \mathcal{R}(s)$ such that $\mathfrak{B}(v)(s, \rho', \sigma) > \mathfrak{B}(v)(s, \rho', \sigma')$.

APPENDIX B
FURTHER EXAMPLES

A. Best Exit

Example 39 (Deflating BECs). Consider the CSG depicted in Fig. 2. At each state both players can choose among two actions. The game contains one MEC $C := \{s_0, s_1, s_2\}$. The matrix games played at the states s_0, s_1 and s_2 with respect to an upper bound U^k (where $k \in \mathbb{N}$) are defined as follows.

$$\begin{aligned} Z_{U^k}(s_0) &:= \begin{pmatrix} d_1 & d_2 \\ U^k(s_0) & U^k(s_1) \\ \frac{1}{2}(U^k(s_0) + \alpha) & \alpha \end{pmatrix} \begin{matrix} a_1 \\ a_2 \end{matrix}, & Z_{U^k}(s_1) &= \begin{pmatrix} e_1 & e_2 \\ \gamma & \beta \\ \frac{1}{2}(U^k(s_1) + U^k(s_2)) & U^k(s_0) \end{pmatrix} \begin{matrix} b_1 \\ b_2 \end{matrix}, \\ Z_{U^k}(s_2) &= \begin{pmatrix} f_1 & f_2 \\ U^k(s_0) & \frac{1}{2}(U^k(s_0) + U^k(s_1)) \\ \gamma & U^k(s_2) \end{pmatrix} \begin{matrix} c_1 \\ c_2 \end{matrix}. \end{aligned}$$

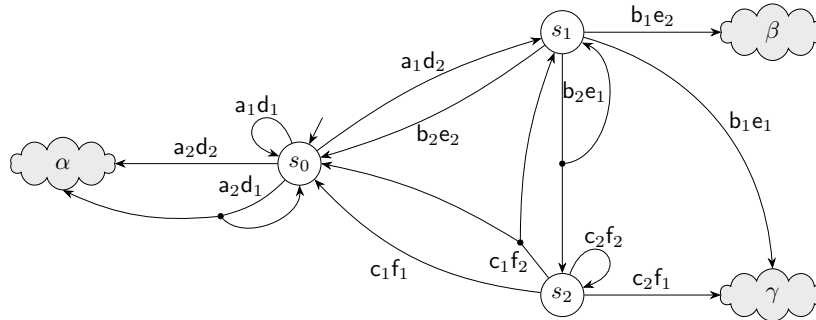


Fig. 2: CSG with non-single-state EC. The clouds represent the irrelevant parts of the game.

TABLE III: Full example of the BVI algorithm for the CSG depicted in Fig. 2. i is the i -th iteration of the main loop in Alg. 1, j is the j -th iteration of the while-loop in Alg. 3. C is the set of states among which a BEC \mathcal{X} is searched. The best exit from \mathcal{X} is underlined.

i	j	$U^i(s_0)$	$U^i(s_1)$	$U^i(s_2)$	C	\mathcal{X}
0	0	1	1	1	$\{s_0, s_1, s_2\}$	$\{s_0, s_1, \underline{s_2}\}$
	1	0.9	0.9	0.9	$\{s_0, s_1\}$	$\{s_0, \underline{s_1}\}$
	2	0.7	0.7	0.9	$\{s_0\}$	$\{\underline{s_0}\}$
	3	0.2	0.7	0.9	\emptyset	\emptyset
1	0	0.2	0.7	0.9	$\{s_0, s_1, s_2\}$	$\{s_0\}$
	1	0.2	0.7	0.9	$\{s_1, s_2\}$	$\{\underline{s_2}\}$
	2	0.2	0.7	0.45	$\{s_1\}$	\emptyset

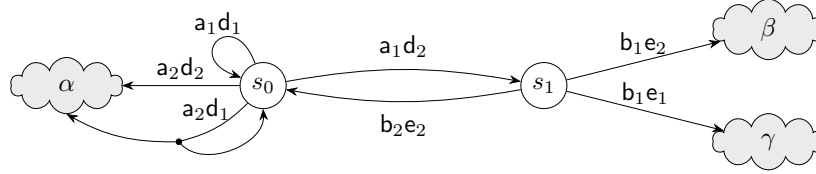


Fig. 3: The resulting game after removing s_2 from the game illustrated in Fig. 2.

For $k = 0$, $\alpha = 0.2$, $\beta = 0.7$ and $\gamma = 0.9$ and the initialization $U^0(s_0) = U^0(s_1) = U^0(s_2) = 1$, the matrix games look as follows.

$$Z_{U^0}(s_0) = \begin{pmatrix} d_1 & d_2 \\ 1 & 1 \\ 0.6 & 0.2 \end{pmatrix} \begin{matrix} a_1 \\ a_2 \end{matrix}, \quad Z_{U^0}(s_1) = \begin{pmatrix} e_1 & e_2 \\ 0.9 & 0.7 \\ 1 & 1 \end{pmatrix} \begin{matrix} b_1 \\ b_2 \end{matrix}, \quad Z_{U^0}(s_2) = \begin{pmatrix} f_1 & f_2 \\ 1 & 1 \\ 0.9 & 1 \end{pmatrix} \begin{matrix} c_1 \\ c_2 \end{matrix}.$$

Since at each state, there exist hazardous and trapping strategies, the three states form a BEC. To estimate the values of each exit, we need to solve three sub-matrix games played at each state of the BEC. The best exit value is then the maximum of the three solutions. The three linear programs that solve the three sub-matrix games are the following.

$$\begin{array}{lll} \max U^1(s_0) \text{ s.t.} & \max U^1(s_1) \text{ s.t.} & \max U^1(s_2) \text{ s.t.} \\ U^1(s_0) \leq 0.6 \cdot x_2 & U^1(s_1) \leq 0.9 \cdot x_1 & U^1(s_2) \leq 0.9 \cdot x_2 \\ U^1(s_0) \leq 0.2 \cdot x_2 & U^1(s_1) \leq 0.7 \cdot x_1 & U^1(s_2) \leq 1 \cdot x_2 \\ x_2 = 1 & x_1 = 1 & x_2 = 1 \end{array}$$

Here x_1 and x_2 are the probabilities that player \mathcal{R} chooses the first or second action available at the corresponding state. The best exit is $\max\{0.2, 0.7, 0.9\} = 0.9$, therefore, the upper bound of all the three states can be safely reduced to 0.9. \triangle

B. BVI Algorithm - Full Example

Example 40 (BVI). Consider the CSG depicted in Fig. 2, where $\alpha = 0.2$, $\beta = 0.7$, and $\gamma = 0.9$. The matrix games played at each state are given by the matrices $Z_{U^k}(s_0)$, $Z_{U^k}(s_1)$ and $Z_{U^k}(s_2)$ defined as in Example 39. We choose $\varepsilon = 0.001$.

Table III summarizes the steps of the algorithm. Initially it holds that $U^0(s_0) = U^0(s_1) = U^0(s_2) = 1$. Since $C = \{s_0, s_1, s_2\}$ is a MEC, it will be found by Alg. 3. Within C the BEC $\mathcal{X} = C$ is found. The best exit value from \mathcal{X} is calculated, i.e. the three sub-matrix games are solved and the best exit value needed for deflating the BEC is the maximum among the solutions:

$$Z_{U^0}^{\text{exit}}(s_0) := \begin{pmatrix} d_1 & d_2 \\ 0.6 & 0.2 \end{pmatrix} \begin{matrix} a_1 \\ a_2 \end{matrix}, \quad Z_{U^0}^{\text{exit}}(s_1) = \begin{pmatrix} e_1 & e_2 \\ 0.9 & 0.7 \end{pmatrix} \begin{matrix} b_1 \\ b_2 \end{matrix}, \quad Z_{U^0}^{\text{exit}}(s_2) = \begin{pmatrix} f_1 & f_2 \\ 0.9 & 1 \end{pmatrix} \begin{matrix} c_1 \\ c_2 \end{matrix}$$

The best exit from \mathcal{X} is s_2 (underlined in Table III) and the value of the exit is 0.9, so in the next step the over-approximations of those three states are deflated to 0.9 and the best exit is removed from C . The resulting sub-game is depicted in Fig. 3. For

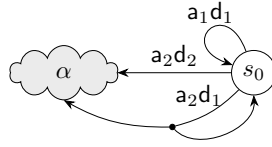


Fig. 4: The resulting game after removing s_1 from the game illustrated in Fig. 3.

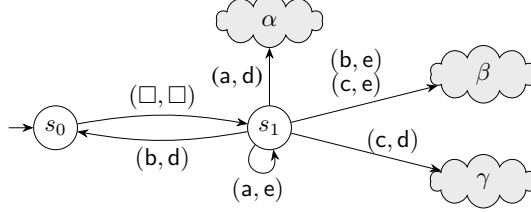


Fig. 5: Monotonicity is not guaranteed in general.

$U'(s_0) = U'(s_1) = U'(s_2) = 0.9$, the set $C = \{s_0, s_1\}$ contains the **BEC** $\mathcal{X} = C$. Now, two linear programs need to be solved to solve the two sub-matrix games.

$$Z_{U'}^{\text{exit}}(s_0) := \begin{matrix} & d_1 & d_2 \\ \begin{matrix} a_2 \end{matrix} & (0.55 & 0.2) \end{matrix} \quad Z_{U'}^{\text{exit}}(s_1) = \begin{matrix} & e_1 & e_2 \\ \begin{matrix} b_1 \end{matrix} & (0.9 & 0.7) \end{matrix}$$

The best exit is s_1 and the value of the exit 0.7, thus, the upper bounds of s_0 and s_1 are reduced to 0.7 and s_1 is removed from C . The resulting sub-game is depicted in Fig. 4. Finally, for $U''(s_0) = 0.7$, the set $C = \{s_0\}$ contains the **BEC** $\mathcal{X} = \{s_0\}$. To deflate it, we need to solve the following exiting sub-game.

$$Z_{U''}^{\text{exit}}(s_0) := \begin{matrix} & d_1 & d_2 \\ \begin{matrix} a_2 \end{matrix} & (0.45 & 0.2) \end{matrix}$$

The best exit value is 0.2. The upper upper bound of s_0 is reduced to 0.2. After removing the state from \mathcal{X} we obtain an empty set and the algorithm can proceed with the next **MEC**. Since, we assumed that the game has only one **MEC**, the deflating phase is finished. As $U^0 - L^0 > \varepsilon$ holds, the next iteration of the algorithm is executed.

The Bellman update returns the same valuation for all states, i.e. $U^1(s_0) = 0.2$, $U^1(s_1) = 0.7$ and $U^1(s_2) = 0.9$. Alg. 3 again finds the **MEC** $C = \{s_0, s_1, s_2\}$ that contains the two **BECs** $\mathcal{X}_1 := \{s_0\}$ and $\mathcal{X}_1 := \{s_2\}$. First $\mathcal{X} = \{s_0\}$ is deflated to 0.2 and next the **BEC** $\mathcal{X} = \{s_2\}$ is deflated. For this, the following exiting sub-game needs to be solved.

$$Z_{U^1}^{\text{exit}}(s_0) := \begin{matrix} & d_2 \\ \begin{matrix} a_2 \end{matrix} & (0.45) \end{matrix}$$

Notice, that here we only consider staying strategies for player \mathcal{S} , which is why we only consider action d_2 . Therefore, the best exit value is 0.45. After removing s_0 and s_2 from \mathcal{X} , no further **BECs** are contained in the **MEC**. Since now $U^1 - L^1 < \varepsilon$ holds, the **BVI** algorithm terminates. \triangle

C. Non-Monotonicity of Deflation

Example 41. Consider the **CSG** depicted in Fig. 5. The variables α, β and γ are placeholders indicating that upon leaving a certain valuation is obtained. The set of states $\{s_0, s_1\}$ is an **EC**. We consider two valuations, U and U' , such that $U \geq U'$ and the **EC** $\{s_0, s_1\}$ is for both upper bounds a **BEC**. The upper bound U assigns the following valuations: $U(s_0) = 0.6$, $U(s_1) = 0.6$, $U(\alpha) = 0.8$, $U(\beta) = 0.5$, $U(\gamma) = 0.55$. The upper bound U' assigns the following valuations: $U'(s_0) = 0.6$, $U'(s_1) = 0.45$, $U'(\alpha) = 0.5$, $U'(\beta) = 0.5$, $U'(\gamma) = 0.55$. Then, the matrix games played for the two valuations at state s_1 are given by the following matrices:

$$Z_U(s_1) := \begin{matrix} & d & e \\ \begin{matrix} a \\ b \\ c \end{matrix} & \begin{pmatrix} 0.8 & 0.6 \\ 0.6 & 0.5 \\ 0.55 & 0.5 \end{pmatrix} \end{matrix}, \quad Z_{U'}(s_1) := \begin{matrix} & d & e \\ \begin{matrix} a \\ b \\ c \end{matrix} & \begin{pmatrix} 0.5 & 0.45 \\ 0.6 & 0.5 \\ 0.55 & 0.5 \end{pmatrix} \end{matrix}$$

Then, under U the strategy $\{a \mapsto 1\}$ is a hazardous strategy and the exit value is 0.5. In contrast for U' the strategy $\{b \mapsto 1\}$ is hazardous and the exit value is 0.55. Thus, for arbitrary $U, U' \in \mathbb{R}^{|S|}$ it might happen that $\text{exitVal}_U(\mathcal{X}) < \text{exitVal}_{U'}(\mathcal{X})$ for some BEC \mathcal{X} although $U \geq U'$ holds. Intuitively, this arises because the sub-EC which forms a BEC changes when the relative ordering of exits is modified. However, the problem cannot occur when considering valid over-approximations (see Lem. 36) because then the upper bounds decrease in a well-behaved way when the relative ordering of the exits changes. \triangle

APPENDIX C PROOFS FOR SEC. III

Throughout the whole Appendix C, when proving convergence of BVI, we utilize definitions and theorems from lattice and fixpoint theory. Thus, we first briefly recall some necessary definitions (adjusting notation to our work) and theorems.

Definition 42 (Ordered set [13, Chapter 1.2]). A set P equipped with a relation $\preceq: P \times P$ is called an *ordered set* if and only if \preceq is reflexive, antisymmetric and transitive.

Definition 43 (Directed set [13, Chapter 7.7]). Let P be an ordered set. A non-empty set $D \subseteq P$ is *directed* if and only if for every pair of elements $x, y \in D$ there exists $z \in D$ that is a lower bound for both, formally $z \preceq x$ and $z \preceq y$.

Definition 44 (Complete partial order (CPO) [13, Chapter 8.1]). An ordered set P is a complete partially ordered set (CPO) if and only if

- (i) P has a top element, $\top := \inf_P \emptyset$, and
- (ii) for every directed set $D \subseteq P$, we have $\inf_P D$ exists.

Definition 45 (Continuity [13, Chapter 8.6]). Let P and Q be two CPOs. A mapping $\varphi: P \rightarrow Q$ is *continuous* if

- (i) for every directed set $D \subseteq P$, the subset $\varphi(D)$ of Q is also directed, and
- (ii) it holds that $\varphi(\inf D) = \inf \varphi(D) := \inf \{\varphi(x) \mid x \in D\}$.

Definition 46 (Order-preserving [13, Chapter 1.34]). Let P and Q be ordered sets. A map $\varphi: P \rightarrow Q$ is *order-preserving* (also called monotone) if $x \preceq y$ in P implies $\varphi(x) \preceq \varphi(y)$ in Q .

Theorem 47 (Fixpoint Theorem [13, Chapter 8.15]). Let P be a complete partial order, let F be an order-preserving and continuous self-map on P and define $\alpha := \sup_{n \geq 0} F^n(\top)$. Then α is the greatest fixpoint of F , i.e. the largest element of P satisfying $F(\alpha) = \alpha$.

We remark that we inverted the definitions and the theorem: This is because we are interested in a greatest fixpoint, whereas the textbook [13] only speaks about least fixpoints. Inverting the comparator and replacing \inf with \sup yields the original definitions. With only these changes, the proof of [13, Chapter 8.15] yields our modified claim Thm. 47.

A. Convergence in the absence of end components

We start by proving a technical lemma that is also useful for several future proofs: The over-approximation computed using only Bellman updates converges to a fixpoint.

Lemma 48 (Upper bound converges to a fixpoint). Let $(U^k)_{k \in \mathbb{N}_0}$ be the sequence of upper bounds computed by applying Eq. (3) on a CSG G . Let $U^* := \lim_{k \rightarrow \infty} U^k$ be the limit of the sequence. This limit is a fixpoint of the Bellman update, i.e. for all $s \in S$, $\mathfrak{B}(U^*)(s) = U^*(s)$.

Proof. This lemma is a consequence of the fixpoint theorem we just recalled. Thus, we proceed as follows: We explain that the domain of \mathfrak{B} is a CPO and prove that \mathfrak{B} is order-preserving and continuous. Then, Thm. 47 yields that U^* is a (namely the greatest) fixpoint.

Complete partial order. The domain of \mathfrak{B} are valuations, i.e. vectors $[0, 1]^{|S|}$ mapping every state to a number. Thus, we define the set P to be the set of all valuations. We use the standard point-wise comparisons as relation, i.e. $v_1 \preceq v_2$ if and only if for all states $s \in S$ we have $v_1(s) \leq v_2(s)$. Thus, the top element \top is the function that maps all states to 1. For every directed set D , a greatest lower bound $\ell = \inf_P D$ exists: Set $\ell(s) = \inf_{d \in D} d(s)$ for all $s \in S$. It is a lower bound, as by point-wise comparison, it is smaller than all valuations in D ; it is the greatest lower bound, since picking a larger number for any state would not be a lower bound any more. Thus, the set consisting of valuations $[0, 1]^{|S|}$ with this relation is a CPO.

Order-preserving. Recall that the Bellman operator on a state is defined as follows: $\mathfrak{B}(v)(s) := \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(v)(s, \rho, \sigma)$, where

$$\mathfrak{B}(v)(s, \rho, \sigma) := \sum_{(a, b) \in A} \sum_{s' \in S} v(s') \cdot \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b).$$

Eq. (3) lifts it to valuations by applying it state-wise. Hence, for every state, we apply an operation consisting of multiplications and summations, which are order-preserving. Thus, overall, the Bellman operator is order-preserving.

Continuous. We just showed that the Bellman operator on valuations is an order-preserving self-map on the set P of valuations. Then, [13, Lemma 8.7 (i)] yields that for every directed subset $D \subseteq P$, we have that $\mathfrak{B}(D) := \{\mathfrak{B}(d) \mid d \in D\}$ is a directed subset, which is Condition (i) of Def. 45. It remains to show Condition (ii): $\inf_{d \in D} \mathfrak{B}(d) = \mathfrak{B}(\inf D)$. Since the comparisons by the relation \preceq are performed point-wise, we have to prove that for all states $s \in S$, we have $\inf_{d \in D} \mathfrak{B}(d)(s) = \mathfrak{B}(\inf D)(s)$. Thus, fix an arbitrary state $s \in S$, and conclude using the following chain of equations.

$$\begin{aligned}
& \mathfrak{B}(\inf D)(s) \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \\
&\quad \left(\inf_{d \in D} d(s') \right) \cdot \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b) \quad (\text{Unfolding definition of Bellman operator}) \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \inf_{d \in D} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \\
&\quad d(s') \cdot \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b) \quad (\text{Claim 1}) \\
&= \inf_{d \in D} \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \\
&\quad d(s') \cdot \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b) \quad (\text{Claim 2}) \\
&= \inf_{d \in D} \mathfrak{B}(d)(s). \quad (\text{Collapsing the Bellman operator definition})
\end{aligned}$$

Claim 1: This step moves the $\inf_{d \in D}$ out of the summation, which is correct, since addition is a continuous operation. \blacktriangle

Claim 2: This step moves $\inf_{d \in D}$ to the front, first utilizing that infima can be switched. Then, to switch $\inf_{d \in D}$ and $\sup_{\rho \in \mathcal{R}(s)}$, we make use of the Minimax Theorem [45] which states that for a function $f : X \times Y \rightarrow \mathbb{R}$ that is concave-convex, it holds that $\sup_{x \in X} \inf_{y \in Y} f(x, y) = \inf_{y \in Y} \sup_{x \in X} f(x, y)$. f is concave-convex if f is concave for a fixed $y \in Y$ and convex for a fixed $x \in X$. This holds, in particular, for bilinear functions, i.e. functions that are linear in both arguments. The function considered at this point is the following:

$$\begin{aligned}
f(\rho, d) &= \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \\
&\quad d(s') \cdot \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b).
\end{aligned}$$

This function is indeed bilinear, since addition and multiplication are linear functions. \blacktriangle

Overall, we have shown that the sequence $(U^k)_{k \in \mathbb{N}_0}$ is the result of applying an order-preserving, continuous function to the top-element of a complete partial order, and thus it converges to a (the greatest) fixpoint. \square

Theorem 9 (Convergence without ECs — Proof in App. C-A). *Let G be a CSG where all ECs are trivial, i.e. for every EC C we have $C \subseteq W_{\mathcal{S}} \cup T$. Then, the over-approximation using only Eq. (3) converges, i.e. $\lim_{k \rightarrow \infty} U^k = V_{\mathcal{R}}$.*

Proof. This proof is an extension of the proof of [18, Theorem 1] for turn-based games to the concurrent setting. The underlying idea is the same, and can be briefly summarized as follows: We assume towards a contradiction that $U^* \neq V_{\mathcal{R}}$, and find a set \mathcal{X} that maximizes the difference between upper bound and value. Every pair of strategies leaving the set decreases the difference. However, $V_{\mathcal{R}}$ and U^* are fixpoints of the Bellman updates, from [17, Theorem 1] and Lem. 48, respectively. Consequently, optimal strategies need to remain in the set. However, in the absence of ECs, optimal strategies have to leave the set, which yields a contradiction and proves that $U^* = V_{\mathcal{R}}$.

Main challenge. The key difference to the proof of [18, Theorem 1] is that we cannot argue about actions anymore, but have to consider mixed strategies. This significantly complicates notation. Additionally, and more importantly, the former proof crucially relied on the fact that for a state of player \mathcal{R} , we know that its valuation is at least as large as that of any action, and dually for a state of player \mathcal{S} , its valuation is at most as large as that of any action. In the concurrent setting, this is not true. The optimal strategies need not be maximizing nor minimizing the valuation and, moreover, they can be maximizing for one valuation and minimizing for another. Thus, we found a more general, and in fact simpler, way of proving that “no state in \mathcal{X} can depend on the outside” [18, Statement 5] and deriving the contradiction. The crucial insight is that we can fix locally optimal strategies and then apply Claim 4.

Notation for Bellman operator. Before we begin the formal proof, we establish a condensed notation for a number of terms in the Bellman operator:

$$\text{rest}(s, a, b, s', \rho, \sigma) = \delta(s, a, b)(s') \cdot \rho(s)(a) \cdot \sigma(s)(b).$$

Thus, the Bellman operator for some valuation v and pair of strategies (ρ, σ) simplifies to

$$\mathfrak{B}(v)(s, \rho, \sigma) = \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} v(s') \cdot \text{rest}(s, a, b, s', \rho, \sigma).$$

The set \mathcal{X} with maximum difference. We define the *difference of a state* $s \in S$ as $\Delta(s) := U^*(s) - V_{\mathcal{R}}(s)$. Recall that $V_{\mathcal{R}}$ is the least fixpoint and U^* the greatest fixpoint of \mathfrak{B} . Hence, we know that $\Delta(s) \geq 0$ for all states. Further, since we assume for contradiction that $U^* \neq V_{\mathcal{R}}$, there exist states with $\Delta(s) > 0$. Thus, we can find a non-empty set of states with maximum difference: $\mathcal{X} := \{s \in S \mid \Delta(s) = \max_{s \in S} \Delta(s)\}$.

A leaving pair of strategies decreases the difference. Let $s \in \mathcal{X}$ be a state in \mathcal{X} . Let $(\rho, \sigma) \in (\mathcal{R}(s) \times \mathcal{S}(s))$ be a pair of strategies such that (s, ρ, σ) *leaves* \mathcal{X} . Then, following this pair of strategies for one step decreases the difference, formally

$$\mathfrak{B}(\Delta)(s, \rho, \sigma) < \Delta(s). \quad (5)$$

We prove this using the following chain of equations:

$$\begin{aligned} & \mathfrak{B}(\Delta)(s, \rho, \sigma) \\ &= \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \Delta(s') \cdot \text{rest}(s, a, b, s', \rho, \sigma) && \text{(Definition of Bellman operator)} \\ &< \Delta(s). && \text{(Claim 3)} \end{aligned}$$

Claim 3: By assumption, we have that there exists a $t \in S \setminus \mathcal{X}$ such that this t is reached with positive probability under the exiting strategies, i.e. $\sum_{(a,b) \in \mathcal{A}} \text{rest}(s, a, b, t, \rho, \sigma) > 0$. For this t outside of \mathcal{X} , we have $\Delta(t) < \Delta(s)$, since \mathcal{X} is defined as the set of all states with maximum difference. Further, no state can have a difference larger than $\Delta(s)$. Furthermore, the remaining terms in the sum that the differences $\Delta(s')$ are multiplied with are a probability distribution, formally $\sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \text{rest}(s, a, b, s', \rho, \sigma) = 1$. Thus, if all differences were equal to $\Delta(s)$, the sum would yield $\Delta(s)$. As one of the summands is smaller than $\Delta(s)$ and all others are at most $\Delta(s)$, we get that the sum has to be smaller than $\Delta(s)$. \blacktriangle

Without non-trivial ECs, \mathcal{X} must be left. We have that $\mathcal{X} \cap (W_{\mathcal{S}} \cup T) = \emptyset$: The difference is 0 for target states because both the value and the upper bound are equal to 1; and the difference is 0 for the sure winning region of player \mathcal{S} , since the upper bound and value are equal to 0 (see Eq. (3)). Thus, since by assumption there are no ECs in $S \setminus (W_{\mathcal{S}} \cup T)$, the set \mathcal{X} cannot contain an EC. Consequently, there exists a state $s \in \mathcal{X}$ such that for all pairs of available actions $(a, b) \in \Gamma_{\mathcal{R}}(s) \times \Gamma_{\mathcal{S}}(s)$, we have $\text{Supp}(\delta(s, a, b)) \cap (S \setminus \mathcal{X}) \neq \emptyset$, i.e. there is a successor state outside of \mathcal{X} . This is the case because, if all states had a pair of actions that stays in \mathcal{X} , then there exists a pair of strategies that keeps the play inside a subset of \mathcal{X} , which would then form an EC. For a formal proof, we refer to [18, Lemma 2]. Note that, while their proof is for turn-based games, the definition of EC is a graph theoretic notion where, intuitively, the players “work together” (formally, it is only about the existence of an edge in the underlying hypergraph), and thus the proof is applicable to CSGs, too. In the following, we let s denote such a state where all strategies are leaving.

Notation for locally optimal strategies. For any state and valuation, locally optimal strategies for both players exist. We establish a shorthand for the locally optimal strategies in state s (the one obtained in the previous step) with respect to U^* and $V_{\mathcal{R}}$. For player \mathcal{R} and valuation U^* , we denote a locally optimal strategy by

$$\rho_U \in \arg \max_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(U^*)(s, \rho, \sigma).$$

Similarly we denote a locally optimal strategy of player \mathcal{S} with respect to U^* by

$$\sigma_U \in \arg \min_{\sigma \in \mathcal{S}(s)} \sup_{\rho \in \mathcal{R}(s)} \mathfrak{B}(U^*)(s, \rho, \sigma).$$

Analogously, we define locally optimal strategies with respect to $V_{\mathcal{R}}$, namely ρ_V and σ_V , obtained by replacing U^* with $V_{\mathcal{R}}$ in the above definition.

Deriving the contradiction. Recall that $s \in \mathcal{X}$ is a state where all available pairs of actions, and thus all strategies, leave \mathcal{X} . We derive the contradiction $\Delta(s) < \Delta(s)$.

$$\begin{aligned}
\Delta(s) &= U^*(s) - V_{\mathcal{R}}(s) && \text{(Definition of } \Delta) \\
&= \mathfrak{B}(U^*)(s) - V_{\mathcal{R}}(s) && (U^* \text{ is fixpoint by Lem. 48}) \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_U) - V_{\mathcal{R}}(s) && ((\rho_U, \sigma_U) \text{ locally optimal w.r.t. } U^*) \\
&\leq \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - V_{\mathcal{R}}(s) && \text{(Claim 4)} \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s) && (V_{\mathcal{R}} \text{ is fixpoint [17, Theorem 1]}) \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s, \rho_V, \sigma_V) && ((\rho_V, \sigma_V) \text{ locally optimal w.r.t. } V_{\mathcal{R}}) \\
&\leq \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s, \rho_U, \sigma_V) && \text{(Claim 4)} \\
&= \mathfrak{B}(\Delta)(s, \rho_U, \sigma_V) && \text{(Claim 5)} \\
&< \Delta(s). && ((s, \rho_U, \sigma_V) \text{ leaves } \mathcal{X} \text{ and Eq. (5)})
\end{aligned}$$

Claim 4: This argument is used in two steps in the above chain of equations. For the first usage, observe that σ_V can be at most as good as the optimal σ_U . More formally, recall σ_U was chosen as the arg min of the Bellman operator with respect to U^* . Thus, $\mathfrak{B}(U^*)(s, \rho_U, \sigma_U) \leq \mathfrak{B}(U^*)(s, \rho_U, \sigma_V)$.

For the second usage, by the analogous argument we have $\mathfrak{B}(V_{\mathcal{R}})(s, \rho_V, \sigma_V) \geq \mathfrak{B}(V_{\mathcal{R}})(s, \rho_U, \sigma_V)$. Since this term is the subtrahend of the subtraction, the overall expression can only become greater. \blacktriangle

Claim 5: This step follows from expanding the definition of the Bellman operator, rearranging the sums and collapsing the definition of Bellman operator. Formally, for all states q and strategy pairs (ρ, σ) it holds that

$$\begin{aligned}
\mathfrak{B}(\Delta)(q, \rho, \sigma) &= \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} \Delta(s') \cdot \text{rest}(s, a, b, s', \rho, \sigma) && \text{(Definition of Bellman operator)} \\
&= \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} (U^*(s') - V_{\mathcal{R}}(s')) \cdot \text{rest}(s, a, b, s', \rho, \sigma) && \text{(Definition of } \Delta) \\
&= \left(\sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} U^*(s') \cdot \text{rest}(s, a, b, s', \rho, \sigma) \right) - \\
&\quad \left(\sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} V_{\mathcal{R}}(s') \cdot \text{rest}(s, a, b, s', \rho, \sigma) \right) && \text{(Splitting the sum)} \\
&= \mathfrak{B}(U^*)(q, \rho, \sigma) - \mathfrak{B}(V_{\mathcal{R}})(q, \rho, \sigma). && \text{(Definition of Bellman operator)}
\end{aligned}$$

\blacktriangle

Summary. Starting from the assumption that $U^* \neq V_{\mathcal{R}}$, we derived that there exists a set of states \mathcal{X} where the difference Δ between upper bound and value is maximized. Further, a pair of strategies leaving \mathcal{X} decreases this difference. However, since there are no ECs in \mathcal{X} , there has to be a state where the optimal strategies for $V_{\mathcal{R}}$ and U^* leave, which allows us to derive a contradiction. Thus, the initial assumption is false, and we have $U^* = V_{\mathcal{R}}$. \square

B. Convergence without Bloated End Components

Lemma 14 (Negating Weak Domination — Proof in App. C-B). *Let v be a valuation, $s \in S$ a state, $\mathcal{R}_1, \mathcal{R}_2, \mathcal{R}' \subseteq \mathcal{R}(s)$ and $\mathcal{S}_1, \mathcal{S}_2, \mathcal{S}' \subseteq \mathcal{S}(s)$ sets of local strategies.*

If for some sets of strategies we do not have $\mathcal{R}_2 \prec_{v, \mathcal{S}'} \mathcal{R}_1$, then we have $\mathcal{R}_1 \preceq_{v, \mathcal{S}'} \mathcal{R}_2$. Analogously, not $\mathcal{S}_2 \prec_{v, \mathcal{R}'} \mathcal{S}_1$ implies $\mathcal{S}_1 \preceq_{v, \mathcal{R}'} \mathcal{S}_2$.

Proof. We only provide the proof for Player \mathcal{R} , as the other one is analogous by exchanging the names of the strategy sets, replacing \mathcal{R} with \mathcal{S} and vice versa.

We assume that we do not have $\mathcal{R}_2 \prec_{v, \mathcal{S}'} \mathcal{R}_1$ under the set of counter-strategies \mathcal{S}' with respect to v . Writing out the definition, this means that $\forall \rho_1 \in \mathcal{R}_1. \exists \rho_2 \in \mathcal{R}_2$:

- (i) $\inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_2, \sigma) > \inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_1, \sigma)$, or
- (ii) $\forall \sigma' \in \mathcal{S}'$ we have $\mathfrak{B}(v)(s, \rho_2, \sigma') \geq \mathfrak{B}(v)(s, \rho_1, \sigma')$.

Our goal is to have that $\mathcal{R}_1 \preceq_{v, \mathcal{S}'} \mathcal{R}_2$, formally: Formally, $\exists \rho_2 \in \mathcal{R}_2. \forall \rho_1 \in \mathcal{R}_1$:

$$\inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_1, \sigma) \leq \inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(v)(s, \rho_2, \sigma).$$

Both conditions of negated weak dominance imply our goal. The only remaining problem is the order of quantifiers. However, the choice of ρ_2 does not depend on ρ_1 , and we can always pick ρ_2 as the strategy that maximizes $\inf_{\sigma \in \mathcal{S}'} \mathfrak{B}(\nu)(s, \rho_2, \sigma)$. Thus, we can exchange the order of quantifiers and prove our goal. \square

Lemma 20 (Negating Bloated — Proof in App. C-B). *If an EC $\mathcal{X} \subseteq \mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ is not bloated for a valuation ν , then there exists a state $s \in \mathcal{X}$ that has a locally optimal strategy that is leaving, formally $\exists \rho \in \mathcal{R}_L(\mathcal{X}, s). \mathcal{R}(s) \preceq_{\nu, \mathcal{S}(s)} \{\rho\}$.*

Proof. Since \mathcal{X}' is not a BEC, we know that at some $s \in \mathcal{X}'$ it holds that $\text{Hazard}_{U^*}(\mathcal{X}', s) = \emptyset$. Fix s to be such a state. Every strategy $\rho' \in \mathcal{R}(\mathcal{S})$ must violate at least one of the 3 conditions of Def. 16. We write out the negations:

- (i) $\rho' \notin \mathcal{R}_L(\mathcal{X}', s)$, i.e. the strategy is leaving $\rho' \in \mathcal{R}_L(\mathcal{X}', s)$.
- (ii) We do not have $\mathcal{R}(s) \setminus \{\rho'\} \preceq_{\mathcal{S}(s)} \{\rho'\}$. By contraposition of Lem. 14, this implies $\{\rho'\} \prec_{\mathcal{S}(s)} \mathcal{R}(s) \setminus \{\rho'\}$, so the strategy is sub-optimal.
- (iii) We do not have $\mathcal{R}_L(\mathcal{X}', s) \prec_{\mathcal{S}(s)} \{\rho'\}$. By Lem. 14, this implies $\{\rho'\} \preceq_{\mathcal{S}(s)} \mathcal{R}_L(\mathcal{X}', s)$, i.e. there are leaving strategies that are not worse than ρ' . Note that in particular, this implies that $\mathcal{R}_L(\mathcal{X}', s)$ is non-empty.

Our assumption gives us the disjunction over the three violated conditions. We proceed by a case distinction, always assuming that all strategies violate a certain condition, which allows us to prove our goal, or, if there exists a strategy satisfying the condition, we continue with the next one.

Case “Not (i)”: If all strategies violate Condition (i), that means all strategies are leaving, i.e. $\mathcal{R}_L(s) = \mathcal{R}(s)$. Thus, since there always exist optimal strategies, by picking an optimal $\rho \in \mathcal{R}$ we naturally have $\mathcal{R}(s) \preceq_{\nu, \mathcal{S}(s)} \{\rho\}$.

Case “(i), but not (iii)”: We assume there exists strategies that satisfies Condition (i), but all strategies that satisfy it violate Condition (iii). This means that for every non-leaving strategy, the set of leaving strategies is not worse than it. As this holds for all non-leaving strategies, we have $\mathcal{R}_L(s) \preceq_{\nu, \mathcal{S}(s)} \mathcal{R}_L(s)$. Using that $\mathcal{R}(s) = \mathcal{R}_L(s) \cup \mathcal{R}_L(s)$ and the set of leaving strategies trivially is not worse than itself, we obtain: $\mathcal{R}(s) \preceq_{\nu, \mathcal{S}(s)} \mathcal{R}_L(s)$. Moreover, the set of leaving strategies is non-empty, since the definition of not worse requires that there exists a strategy in the right-hand set. This proves our goal.

Case “(i) and (iii), but not (ii)”: We assume there exists strategies that satisfy Condition (i) and (iii), but all these strategies violate Condition (ii). This case cannot happen, and below we derive a contradiction. This completes our cast distinction, since every strategy has to violate at least one of the three conditions.

Our assumption is that there exists a non-leaving strategy that weakly dominates the set of all leaving strategies. However, every such strategy is suboptimal, as by violating Condition (ii) it is weakly dominated by all other strategies. This is a contradiction, because then there are no optimal strategies. More formally, if we assume the optimal strategy is leaving, this is a contradiction, because there exists a non-leaving strategy dominating the set of leaving strategies. And if we assume the optimal strategy is non-leaving, this is a contradiction, because every non-leaving strategy is weakly dominated by the set of others. \square

Theorem 21 (Non-convergence implies BECs — Proof in App. C-B). *Let $U^* := \lim_{k \rightarrow \infty} U^k$ be the limit of the naïve upper bound iteration (Eq. (3)) on the CSG G . If VI from above does not converge to the value in the limit, i.e. $U^* > V_{\mathcal{R}}$, then the CSG G contains a BEC in $\mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$ with respect to U^* .*

Proof. Intuition and outline: This proof builds on the proof of Thm. 9. There, we constructed a set \mathcal{X} maximizing the difference between U^* and $V_{\mathcal{R}}$ and showed that if there is a pair of optimal strategies leaving leaving \mathcal{X} , then we can derive a contradiction: The upper bound decreases, which contradicts the fact that it is a fixpoint. In the context of the other proof, that allowed us to show that without ECs, VI converges, because without ECs it is impossible to have a set of states where all optimal strategies stay in that set.

In the presence of ECs, states can indeed have a positive difference between U^* and $V_{\mathcal{R}}$, see e.g. Ex. 6. Our goal is to prove that at least one of these ECs is bloated. Thus, we assume for contradiction that no EC is bloated under U^* . Thus, by Lem. 20, there is an optimal leaving strategy for player \mathcal{R} . Using that, we can repeat the argument from Thm. 9, showing that in this case U^* would decrease. Again, this is a contradiction because it is a fixpoint of applying Bellman updates (Lem. 48). Thus, the initial assumption that no EC is bloated is false, and we can conclude that there exists a BEC.

Establishing the Context. As in the proof of Thm. 9, let $\mathcal{X} := \{s \in \mathcal{S} \mid \Delta(s) = \max_{s \in \mathcal{S}} \Delta(s)\}$ be the set of states with maximum difference. We denote the maximum by Δ^{\max} , and our assumption yields that $\Delta^{\max} > 0$. Note that this implies that $\mathcal{X} \cap (\mathcal{W}_{\mathcal{S}} \cup \mathcal{T}) = \emptyset$, since for those states, their value is set correctly by initialization, and their difference is 0. Thus, if we find a BEC that is a subset of \mathcal{X} , it also satisfies the additional condition of being non-trivial, i.e. not in $(\mathcal{W}_{\mathcal{S}} \cup \mathcal{T})$. The contraposition of Thm. 9 yields that there has to be an EC in \mathcal{X} .

Bottom MECs. To derive the contradiction, in the following, we consider ECs with a particular property, namely ECs that are bottom in \mathcal{X} . A MEC \mathcal{X}' is bottom in \mathcal{X} if the successors of a pair of strategies that leaves the MEC reaches states outside of \mathcal{X} with positive probability. Intuitively, a bottom MEC in \mathcal{X} is a MEC, such that after leaving it, none of the successors is

part of another MEC in \mathcal{X} . One can compute such ECs using the MEC decomposition of \mathcal{X} , ordering them topologically and picking one at the end of a chain.

Let \mathcal{X}' be a bottom MEC with $\mathcal{X}' \subseteq \mathcal{X}$. \mathcal{X}' exists because by assumption \mathcal{X} contains at least one EC, thus, there also has to exist an EC that is bottom in \mathcal{X} .

Optimal Leaving Strategies in Non-BECs. We use the assumption for contradiction to say that \mathcal{X}' is not bloated with respect to U^* . Then, using Lem. 20, we know that there exists a state $s \in \mathcal{X}'$ where an optimal strategy ρ_U exists that is leaving \mathcal{X}' . Moreover, since \mathcal{X}' is a bottom MEC in \mathcal{X} , we also have that it is leaving with respect to \mathcal{X} .

Deriving the Contradiction. Using these facts, we can exactly repeat the argument used in the proof of Thm. 9 under the paragraph-heading “Deriving the Contradiction”. Recall we denote locally optimal strategies with respect to U^* by ρ_U, σ_U (and we just proved ρ_U is leaving for all counter-strategies), and analogously locally optimal strategies with respect to $V_{\mathcal{R}}$, by ρ_V and σ_V . We highlight that Eq. (5), Claim 4 and Claim 5 from Thm. 9 are applicable in the context of this proof, too.

$$\begin{aligned}
\Delta(s) &= U^*(s) - V_{\mathcal{R}}(s) && \text{(Definition of } \Delta) \\
&= \mathfrak{B}(U^*)(s) - V_{\mathcal{R}}(s) && (U^* \text{ is fixpoint by Lem. 48}) \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_U) - V_{\mathcal{R}}(s) && ((\rho_U, \sigma_U) \text{ locally optimal w.r.t. } U^*) \\
&\leq \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - V_{\mathcal{R}}(s) && \text{(Claim 4 in Thm. 9)} \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s) && (V_{\mathcal{R}} \text{ is fixpoint [17, Theorem 1]}) \\
&= \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s, \rho_V, \sigma_V) && ((\rho_V, \sigma_V) \text{ locally optimal w.r.t. } V_{\mathcal{R}}) \\
&\leq \mathfrak{B}(U^*)(s, \rho_U, \sigma_V) - \mathfrak{B}(V_{\mathcal{R}})(s, \rho_U, \sigma_V) && \text{(Claim 4 in Thm. 9)} \\
&= \mathfrak{B}(\Delta)(s, \rho_U, \sigma_V) && \text{(Claim 5 in Thm. 9)} \\
&< \Delta(s). && ((s, \rho_U, \sigma_V) \text{ leaves } \mathcal{X} \text{ and Eq. (5)})
\end{aligned}$$

Now that we have derived a contradiction, our initial assumption that all ECs are not BECs is wrong, so we know there exists a BEC $\mathcal{X}' \subseteq (S \setminus (T \cup W_{\mathcal{J}}))$ with respect to U^* . This concludes the proof.

As a side note, we remark that it is indeed possible that there is only one BEC that causes many states, even some not in an EC, to have a positive difference Δ . For an example, we refer to [18, Fig. 4].

□

C. Soundness of DEFLATE

The following is a technical lemma that is needed to show the soundness of deflation. Intuitively, it says that the value of all states in an EC needs to depend on some exit. Note that there can be states whose value is higher than their own exit value, namely if they can reach a better exit. However, this cannot be the case for all states, but there must be some whose value is less than or equal to their exit value (first condition), and in fact no state can have a higher value than these states that actually depend on exiting (second condition). The proof is very technical, as essentially it requires unfolding all the definitions, and thereby also unfolding all the included case distinctions.

Lemma 49 (No state has a larger value than that of an exit from its EC). *Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{J}})$ be an EC. Then, it holds that*

- (i) $\mathcal{X}' := \{s \in \mathcal{X} \mid V_{\mathcal{R}}(s) \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)\} \neq \emptyset$, and
- (ii) $\max_{s \in \mathcal{X}'} V_{\mathcal{R}}(s) \geq \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$.

Proof. We prove the lemma by contradiction, i.e. we assume that one of the two conditions posed by the lemma is violated. We make the following case distinction:

- (i) $\mathcal{X}' = \emptyset$; and
- (ii) $\mathcal{X}' \neq \emptyset$ but $\max_{s \in \mathcal{X}'} V_{\mathcal{R}}(s) < \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$.

Case (i)

In this case it holds that $\mathcal{X}' = \emptyset$, i.e. for all $s \in \mathcal{X}$ we have $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$. Recall that in Def. 27 if $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ is true at some state $s \in \mathcal{X}$, then the exit value is given by the value of the matrix game played at that state.

Consequently, as for all $s \in \mathcal{X}$ it holds that $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$, at all states $s \in \mathcal{X}$ it must hold that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$ because otherwise we would obtain the following contradiction: $\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) = V_{\mathcal{R}}(s)$.

We proceed with the assumption that for all $s \in \mathcal{X}$ it holds that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$. Since by the case assumption at all $s \in \mathcal{X}$ it holds that $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$ the true values of the states in \mathcal{X} are attainable with the hazardous and trapping strategies. More formally, at each $s \in \mathcal{X}$ the following chain of equations holds.

$$\begin{aligned}
V_{\mathcal{R}}(s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (V_{\mathcal{R}} \text{ is a fixpoint \& case assumption: } V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \text{ are optimal by Def. 16}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma). && (\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \text{ are optimal by Def. 24})
\end{aligned}$$

Thus, for player \mathcal{R} staying in \mathcal{X} is optimal and since no target state is contained in \mathcal{X} it has to hold that $V_{\mathcal{R}}(s) = 0$ for all $s \in \mathcal{X}$. However, this is a contradiction to the assumption that $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$ since $\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) \geq 0$ (as all possible valuations are non-negative and the value of the exiting sub-game is given by the maximum between 0 and exitVal - see Def. 27 and Def. 26). Thus, the case assumption that $\mathcal{X}' = \emptyset$ must be false.

Case (ii)

In this case it holds that $\mathcal{X}' \neq \emptyset$ but $\max_{s \in \mathcal{X}'} V_{\mathcal{R}}(s) < \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ is true.

We make the following case distinction: (ii.a) $\mathcal{X} \setminus \mathcal{X}'$ is a **BEC**; and (ii.b) $\mathcal{X} \setminus \mathcal{X}'$ is not a **BEC**.

Case (ii.a)

In this case $\mathcal{X} \setminus \mathcal{X}'$ is a **BEC**, i.e. for all $s \in \mathcal{X} \setminus \mathcal{X}'$ it holds that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s) \neq \emptyset$.

In case it holds that $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s)$ for all $s \in \mathcal{X} \setminus \mathcal{X}'$, then since no target state is contained in \mathcal{X} and therefore neither in $\mathcal{X} \setminus \mathcal{X}'$, it must be true that $V_{\mathcal{R}}(s) = 0$ for all $s \in \mathcal{X} \setminus \mathcal{X}'$. However, similarly as in Case (i), this is a contradiction to the assumption that $V_{\mathcal{R}}(s) > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s)$ for all $s \in \mathcal{X} \setminus \mathcal{X}'$ as $\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s) \geq 0$ for all $s \in \mathcal{X} \setminus \mathcal{X}'$.

Consequently, there has to exist $s \in \mathcal{X} \setminus \mathcal{X}'$ such that $V_{\mathcal{R}}(s) \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s)$.

We make another case distinction:

- (ii.a.1) for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $V_{\mathcal{R}}(s') > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')$; and
- (ii.a.2) there exists $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ such that $V_{\mathcal{R}}(s') \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')$.

Case (ii.a.1) In this case for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $V_{\mathcal{R}}(s') > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')$.

Thus, at each $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s')$ the following chain of equations holds.

$$\begin{aligned}
V_{\mathcal{R}}(s') &= \sup_{\rho \in \mathcal{R}(s')} \inf_{\sigma \in \mathcal{S}(s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma) && (V_{\mathcal{R}} \text{ is a fixpoint}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \mathcal{S}(s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma) && (\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \text{ are optimal by Def. 16}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma). && (\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \text{ are optimal by Def. 24})
\end{aligned}$$

Thus, at all states that attain the highest value among $\mathcal{X} \setminus \mathcal{X}'$, it is optimal for both players to choose strategies that together are staying in $\mathcal{X} \setminus \mathcal{X}'$. However, since \mathcal{X} and thus $\mathcal{X} \setminus \mathcal{X}'$ does not belong to the winning region of Player \mathcal{S} , leaving $\mathcal{X} \setminus \mathcal{X}'$ has to be possible. Thus, from each state $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$, Player \mathcal{R} has to possess an optimal strategy that leads to a state $s'' \in \mathcal{X} \setminus \mathcal{X}'$ where leaving $\mathcal{X} \setminus \mathcal{X}'$ is possible, i.e. where $V_{\mathcal{R}}(s'') \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s'')$ holds. Thus, at state s'' the highest value also has to be attainable. However, this is a contradiction to the case assumption that for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $V_{\mathcal{R}}(s') > \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')$.

Case (ii.a.2) In this case there exists $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ such that $V_{\mathcal{R}}(s') \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')$. Let $\mathcal{X}'' := \{s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s) \mid V_{\mathcal{R}}(s') \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')\}$. Then, we need to make another case distinction:

- (ii.a.2.1) for all $s \in \mathcal{X}''$ it holds that $\text{Deflv}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s) = \emptyset$; and
- (ii.1.2.2) there exists $s' \in \mathcal{X}''$ such that $\text{Deflv}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset$.

Case (ii.a.2.1) In this case for all $s \in \mathcal{X}''$ it holds that $\text{Deflv}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s) = \emptyset$. Then, since \mathcal{X} and so $\mathcal{X} \setminus \mathcal{X}'$ do not belong to the winning region of Player \mathcal{S} , there has to exist a state $s'' \in \mathcal{X} \setminus \mathcal{X}'$ such that $\text{Deflv}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s'') \neq \emptyset$

so leaving $\mathcal{X} \setminus \mathcal{X}'$ is possible. However, then $s'' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ must hold which is a contradiction to the case assumption.

Case (ii.a.2.2) In this case there exists $s' \in \mathcal{X}''$ such that $\text{Deflv}_{\mathcal{R}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset$. Then, the following chain of equations holds.

$$\begin{aligned}
V_{\mathcal{R}}(s') &\leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') && \text{(Case assumption)} \\
&= \max(0, V(\hat{Z}(s'))) && \text{(By case assumption (ii.a): Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset) \\
&= V(\hat{Z}(s')) && \text{(Case assumption: Deflv}_{\mathcal{R}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset) \\
&= \sup_{\rho \in \text{Deflv}_{\mathcal{R}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma) \\
&\quad \text{(Value of the exiting sub-game — see Def. 27)} \\
&= \sup_{\rho \in \text{Deflv}_{\mathcal{R}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \sum_{(a,b) \in \mathcal{A}} \sum_{s \in \mathcal{X} \setminus \mathcal{X}' \atop s \leq \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \underbrace{V_{\mathcal{R}}(s)}_{\leq \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \cdot \delta(s', a, b)(s) \cdot \rho(a) \cdot \sigma(b) \\
&+ \sum_{s \in \mathcal{X}' \atop s < \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \underbrace{V_{\mathcal{R}}(s)}_{< \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \cdot \delta(s', a, b)(s) \cdot \rho(a) \cdot \sigma(b) \\
&\quad \text{(Def. of } \mathfrak{B} \text{ and case assumption (ii): } \max_{s \in \mathcal{X}'} V_{\mathcal{R}}(s) < \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)) \\
&< \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s). && \text{(Everything sums up to 1)}
\end{aligned}$$

Thus, from the case assumption we derived a contradiction because $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$.

Case (ii.b)

In this case $\mathcal{X} \setminus \mathcal{X}'$ is not a BEC, i.e. there exists $s' \in \mathcal{X} \setminus \mathcal{X}'$ such that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') = \emptyset$.

We need another case distinction:

(ii.b.1) for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset$; and

(ii.b.2) there exists $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ such that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') = \emptyset$.

Case (ii.b.1) In this case we have that for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset$.

Then, at each $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ the following chain of equations holds.

$$\begin{aligned}
V_{\mathcal{R}}(s') &= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \mathcal{S}(s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma) && \text{(Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \text{ are optimal by Def. 16)} \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma). && \text{(Trap}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \text{ are optimal by Def. 24)}
\end{aligned}$$

Thus, at all states that attain the highest value among $\mathcal{X} \setminus \mathcal{X}'$ it is optimal for both players to choose strategies that together are staying in $\mathcal{X} \setminus \mathcal{X}'$. However, since \mathcal{X} and thus $\mathcal{X} \setminus \mathcal{X}'$ does not belong to the winning region of Player \mathcal{S} , leaving $\mathcal{X} \setminus \mathcal{X}'$ has to be possible. Thus, from each state $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$, Player \mathcal{R} has to possess an optimal strategy that leads to a state $s'' \in \mathcal{X} \setminus \mathcal{X}'$ where leaving $\mathcal{X} \setminus \mathcal{X}'$ is possible, i.e. where $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s'') = \emptyset$ holds. Thus, at state s'' the highest value also has to be attainable. However, this is a contradiction to the case assumption that for all $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ it holds that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') \neq \emptyset$.

Case (ii.b.2) In this case there exists $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$ such that $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') = \emptyset$.

Then, the following chain of equations holds.

$$\begin{aligned}
V_{\mathcal{R}}(s') &= \sup_{\rho \in \mathcal{R}(s')} \inf_{\sigma \in \mathcal{S}(s')} \mathfrak{B}(V_{\mathcal{R}})(s', \rho, \sigma) && (V_{\mathcal{R}} \text{ is a fixpoint}) \\
&= \sup_{\rho \in \mathcal{R}(s')} \inf_{\sigma \in \mathcal{S}(s')} \sum_{(a,b) \in \mathcal{A}} \sum_{s \in \mathcal{X} \setminus \mathcal{X}' \atop s \leq \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \underbrace{V_{\mathcal{R}}(s)}_{\leq \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \cdot \delta(s', a, b)(s) \cdot \rho(a) \cdot \sigma(b) \\
&+ \sum_{s \in \mathcal{X}' \atop s < \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \underbrace{V_{\mathcal{R}}(s)}_{< \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'')} \cdot \delta(s', a, b)(s) \cdot \rho(a) \cdot \sigma(b) \\
&< \max_{s'' \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s'').
\end{aligned}$$

(Everything sums up to 1 and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X} \setminus \mathcal{X}', s') = \emptyset$ thus at least one successor state is in \mathcal{X}')

Thus, from the case assumption we derived a contradiction because $s' \in \arg \max_{s \in \mathcal{X} \setminus \mathcal{X}'} V_{\mathcal{R}}(s)$.

Thus, every case leads to a contradiction which concludes the proof. \square

Using Lem. 49, we now prove that for every EC \mathcal{X} , we have that for an upper bound U , the best exit value of the EC is also an upper bound. This is the crucial ingredient for showing that deflation cannot decrease a valuation below the value $V_{\mathcal{R}}$.

Lemma 50. *Let $\mathcal{X} \subseteq S \setminus (T \cup W_{\mathcal{J}})$ be an EC, and $U \in [0, 1]^{|S|}$ be a valuation with $U \geq V_{\mathcal{R}}$. Then, for all states $s \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$, we have $\text{exitVal}_U(\mathcal{X}, s) \geq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$.*

Proof. Let $s \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$. Our goal is to show that $\text{exitVal}_U(\mathcal{X}, s) \geq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$. Since for the estimation of $\text{exitVal}_U(\mathcal{X}, s)$ and $\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$ the sets of strategies, $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s)$, $\text{Trap}_U(\mathcal{X}, s)$, $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)$, and $\text{Hazard}_U(\mathcal{X}, s)$, which can be empty or non-empty, we have to consider all possible combinations. We consider the following four main cases, which we further analyse with respect to the U sets of strategies where necessary:

- Case (I)** $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$,
- Case (II)** $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$,
- Case (III)** $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$, and
- Case (IV)** $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$.

For each case we show that either the case is impossible or that the statement of the lemma holds.

Case (I) $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$.

Let $\sigma' \in \mathcal{S}(s)$ be an optimal Player \mathcal{S} strategy under $V_{\mathcal{R}}$. Since $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$, at least one of the two conditions posed by the definition of trapping strategies (Def. 24) must be violated under $V_{\mathcal{R}}$.

Since $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$, Condition (ii) of Def. 24, i.e. $\forall \rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) : (s, \rho, \sigma') \text{ staysIn } \mathcal{X}$, is trivially satisfied. Consequently, Condition (i) must be violated, i.e., it must hold that there exists no optimal strategy for Player \mathcal{S} which cannot be true. Thus, this case is impossible.

Case (II) $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$.

At state s it holds that $V_{\mathcal{R}}(s) = \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$, since we chose it to be in $\text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$. Thus, it is in the set \mathcal{X} , constructed in Condition (i) of Lem. 49. By Condition (ii) of Lem. 49, we know that all states in $\mathcal{X} \setminus \{s\}$ attain a value that is either smaller or equal $V_{\mathcal{R}}(s)$. Using this and the fact that no staying strategy can be optimal (as the set of trapping strategies is empty), we can derive a contradiction as follows.

$$\begin{aligned}
V_{\mathcal{R}}(s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (V_{\mathcal{R}} \text{ is fixpoint of } \mathfrak{B}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \text{ are optimal under } V_{\mathcal{R}}) \\
&< \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}_L(\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s), \mathcal{X}, s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}_L(\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s), \mathcal{X}, s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \mathcal{X}} \underbrace{V_{\mathcal{R}}(s') \cdot \delta(s, \rho, \sigma) \cdot \rho(a) \cdot \sigma(b)}_{\leq V_{\mathcal{R}}(s)} \\
&&& (\text{Unfolding definition of } \mathfrak{B}) \\
&\leq \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}_L(\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s), \mathcal{X}, s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \mathcal{X}} V_{\mathcal{R}}(s) \cdot \delta(s, \rho, \sigma) \cdot \rho(a) \cdot \sigma(b) \\
&&& (V_{\mathcal{R}}(s') \leq V_{\mathcal{R}}(s) \text{ by (ii) in Lem. 49}) \\
&= \mathfrak{B}(V_{\mathcal{R}})(s) \\
&= V_{\mathcal{R}}(s). && (V_{\mathcal{R}} \text{ is fixpoint of } \mathfrak{B})
\end{aligned}$$

Thus, overall $V_{\mathcal{R}}(s) < V_{\mathcal{R}}(s)$, a contradiction.

Case (III): $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$.

Since $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset$ we know by Lem. 14 that there exists $\rho_L \in \mathcal{R}_L(\mathcal{X}, s) \cdot \mathcal{R}(s) \preceq_{V_{\mathcal{R}}, \mathcal{S}(s)} \{\rho_L\}$. Our goal is to show that there exists $\rho^* \in \mathcal{R}(s)$ that is optimal under $V_{\mathcal{R}}$ and $\rho^* \in \text{Defl}_U(\mathcal{X}, s)$. We distinguish two cases:

- (III.a): For all $\rho_M \in \mathcal{R}_L(\mathcal{X}, s)$ that are optimal under $V_{\mathcal{R}}$, there exists $\rho_H \in \text{Hazard}_U(\mathcal{X}, s)$ such that $\text{Supp}(\rho_M) \cap \text{Supp}(\rho_H) \neq \emptyset$, and
- (III.b): There exists $\rho_L \in \mathcal{R}_L(\mathcal{X}, s)$ that is optimal under $V_{\mathcal{R}}$ and it holds that $\text{Supp}(\rho_L) \cap \bigcup_{\rho \in \text{Hazard}_U(\mathcal{X}, s)} \text{Supp}(\rho) = \emptyset$.

Case (III.a) In this case, for all $\rho_M \in \mathcal{R}_L(\mathcal{X}, s)$ that are optimal under $V_{\mathcal{R}}$, there exists $\rho_H \in \text{Hazard}_U(\mathcal{X}, s)$ such that $\text{Supp}(\rho_M) \cap \text{Supp}(\rho_H) \neq \emptyset$.

Let $\rho_H \in \text{Hazard}_U(\mathcal{X}, s)$ and $\rho_L \in \{\rho \in \mathcal{R}_L(\mathcal{X}, s) \mid \text{Supp}(\rho_H) \cap \text{Supp}(\rho_L) = \emptyset\}$ such that $\text{Supp}(\rho_M) \cap \text{Supp}(\rho_H) \neq \emptyset$ and $\text{Supp}(\rho_M) \cap \text{Supp}(\rho_L) \neq \emptyset$. In other words, only a strategy that mixes the supports of ρ_H and ρ_L is optimal under $V_{\mathcal{R}}$.

Due to the case assumption that only strategies that mix with some hazardous strategy are optimal, ρ_L must be sub-optimal under $V_{\mathcal{R}}$, as well as ρ_H .

More formally, for all $\rho_M \in \mathcal{R}_L(\mathcal{X}, s)$, such that $\mathcal{R}(s) \setminus \{\rho_M\} \preceq \{\rho_M\}$ it holds that

$$\begin{aligned} & \exists \rho_H \in \text{Hazard}_U(\mathcal{X}, s). \exists \rho_L \in \{\rho \in \mathcal{R}_L(\mathcal{X}, s) \mid \text{Supp}(\rho_H) \cap \text{Supp}(\rho_L) = \emptyset\} : \\ & \sup_{\text{Dist}(\text{Supp}(\rho_M) \setminus \text{Supp}(\rho_L))} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) < \{\rho_M\} \\ & \text{and} \\ & \sup_{\text{Dist}(\text{Supp}(\rho_M) \setminus \text{Supp}(\rho_H))} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) < \{\rho_M\}. \end{aligned}$$

Let $\mathcal{S}^* \subseteq \mathcal{S}(s)$ be the optimal Player \mathcal{S} strategies with respect to ρ_M . Since ρ_M has to mix ρ_L and ρ_H , the following must be true:

- there exists $\sigma_1 \in \mathcal{S}^*$ such that

$$\mathfrak{B}(V_{\mathcal{R}})(s, \rho_L, \sigma_1) < \mathfrak{B}(V_{\mathcal{R}})(s, \rho_H, \sigma_1), \text{ and} \quad (6)$$

- there exists $\sigma_2 \in \mathcal{S}^*$ such that

$$\mathfrak{B}(V_{\mathcal{R}})(s, \rho_L, \sigma_2) > \mathfrak{B}(V_{\mathcal{R}})(s, \rho_H, \sigma_2), \quad (7)$$

because otherwise either ρ_L or ρ_H would be optimal under $V_{\mathcal{R}}$ thus mixing would not be necessary. Further, Player \mathcal{S} also needs to mix the two strategies σ_1 and σ_2 to ensure optimality.

Our goal now is to show that $s \notin \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$ which would lead to a contradiction to the assumption that $s \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$, showing that at a best exit there cannot exist only optimal strategies that fulfill the same properties as ρ_M .

To estimate the best exit value from \mathcal{X} under $V_{\mathcal{R}}$, all exists from $s \in \mathcal{X}$ have to be estimated. By Eq. (6) and Eq. (7), we know that the part of the EC \mathcal{X} that is reachable with (ρ_H, σ_1) , say $\mathcal{X}' \subset \mathcal{X}$, attains a higher value than the one that is reachable with (ρ_H, σ_2) .

More formally, the following holds. Let $\hat{s} \in \max_{s' \in \mathcal{X}'} \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s')$. At \hat{s} leaving \mathcal{X}' has to be possible, because otherwise, the values at all $s' \in \mathcal{X}'$ would be equal to 0, which would be a contradiction to the assumption that $\mathcal{X} \subseteq \mathcal{S} \setminus (\mathcal{T} \cup \mathcal{W}_{\mathcal{S}})$. Then, the following chain of equations holds.

$$\begin{aligned} \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset) \\ &= \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho_M, \sigma) && (\rho_M \text{ is optimal under } V_{\mathcal{R}}) \\ &= \mathfrak{B}(V_{\mathcal{R}})(s, \rho_M, \sigma_1) && (\sigma_1 \text{ is optimal under } V_{\mathcal{R}} \text{ with respect to } \rho_M) \\ &< \mathfrak{B}(V_{\mathcal{R}})(s, \rho_H, \sigma_1). && (\text{Under } \sigma_1, \rho_H \text{ is optimal by Eq. (6)}) \\ &= \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \text{Post}(s, \rho_H, \sigma_1)} V_{\mathcal{R}}(s') \cdot \delta(s, a, b) \cdot \rho(a) \cdot \sigma(b) && (\text{Unfolding the def. of } \mathfrak{B}) \\ &\leq \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \text{Post}(s, \rho_H, \sigma_1)} \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, \hat{s}) \cdot \delta(s, a, b) \cdot \rho(a) \cdot \sigma(b) \\ & && (\text{At } \hat{s} \text{ one can attain the highest value upon leaving}) \\ &= \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, \hat{s}). && (\text{Everything sums-up to 1}) \end{aligned}$$

Thus, under $V_{\mathcal{R}}$, the state \hat{s} attains a higher value upon leaving than s . This is a contradiction to the assumption that $s \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$.

Case (III.b) There exists $\rho_L \in \mathcal{R}_L(\mathcal{X}, s)$ that is optimal under $V_{\mathcal{R}}$ and it holds that $\text{Supp}(\rho_L) \cap \bigcup_{\rho \in \text{Hazard}_U(\mathcal{X}, s)} \text{Supp}(\rho) = \emptyset$.

Then, ρ_L must belong to the set of deflating strategies under U , as it satisfies both conditions posed by the definition of deflating strategies (see Def. 25).

If $\text{Trap}_U(\mathcal{X}, s) \neq \emptyset$, then the following chain of equations holds.

$$\begin{aligned}
\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \sigma(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{By Def. 16 in case } \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset) \\
&= \inf_{\sigma \in \sigma(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho_L, \sigma) && (\rho_L \text{ is optimal}) \\
&\leq \sup_{\rho \in \text{Defl}_U(\mathcal{X}, s)} \inf_{\sigma \in \sigma(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\{\rho_L\} \subseteq \text{Defl}_U(\mathcal{X}, s)) \\
&\leq \sup_{\rho \in \text{Defl}_U(\mathcal{X}, s)} \inf_{\sigma \in \text{Trap}_U(\mathcal{X}, s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Trap}_U(\mathcal{X}, s) \subseteq \mathcal{S}(s)) \\
&\leq \sup_{\rho \in \text{Defl}_U(\mathcal{X}, s)} \inf_{\sigma \in \text{Trap}_U(\mathcal{X}, s)} \mathfrak{B}(U)(s, \rho, \sigma) && (\mathfrak{B} \text{ is order-preserving}) \\
&= \text{exitVal}_U(\mathcal{X}, s).
\end{aligned}$$

If $\text{Trap}_U(\mathcal{X}, s) = \emptyset$, then the following chain of equations holds.

$$\begin{aligned}
\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \sigma(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{By Def. 16 in case } \text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \emptyset) \\
&= \inf_{\sigma \in \sigma(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho_L, \sigma) && (\rho_L \text{ is optimal}) \\
&\leq \sup_{\rho \in \text{Defl}_U(\mathcal{X}, s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(U)(s, \rho, \sigma) && (\mathfrak{B} \text{ is order-preserving}) \\
&= \text{exitVal}_U(\mathcal{X}, s).
\end{aligned}$$

Case (IV): $\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$ and $\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \neq \emptyset$.

It holds that $s \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$, thus there exist no other state at which leaving \mathcal{X} attains better value. By Lem. 49 we have that for all $s' \in \text{bestExits}_{V_{\mathcal{R}}}(\mathcal{X})$ it holds that $V_{\mathcal{R}}(s') \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s')$. Consequently, at a best exit it can only hold that $V_{\mathcal{R}}(s') = \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s')$ because otherwise $V_{\mathcal{R}}$ would not be a fixpoint of \mathfrak{B} . Therefore, the following chain of equations holds at s .

$$\begin{aligned}
V_{\mathcal{R}}(s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (V_{\mathcal{R}} \text{ is fixpoint of } \mathfrak{B}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Hazard}_{V_{\mathcal{R}}}(\mathcal{X}, s) \text{ are optimal under } V_{\mathcal{R}}) \\
&= \sup_{\rho \in \text{Hazard}_{V_{\mathcal{R}}}(s)} \inf_{\sigma \in \text{Trap}_{V_{\mathcal{R}}}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma) && (\text{Trap}_{V_{\mathcal{R}}}(\mathcal{X}, s) \text{ are optimal under } V_{\mathcal{R}}) \\
&= \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s).
\end{aligned}$$

(By Lem. 49 we have $V_{\mathcal{R}}(s) \leq \text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s)$ however only “=” can hold as $V_{\mathcal{R}}$ is fixpoint of \mathfrak{B})

Thus, this case is equal to Case (III) where $\text{exitVal}_{V_{\mathcal{R}}}(\mathcal{X}, s) = \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(V_{\mathcal{R}})(s, \rho, \sigma)$, as the assumptions of cases (I) and (II) are not possible. □

Lemma 51 (Existence of maximal BECs). *Given a CSG G and let $C \subseteq (S \setminus (T \cup W_{\mathcal{J}}))$ be an EC and let U be a valuation. If under U there exists a BEC $\mathcal{X} \subseteq C$, then there also exists a maximal BEC, i.e. there exists a BEC $\mathcal{X}^{\max} \subseteq C$ such that $\mathcal{X} \subseteq \mathcal{X}^{\max}$ and it holds that $\mathcal{X}^{\max} \cup \{s\}$ is not a BEC for all $s \in C \setminus \mathcal{X}^{\max}$.*

Proof. To prove the lemma it suffices to show the following: Let $\mathcal{X}_1 \subseteq (S \setminus (T \cup W_{\mathcal{J}}))$ and $\mathcal{X}_2 \subseteq (S \setminus (T \cup W_{\mathcal{J}}))$ be two ECs that are BEC under U with $\mathcal{X}_1 \cap \mathcal{X}_2 \neq \emptyset$. Then, it holds that $\mathcal{X}_1 \cup \mathcal{X}_2$ is also a BEC.

Let $s^\cap \in (\mathcal{X}_1 \cap \mathcal{X}_2)$, $s' \in \mathcal{X}_1$, and $s'' \in \mathcal{X}_2$. Since \mathcal{X}_1 is a BEC at each state $s \in \mathcal{X}_1$ it holds that $\text{Hazard}_U(\mathcal{X}_1, s) \neq \emptyset$ and $\text{Trap}_U(\mathcal{X}_1, s) \neq \emptyset$ for all $s \in \mathcal{X}_1$. Further, since \mathcal{X}_1 is an EC, there exists a play $s_0 s_1 \dots$ such that $s_0 = s'$ and $s_n = s^\cap$ for some n and for all $0 \leq i < n$ it holds that $s_{i+1} \in \text{Post}(s_i, \rho_i^*, \sigma_i^*)$ where $\rho_i^* \in \text{Hazard}_U(\mathcal{X}_1, s_i)$ and $\sigma_i^* \in \text{Trap}_U(\mathcal{X}_1, s_i)$.

Similarly, since \mathcal{X}_2 is also a BEC, there exists a play $s'_0 s'_1 \dots$ such that $s'_0 = s^\cap$ and $s'_m = s''$ for some m and for all $0 \leq j < m$ it holds that $s'_{j+1} \in \text{Post}(s'_j, \rho'_j, \sigma'_j)$ where $\rho'_j \in \text{Hazard}_U(\mathcal{X}_2, s'_j)$ and $\sigma'_j \in \text{Trap}_U(\mathcal{X}_2, s'_j)$.

Therefore, for any $s, s' \in \mathcal{X}_1 \cup \mathcal{X}_2$ there exist a play $s''_0 s''_1 \dots$ such that $s''_0 = s'$ and $s''_l = s''$ for some l and for all $0 \leq k < l$ it holds that $s''_{k+1} \in \text{Post}(s''_k, \rho''_k, \sigma''_k)$ where $\rho''_k \in \text{Hazard}_U(\mathcal{X}_1 \cup \mathcal{X}_2, s''_k)$ and $\sigma''_k \in \text{Trap}_U(\mathcal{X}_1 \cup \mathcal{X}_2, s''_k)$.

Thus, we have shown that the $\mathcal{X}_1 \cup \mathcal{X}_2$ is an EC. Further, since at all states of \mathcal{X}_1 and \mathcal{X}_2 the two conditions (i) and (ii) of Def. 18 are fulfilled, also $\mathcal{X}_1 \cup \mathcal{X}_2$ fulfills them. Consequently, $\mathcal{X}_1 \cup \mathcal{X}_2$ is a BEC which in turn proves that there exist maximal BECs. □

Now that we have proven that maximal **BECs** indeed exist, the next step is to prove the correctness of **FIND_MBECs**, i.e. the algorithm that can find maximal **BECs**.

Lemma 30 (**FIND_MBECs** is correct— Proof in App. C-C). *For a **CSG**, a **MEC** C and a valid upper bound U , it holds that $\mathcal{X} \in \text{FIND_MBECs}(G, C, U)$ if and only if \mathcal{X} is a **BEC** in C and there exists no $T \subseteq C$ that is a **BEC** and $\mathcal{X} \subsetneq T$.*

Proof. We prove the lemma in two steps. First, if for a set of states $\mathcal{X} \subseteq C$ it holds that $\mathcal{X} \in \text{FIND_MBECs}(G, C, U)$ then \mathcal{X} is a maximal **BEC**. Second, we will show that if a set of states $\mathcal{X} \subseteq C$ is a maximal **BEC**, then $\mathcal{X} \in \text{FIND_MBECs}(G, C, U)$.

1. Direction “ \Rightarrow ” Let $\mathcal{X} \in \text{FIND_MBEC}(G, C, U)$. We need to show that \mathcal{X} is a maximal **BEC**.

Let $B := \{s \in C \mid \text{Hazard}_U(C, s) \neq \emptyset\}$. Since, $\mathcal{X} \in \text{FIND_MBEC}(G, C, U)$, then $\mathcal{X} \subseteq B$.

We show that \mathcal{X} is a maximal **BEC** in C via a contradiction. In particular, we consider the following two cases:

Case (i) \mathcal{X} is maximal **EC** in B but not **BEC** in C ; and

Case (ii) \mathcal{X} is a **BEC** in C but not a maximal one.

Case (i) Assume towards a contradiction that \mathcal{X} is a maximal **EC** in B but not **BEC** in C . Then there exists $s' \in \mathcal{X}$ such that $\text{Hazard}_U(\mathcal{X}, s') = \emptyset$. Thus, every strategy $\rho' \in \mathcal{R}(s')$ must violate at least one of the three conditions of Def. 16. We write out the negations.

(a) $\rho' \notin \mathcal{R}_L(\mathcal{X}, s')$, i.e. the strategy is leaving $\rho' \in \mathcal{R}_L(\mathcal{X}, s')$.

(b) We do not have $\mathcal{R}(s') \setminus \{\rho'\} \preceq_{S(s')} \{\rho'\}$. By contraposition of Lem. 14, this implies $\{\rho'\} \prec_{S(s')} \mathcal{R}(s') \setminus \{\rho'\}$, so the strategy is sub-optimal.

(c) We do not have $\mathcal{R}_L(\mathcal{X}, s') \prec_{S(s')} \{\rho'\}$. By Lem. 14, this implies $\{\rho'\} \preceq_{S(s')} \mathcal{R}_L(\mathcal{X}, s')$, i.e. there are leaving strategies that are not worse than ρ' . Note that in particular, this implies that $\mathcal{R}_L(\mathcal{X}, s')$ is non-empty.

Our assumption gives us the disjunction over the three violated conditions. We proceed by a case distinction, always assuming that each strategy violates a certain condition, which allows us to prove our goal, or, if there exists a strategy satisfying the condition, we continue with the next one.

Case “all strategies satisfy (a)”. In this case all strategies of player \mathcal{R} must be leaving \mathcal{X} , i.e. $\mathcal{R}_L(\mathcal{X}, s') = \mathcal{R}(s')$. Thus, for all $\sigma \in S(s')$ and for all $\rho \in \mathcal{R}(s')$ there exists $\hat{s} \in \text{Post}(s', \rho, \sigma)$ such that $\hat{s} \notin \mathcal{X}$. Consequently, \mathcal{X} is not an **EC** which is a contradiction to the assumption that \mathcal{X} is an **EC**.

Case “some strategy violates (a) and satisfies (c)”. We assume that ρ' violates (a) but satisfies condition (c). Thus, Player \mathcal{R} has a non-leaving strategy, however it is not optimal. In other words, Player \mathcal{R} has an optimal strategy $\rho \in \mathcal{R}(s')$ such that there exists $\hat{s} \in \text{Post}(s', \rho, \sigma)$ such that $\hat{s} \notin \mathcal{X}$. Then, using the same argumentation as in the previous case we can derive a contradiction to the assumption that $\mathcal{X} \in \text{FIND_MBEC}(G, C, U)$.

Case “some strategy violate (a) and (c) but satisfy (b)”. We assume that ρ' violates (a) and (c) but satisfies condition (b). Then, strategy ρ' is dominated by another non-leaving strategy, say $\rho \in \mathcal{R}_L(\mathcal{X}, s')$, that in turn has to violate one of the two previous cases.

Case (ii) Now, assume towards a contradiction that \mathcal{X} is a **BEC** in C but not maximal. Since $\mathcal{X} \in \text{FIND_MECs}(G, B)$, \mathcal{X} is maximal **BEC** in B . Then, due to the case assumption that \mathcal{X} is not maximal, there has to exist a set of states $\hat{S} \subseteq C \setminus B$ such that $\mathcal{X} \cup \hat{S}$ is a maximal **BEC** in C . However, since for all $s \in C \setminus B$ it holds that $\text{Hazard}_U(C, s) = \emptyset$, $\mathcal{X} \cup \hat{S}$ cannot be an **EC** in C which is a contradiction to the assumption that \mathcal{X} is an **EC** in C .

Overall, we have shown that \mathcal{X} is a **BEC** and it is also maximal.

2. Direction “ \Leftarrow ” Let $\mathcal{X} \subseteq C$ be a maximal **BEC** in the **MEC** C . We need to show that $\mathcal{X} \in \text{FIND_MBEC}(G, C, U)$ holds. Let $B := \{s \in C \mid \text{Hazard}_U(C, s) \neq \emptyset\}$. Assume towards a contradiction $\mathcal{X} \notin \text{FIND_MBEC}(G, C, U)$. Then, we distinguish two cases: (i) $\exists s \in \mathcal{X}$ such that $s \notin B$, or (ii) $\mathcal{X} \notin \text{FIND_MECs}(G, B)$.

Case (i) If $\exists s \in \mathcal{X}$ such that $s \notin B$, then it holds that $\text{Hazard}_U(C, s) = \emptyset$, by Claim 6 we also have $\text{Hazard}_U(\mathcal{X}, s) = \emptyset$ as $\mathcal{X} \subseteq C$, a contradiction.

Case (ii) In this case it holds that $\forall s \in \mathcal{X}$ we have that $s \in B$, however, $\mathcal{X} \notin \text{FIND_MECs}(G, B)$. We can again distinguish three cases: (ii.a) $\exists \mathcal{X}' \in \text{FIND_MECs}(G, B)$ such that $\mathcal{X}' \subsetneq \mathcal{X}$; and (ii.b) $\exists \mathcal{X}' \in \text{FIND_MECs}(G, B)$ such that $\mathcal{X}' \cap \mathcal{X} \neq \emptyset$ and $\exists s' \in \mathcal{X}' \setminus \mathcal{X}$; and (ii.c) $\forall \mathcal{X}' \in \text{FIND_MECs}(G, B)$ it holds that $\mathcal{X}' \cap \mathcal{X} = \emptyset$.

Case (ii.a) Then there is a state in $\mathcal{X}' \setminus B$, a contradiction to the assumption that \mathcal{X} is a maximal **BEC** in B .

Case (ii.b) By “1. Direction “ \Rightarrow ”” we know that \mathcal{X}' is a maximal **BEC**. In particular, for all $s \in \mathcal{X}'$ it holds that $\text{Hazard}_U(\mathcal{X}', s) \neq \emptyset$. Then, by Lem. 51 we also know that since $\mathcal{X}' \cap \mathcal{X} \neq \emptyset$, then $\mathcal{X}' \cup \mathcal{X}$ has to be also a **BEC**. However, this is a contradiction to the assumption that \mathcal{X} is already a maximal **BEC**.

Case (ii.c) Since $\forall \mathcal{X}' \in \text{FIND_MECs}(G, B)$ it holds that $\mathcal{X}' \cap \mathcal{X} = \emptyset$ then, for all $s' \in \mathcal{X}$ we have $s' \notin B$. However, \mathcal{X} is a maximal **BEC** in the **MEC** C by assumption and therefore for all $s \in \mathcal{X}$ it holds that $\text{Hazard}_U(\mathcal{X}, s) \neq \emptyset$. Further, since $\mathcal{X} \subseteq C$, by Claim 6 it also holds that $\text{Hazard}_U(C, s) \neq \emptyset$. Consequently, such s' cannot exist and therefore such \mathcal{X}' cannot exist, a contradiction.

Claim 6: For any ECs $C \subseteq C'$ we have $\text{Hazard}_U(C, s) \subseteq \text{Hazard}_U(C', s)$.

Proof. Let $\rho \in \mathcal{R}_L(C, s)$. We need to show that $\rho \in \mathcal{R}_L(C', s)$ holds.

First, strategy ρ satisfies condition (i) of Def. 16, thus, there has to exist $\hat{\sigma} \in \mathcal{S}(s)$, such that $(s, \rho, \hat{\sigma})$ **staysIn** C holds. Since $C \subseteq C'$, it also has to be true that $(s, \rho, \hat{\sigma})$ **staysIn** C' . Consequently, ρ satisfies condition (i) of Def. 16 with respect to C' .

Second, ρ satisfies condition (ii) of Def. 16, i.e. ρ is optimal. As the optimality of a strategy is independent of any set of states this condition is satisfied anyway.

Third, ρ satisfies condition (iii) of Def. 16, i.e. all strategies that are leaving with respect to C are sub-optimal. More precisely, it holds that $\mathcal{R}_L(C, s) \prec_{\mathcal{S}(s)} \{\rho\}$. Since $C \subseteq C'$ holds, there cannot exist an optimal strategy that is leaving with respect to the greater set of states C' . Consequently, ρ also satisfies condition (iii) of Def. 16 with respect to C' . Consequently, ρ satisfies all conditions posed by Def. 16 with respect to C' , so it holds that $\text{Hazard}_U(C, s) \subseteq \text{Hazard}_U(C', s)$. \square

▲
 \square

D. Soundness and Completeness

In order to prove the soundness and completeness of Alg. 1 with Alg. 3 as DEFLATE routine, we need to prove that the sequence of lower and upper bounds converges to a fixpoint. Therefore, before we can prove Thm. 37, first we need to prove the following lemma.

Lemma 52 (Alg. 1 converges to a fixpoint). *The BVI algorithm (Alg. 1) converges to a fixpoint, i.e. $\lim_{k \rightarrow \infty} (\mathfrak{B}^k(L^0), (\mathfrak{D} \circ \mathfrak{B})^k(U^0)) = (\mathfrak{B}(\lim_{k \rightarrow \infty} \mathfrak{B}^k(L^0)), (\mathfrak{D} \circ \mathfrak{B})(\lim_{k \rightarrow \infty} (\mathfrak{D} \circ \mathfrak{B})^k(U^0)))$.*

Proof. We consider the domain $\mathbb{V} := [0, 1]^{|S|} \times [0, 1]^{|S|}$, i.e. every element consists of two vectors of real numbers, the under- and over-approximation. The bottom element of the domain, denoted by \perp , is $(\vec{0}, \vec{1})$, where for $a \in [0, 1]$, \vec{a} denotes the function that assigns a to all states. We further restrict the domain to exclude elements of the domain that are trivially irrelevant for the computation. In particular, we exclude all tuples (L, U) where $L(s) < 1$ for a target state $s \in T$ or $U(s) > 0$ for a state with no path to the target state $s \in W_{\mathcal{S}}$. Then the bottom element is $\perp = (L^0, U^0)$, i.e. the vector that we have before the first iteration of the main loop of Alg. 1. Concretely, $L^0(s)$ is 1 for all $s \in T$, i.e. target states, and 0 everywhere else, and $U^0(s)$ is 0 for all $s \in W_{\mathcal{S}}$, i.e. states where \mathcal{S} can surely win, and 1 everywhere else.

We define a comparator \sqsubseteq on \mathbb{V} , to compare two elements of the domain. We write $(L^k, U^k) \sqsubseteq (L^{k+1}, U^{k+1})$ if and only if both $L^k \leq L^{k+1}$ and $U^k \geq U^{k+1}$ hold with component-wise comparison. Intuitively, $(L^k, U^k) \sqsubseteq (L^{k+1}, U^{k+1})$ holds if (L^{k+1}, U^{k+1}) is a more precise approximation than (L^k, U^k) . The comparator \sqsubseteq induces a complete partial order over the domain, since we have a bottom element and every direct subset has a supremum; the latter claim holds, because \sqsubseteq reduces to component-wise comparison between real numbers from $[0, 1]$, where suprema exist. For more details on the definition of directed set and complete partial orders, we refer to [13].

Alg. 1 first applies the Bellman operator on the over- and under-approximation and subsequently applies the deflate operator on the over-approximation (i.e. upper bound). Thus, the operator that mimics the behavior of the algorithm is the following $\text{BVI}(L^k, U^k) = (\mathfrak{B}(L^k), (\mathfrak{D} \circ \mathfrak{B})(U^k))$.

From Thm. 5 we know that the under-approximation converges. Also, by Lem. 36 we know that $(\mathfrak{D} \circ \mathfrak{B})$ is order-preserving. Thus, for the final argument it remains to show that \mathfrak{D} is (Scott-)continuous.

A map is (Scott-)continuous if, for every directed set D in $(\text{proj}_2(\mathbb{V}), \sqsubseteq)$, the subset $\mathfrak{D}(D)$ of $(\text{proj}_2(\mathbb{V}), \sqsubseteq)$ is directed, and $\mathfrak{D}(\sup D) = \sup \mathfrak{D}(D)$. Let $s \in S$, if s , under $\sup D$, does not belong to any BEC, then $\mathfrak{D}(\sup D)(s) = \sup D(s)$. Thus, we proceed with the assumption that under the valuation $\sup D$, $s \in \mathcal{X} \subseteq S$ such that \mathcal{X} is a BEC.

Let $\text{Hazard}_{\sup D}(\mathcal{X}, s)$ be the set of hazardous strategies and let $\text{Trap}_{\sup D}(\mathcal{X}, s)$ be the set of suitable counter strategies for player \mathcal{S} (see Def. 18). Since \mathcal{X} is a BEC both sets are non-empty, i.e. $\text{Hazard}_{\sup D}(\mathcal{X}, s) \neq \emptyset$ and $\text{Trap}_{\sup D}(\mathcal{X}, s) \neq \emptyset$, at all $s \in \mathcal{X}$. Thus, $\mathfrak{D}(\sup D)(s) = \min(\sup D(s), \text{bestExitVal}_{\sup D}(\mathcal{X}))$. For the sake of readability let $\mathcal{R}'(s) := \text{Dist}(\Gamma_{\mathcal{S}}(s) \setminus \bigcup_{\rho'' \in \text{Hazard}_{\sup D}(\mathcal{X}, s)} \text{Supp}(\rho''))$ and $\mathcal{S}'(s) := \text{Trap}_{\sup D}(\mathcal{X}, s)$ for some $s \in \mathcal{X}$. Then, the following chain of equations holds for $s \in \mathcal{X}$.

$$\begin{aligned}
\mathfrak{D}(\sup D)(s) &= \min(\sup D(s), \text{bestExitVal}_{\sup D}(\mathcal{X})) \\
&= \min(\sup D(s), \max_{s' \in \mathcal{X}} \text{exitVal}_{\sup D}(\mathcal{X}, s')) && \text{(By Def. 28 (best exit value))} \\
&= \min(\sup D(s), \max_{s' \in \mathcal{X}} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(\sup d)(s', \rho, \sigma)) && \text{(By Def. 27 (exit value))} \\
&= \min(\sup D(s), \max_{s' \in \mathcal{X}} \sup_{d \in D} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma)) && \text{(Bellman operator is Scott-continuous (see proof Theorem 9))}
\end{aligned}$$

$$\begin{aligned}
&= \min(\sup D(s), \sup_{d \in D} \max_{s' \in \mathcal{X}} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma)) && \text{(By Claim 7)} \\
&= \min(\sup D(s), \sup_{d \in D} \max_{s' \in \mathcal{X}} \text{exitVal}_d(\mathcal{X}, s')) && \text{(By Def. 27 (exit value))} \\
&= \min(\sup D(s), \sup_{d \in D} \text{bestExitVal}_d(\mathcal{X})) && \text{(By Def. 28 (best exit value))} \\
&= \sup_{d \in D} \min(d(s), \text{bestExitVal}_d(\mathcal{X})) && \text{(min is Scott-continuous)} \\
&= \sup_{d \in D} \mathfrak{D}(d)(s). && \text{(Since } \mathfrak{D} \text{ is the deflate operator)}
\end{aligned}$$

Claim 7: Since $\max_{s' \in \mathcal{X}} \text{exitVal}_{\sup D}(\mathcal{X}, s')$ exists the following chain of equations holds:

$$\begin{aligned}
&\sup_{d \in D} \max_{s' \in \mathcal{X}} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma) \\
&= \sup_{d \in D} \sup_{s' \in \mathcal{X}} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma) && \text{(As the maximum exists)} \\
&= \sup_{(d, s') \in (D \times \mathcal{X})} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma) \\
&= \sup_{s' \in \mathcal{X}} \sup_{d \in D} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma) \\
&= \max_{s' \in \mathcal{X}} \sup_{d \in D} \sup_{\rho \in \mathcal{R}'(s')} \inf_{\sigma \in \mathcal{S}'(s')} \mathfrak{B}(d)(s', \rho, \sigma).
\end{aligned}$$

▲

We have shown that \mathfrak{B} and \mathfrak{D} are both Scott-continuous, thus the composition $(\mathfrak{D} \circ \mathfrak{B})$, i.e. the sequential application of the operators \mathfrak{B} and \mathfrak{D} , is also continuous [44]. Further, \mathbb{V} is a complete partial order. Therefore, Kleen's Fixpoint Theorem [13, Theorem 8.15] is applicable.

Then, we know that

$$\lim_{k \rightarrow \infty} (\mathfrak{D} \circ \mathfrak{B})^k(\perp) = \sup_{k \geq 0} (\mathfrak{D} \circ \mathfrak{B})^k(\perp) \quad (\star)$$

holds and by the theorem we know that the fixpoint exists and is given by $\sup_{k \geq 0} (\mathfrak{D} \circ \mathfrak{B})^k(\perp)$. Now, we can finally conclude:

$$\begin{aligned}
&(\mathfrak{D} \circ \mathfrak{B})(\lim_{k \rightarrow \infty} (\mathfrak{D} \circ \mathfrak{B})^k(\perp)) \stackrel{(\star)}{=} (\mathfrak{D} \circ \mathfrak{B})(\sup_{k \geq 0} (\mathfrak{D} \circ \mathfrak{B})^k(\perp)) \\
&= \sup_{k \geq 0} (\mathfrak{D} \circ \mathfrak{B})((\mathfrak{D} \circ \mathfrak{B})^k(\perp)) && \text{(since } (\mathfrak{D} \circ \mathfrak{B}) \text{ is continuous)} \\
&= \sup_{k \geq 1} (\mathfrak{D} \circ \mathfrak{B})^k(\perp) \\
&= \sup_{k \geq 0} (\mathfrak{D} \circ \mathfrak{B})^k(\perp) && \text{(since } \perp \sqsubseteq (\mathfrak{D} \circ \mathfrak{B})^k(\perp) \text{ for all } k) \\
&= \lim_{k \rightarrow \infty} (\mathfrak{D} \circ \mathfrak{B})^k(\perp). && \text{by } (\star)
\end{aligned}$$

□

Now we can prove the following theorem.

Theorem 37 (Soundness and completeness - Proof in App. C-D). *For CSGs Alg. 1, using Alg. 3 as DEFLATE, produces monotonic sequences \mathbf{L} under- and \mathbf{U} over-approximating $\mathbf{V}_{\mathcal{A}}$, and terminates for every $\varepsilon > 0$.*

Proof. We denote by \mathbf{L}^k and \mathbf{U}^k the lower and upper bound function after the k -th call of DEFLATE. \mathbf{L}^k and \mathbf{U}^k are monotonic under-, respectively, over-approximation of $\mathbf{V}_{\mathcal{A}}$ because they are updated via Bellman updates respectively $(\mathfrak{D} \circ \mathfrak{B})$ -updates, which are order-preserving as shown in Lem. 36 and soundness as shown in Lem. 35 .

Since DEFLATE iterates over finite sets, the computations take a finite time. Thus, it remains to prove that the main loop of Alg. 1 terminates, i.e., for all $\varepsilon > 0$, there exists an $n \in \mathbb{N}$ such that for all $s \in \mathbf{S}$ $\mathbf{U}^n(s) - \mathbf{L}^n(s) \leq \varepsilon$. It suffices to show that $\lim_{k \rightarrow \infty} \mathbf{U}^k - \mathbf{V}_{\mathcal{A}} = 0$, because $\lim_{k \rightarrow \infty} \mathbf{L}^k = \mathbf{V}_{\mathcal{A}}$ (from e.g. [42]).

In the following let $\mathbf{U}^* := \lim_{k \rightarrow \infty} \mathbf{U}^k$, that exists by Lemma 48, and $\Delta(s) := \mathbf{U}^*(s) - \mathbf{V}_{\mathcal{A}}(s)$. Assume towards a contradiction that the algorithm does not converge, i.e., there exists a state $s \in \mathbf{S}$ with $\Delta(s) > 0$.

The proof is structured as follows.

- From $\Delta > 0$ we derive that there has to exist a BEC.

- The states with the maximal Δ contain **BECs**.
- Alg. 1 will find a **BEC** contained in $\mathcal{X} \subseteq S$ and deflate it.
- Deflating will decrease the upper bound of that states contained in the **BEC**, which is a contradiction because by Lemma 52, U^* is a fixpoint.

Let $\Delta^{\max} := \max_{s \in S} \Delta(s)$ and let $C := \{s \in S \mid \Delta(s) = \Delta^{\max}\}$. If C does not contain any **BECs**, then the contraposition of Thm. 9 proves our goal. Thus, we continue with the assumption that C contains **BECs**.

Let $\mathcal{X}' \subseteq C$ be a **BEC** contained in C that Alg. 3 will eventually find and deflate. We now consider *bottom* **BECs**. A **BEC** \mathcal{X}' is called bottom in \mathcal{X} if none of the successors of a strategy that leaves the **BEC** \mathcal{X}' , is part of another **BEC** in \mathcal{X} . A bottom **BEC** can be computed by first finding the bottom **MEC** within \mathcal{X} , then identifying all **BECs**, ordering them topologically and finally picking the one at the end of a chain.

Let $\mathcal{X}' \subseteq \mathcal{X}$ be a bottom **BEC** that Algorithm 3 eventually finds. In order to deflate the **BEC**, we need to estimate the exit value for each $s \in \mathcal{X}'$, as defined in Definition 27. Further, $\text{Hazard}_{U^*}(\mathcal{X}', s)$ is the set of hazardous strategies and since \mathcal{X}' is a **BEC**, at all states $s \in \mathcal{X}'$, we have $\text{Hazard}_{U^*}(\mathcal{X}', s) \neq \emptyset$.

We distinguish two cases: (i) $\text{Trap}_{U^*}(\mathcal{X}', s) \neq \emptyset$; and (ii) $\text{Trap}_{U^*}(\mathcal{X}', s) = \emptyset$.

Case (i) In this case we have that $\text{Trap}_{U^*}(\mathcal{X}', s) \neq \emptyset$. Then $\text{Defl}_{U^*}(\mathcal{X}', s) \neq \emptyset$ (and because $\mathcal{X} \cap W_{\mathcal{S}} = \emptyset$), thus the following chain of equations holds.

$$\begin{aligned}
U^*(s) &= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(U^*)(s, \rho, \sigma) && (U^* \text{ is a fixpoint}) \\
&= \sup_{\rho \in \text{Hazard}_{U^*}(\mathcal{X}', s)} \inf_{\sigma \in \text{Trap}_{U^*}(\mathcal{X}', s)} \mathfrak{B}(U^*)(s, \rho, \sigma) && (\text{Hazard}_{U^*}(\mathcal{X}', s) \neq \emptyset) \\
&> \sup_{\rho \in \text{Defl}_{U^*}(\mathcal{X}', s)} \inf_{\sigma \in \text{Trap}_{U^*}(\mathcal{X}', s)} \mathfrak{B}(U^*)(s, \rho, \sigma) && (\text{Defl}_{U^*}(\mathcal{X}', s) \text{ are sub-optimal with respect to } U^*) \\
&= \text{exitVal}_{U^*}(\mathcal{X}', s) && (\text{By Def. 27})
\end{aligned}$$

Since $\text{bestExitVal}_{U^*}(\mathcal{X}') = \max_{s' \in \mathcal{X}'} \text{exitVal}_{U^*}(\mathcal{X}', s')$, and above we have shown that for all $s \in \mathcal{X}'$ it holds $U^*(s) > \text{exitVal}_{U^*}(\mathcal{X}', s)$, we obtain a contradiction to the assumption that U^* is a fixpoint.

Case (ii) In this case it holds that $\text{Trap}_{U^*}(\mathcal{X}', s) = \emptyset$. Thus, Player \mathcal{S} prefers strategies that are leaving with respect to $\text{Hazard}_{U^*}(\mathcal{X}', s)$. Since \mathcal{X}' is a bottom **BEC** at least one successor state does not belong to \mathcal{X} . Therefore, the following chain of equations holds.

$$\begin{aligned}
\Delta^{\max} - V_{\mathcal{R}}(s) &= U^*(s) && (\text{By definition of } \Delta^{\max}) \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \mathfrak{B}(U^*)(s, \rho, \sigma) && (U^* \text{ is a fixpoint}) \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in S} U^*(s') \cdot \delta(s, a, b)(s') \cdot \rho(a) \cdot \sigma(b) && (\text{Unfolding definition of } \mathfrak{B}) \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \mathcal{X}'} U^*(s') \cdot \delta(s, a, b)(s') \cdot \rho(a) \cdot \sigma(b) \\
&\quad + \sum_{s'' \in S \setminus \mathcal{X}'} U^*(s'') \cdot \delta(s, a, b)(s'') \cdot \rho(a) \cdot \sigma(b) && (\text{Player } \mathcal{S} \text{ prefers leaving } \mathcal{X}' \text{ and thus } \mathcal{X}') \\
&= \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \mathcal{X}'} U^*(s') \cdot \delta(s, a, b)(s') \cdot \rho(a) \cdot \sigma(b) \\
&\quad + \sum_{s'' \in S \setminus \mathcal{X}'} \left(\underbrace{\Delta(s'') + V_{\mathcal{R}}(s'')}_{< \Delta^{\max}} \right) \cdot \delta(s, a, b)(s'') \cdot \rho(a) \cdot \sigma(b) && (\text{By def. of } \Delta) \\
&< \sup_{\rho \in \mathcal{R}(s)} \inf_{\sigma \in \mathcal{S}(s)} \sum_{(a,b) \in \mathcal{A}} \sum_{s' \in \mathcal{X}'} U^*(s') \cdot \delta(s, a, b)(s') \cdot \rho(a) \cdot \sigma(b) \\
&\quad + \sum_{s'' \in S \setminus \mathcal{X}'} (\Delta^{\max} + V_{\mathcal{R}}(s'')) \cdot \delta(s, a, b)(s'') \cdot \rho(a) \cdot \sigma(b) \\
&\hspace{15em} (\Delta \text{ for states outside } \mathcal{X} \text{ is strictly smaller than } \Delta^{\max}) \\
&= U^*(s) && (\text{Everything sums-up to } 1)
\end{aligned}$$

Thus, we get a contradiction to the assumption that U^* is a fixpoint. \square

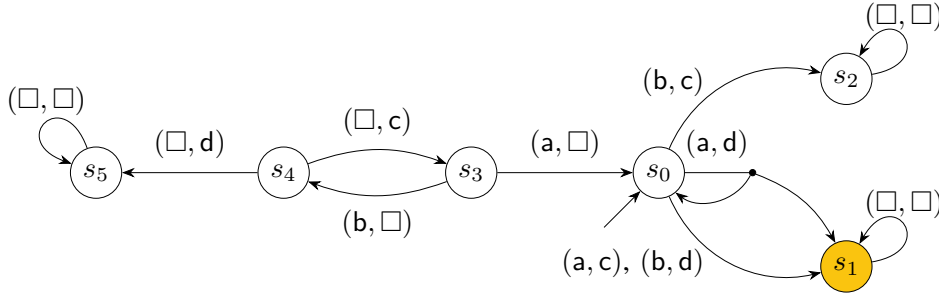


Fig. 6: Counter example for the **BVI** provided in [7].

APPENDIX D COUNTER-EXAMPLES FOR PREVIOUS WORKS

A. Mistake in [9]

In [9], the authors provide an exemplary **CSG** that illustrates the incorrectness of the **BVI** algorithm presented in [7]. Fig. 6 shows the **CSG**. The value attainable at s_5 is 0.6 while the value attainable at s_0 is $2 - \sqrt{2}$. In our algorithm (Alg. 3) the set $\{s_4, s_3\}$ does not constitute a **BEC** for the valuation $U(s_4) = U(s_3) = U(s_1) = 1$ and $U(s_2) = 0$. Consequently, the while-loop of Alg. 3, which is responsible for deflating, is not executed and thus the values of all states are updated only using the Bellman operator. This yields the correct values by Thm. 9. Therefore, Alg. 3 correctly sets the valuation of states s_4 and s_3 to the value 0.6. In contrast, the algorithm presented in [7] correctly sets the value of state s_4 to 0.6 but reduces the value of state s_3 to $2 - \sqrt{2} < 0.6$.

B. Mistake in [19]

There was an attempt to fix the above problem in [19], however, also this solution is incorrect. An exemplary **CSG** that illustrates the incorrectness of the approach presented in [19] is the **CSG** illustrated in Fig. 2. While our approach correctly deflates the value of state s_0 , the best exit as defined in [19] is falsely 1, i.e. the value of s_0 is never reduced.