

Solving Zero-Sum One-Sided Partially Observable Stochastic Games

Karel Horák¹, Branislav Bošanský¹, Vojtěch Kovařík¹, Christopher Kiekintveld²

¹ *Artificial Intelligence Center,
Department of Computer Science,
Faculty of Electrical Engineering,
Czech Technical University in Prague*
`{karel.horak, branislav.bosansky, kovarvo1}@fel.cvut.cz`

² *Computer Science Department,
University of Texas at El Paso,
cdkiekintveld@utep.edu*

Abstract

Many security and other real-world situations are dynamic in nature and can be modelled as strictly competitive (or zero-sum) dynamic games. In these domains, agents perform actions to affect the environment and receive observations – possibly imperfect – about the situation and the effects of the opponent’s actions. Moreover, there is no limitation on the total number of actions an agent can perform — that is, there is no fixed horizon. These settings can be modelled as partially observable stochastic games (POSGs). However, solving *general* POSGs is computationally intractable, so we focus on a broad subclass of POSGs called *one-sided POSGs*. In these games, only one agent has imperfect information while their opponent has full knowledge of the current situation. We provide a full picture for solving one-sided POSGs: we (1) give a theoretical analysis of one-sided POSGs and their value functions, (2) show that a variant of a value-iteration algorithm converges in this setting, (3) adapt the heuristic search value-iteration algorithm for solving one-sided POSGs, (4) describe how to use approximate value functions to derive strategies in the game, and (5) demonstrate that our algorithm can solve one-sided POSGs of non-trivial sizes and analyze the scalability of our algorithm in three different domains: pursuit-evasion, patrolling, and search games.

Keywords: zero-sum partially observable stochastic games, one-sided information, value iteration, heuristic search value iteration

1. Introduction

Non-cooperative game theory models the interaction of multiple agents in a joint environment. Rational agents perform actions in the environment to achieve their own, often conflicting goals. The interaction of agents is typically very

complex in real-world dynamic scenarios — the agents can perform sequences of multiple actions while only having partial information about the actions of others and the events in the environment.

Finding out (approximate) optimal strategies for agents in dynamic environment with imperfect information is a long-standing problem in Artificial Intelligence. Its applications range from recreational games, such as poker [32, 11], to uses in security such as patrolling [4, 50, 5] and pursuit-evasion games [25, 24, 1].

For tackling this problem, game theory can provide appropriate mathematical models and algorithms for computing (approximate) optimal strategies according to some game-theoretic solution concept. Among all existing game-theoretic models suitable for modelling dynamic interaction with imperfect information, *partially observable stochastic games (POSGs)* are one of the most general ones. POSGs model situations where all players have only partial information about the state of the environment, agents perform actions and receive observations, and the length of the interaction among agents is not a priori bounded. As such, the expressive possibilities of POSGs are broad. In particular, they can model all considered security scenarios as well as recreational games.

Despite having high expressive power, POSGs have limited applications due to the complexity of computing (approximate) optimal strategies. There are two main reasons for this. First, the imperfect information provides challenges for sequential decision-making even in the single-agent case – partially observable Markov decision processes (POMDPs). Theoretical results show that various exact and approximate problems in POMDPs are undecidable [31]. Focused research effort has yielded several approximate algorithms with convergence guarantees [39, 27] scalable even to large POMDPs [37]. The main step when solving a POMDP is to reason about *belief states* – probability distributions over possible states. Note that an agent can easily deduce a belief state in a POMDP since the environment changes only as a result of the agent’s actions or because of the environment’s stochasticity (which is known). In POSGs, however, the presence of another agent(s) changing the environment generates another level of complexity. Suppose all players have partial information about the environment. In that case, each player needs to reason not only about their belief over environment states, but also about opponents’ beliefs, their beliefs over beliefs, and so on. This issue is called the problem with *nested beliefs* [30] and cannot be avoided in general unless we pose additional assumptions on the game model. This is primarily because, in general POSGs, the choice of the optimal action (strategy) of a player depends on these nested beliefs. To avoid this issue, we will focus on a subclass of POSGs that does not suffer from the problem of nested beliefs while still being expressive enough to contain many existing real-world games and scenarios.

One such sub-class of POSGs are two-player concurrent-move games where one player is assumed to have full knowledge about the environment and only one player has partial information. In this case, the player with partial information (player 1 from now on) does not have to reconstruct the belief of the opponent (player 2) since player 2 always has full information about the true state of the environment. Similarly, player 2 can always reconstruct the belief of player 1 by

using the full information about the environment, which includes information about the action-history of player 1. The game is played over stages where both players independently choose their next action (i.e., albeit player 2 has full knowledge about the current history and state, he does not know the action player 1 is about to play in the current stage). The state of the game with the joint action of the players determines the next state and the next observation generated for player 1. We term this class of games as *one-sided POSG*. While this class of games has appeared before in the literature (e.g., in [44] as *Level-1 stochastic games*, or in [13] as *semiperfect-information stochastic games*¹) we are the first to focus on designing a practical algorithm for computing (approximately) optimal strategies.

Despite the seemingly-strong assumption on the perfect information for player 2, the studied class of one-sided POSG has broad application possibilities, especially in security. In particular, this model subsumes patrolling games [4, 50, 5] or pursuit-evasion games [25, 24, 1]. In many security-related problems, the defender is protecting an area (or a computer network) against the attacker that wants to attack it (e.g., by intruding into the area or infiltrating the network). The defender does not have full information about the environment since he does not know which actions the attacker performed (e.g., which hosts in the computer network have been compromised by the attacker). At the same time, it is difficult for the defender to exactly know what information the attacker has since the attacker can infiltrate the system or use insider information, and can thus have substantial knowledge about the environment. Hence, as the worst-case assumption, the defender can assume that the attacker has full knowledge about the environment. From this perspective, one-sided POSG can be used to compute robust defense strategies. We restrict to the strictly competitive (or zero-sum) setting. In this case, the defender has guaranteed expected outcome when using such robust strategies even against attackers with less information. Finally, we use the standard assumption that payoffs are computed as discounted sums of immediate rewards. However, our approach could be generalized to the non-discounted version to some extent (by proceeding similarly to [21]).

Our main contribution is the description of the first practical algorithm for computing (approximate) optimal solution for two-player zero-sum one-sided POSG with discounted rewards.² The contribution is threefold: (1) the theoretical contribution proving that our proposed algorithm has guarantees for

¹In this work, however, the authors assumed that the game is turn-taking. In contrast, we consider a more general case where at each timestep, both players choose simultaneously next action to be played.

²Parts of this work appeared in conference publications [20]. This submission is significantly extended from the published works by (1) containing all the proofs and all the technical details regarding the algorithm, (2) full description of the procedure for extracting strategies computed by the algorithm, and (3) new experiments with improved implementation of the algorithm. Finally, we acknowledge that a modification of presented algorithm has been provided in [22, 23] where a compact representation of belief space was proposed for a specific cybersecurity domain and demonstrate that proposed algorithm can scale even beyond experiments however at the cost of losing theoretical guarantees.

approximating the value of any one-sided POSG, (2) showing how to extract strategies from our algorithm and use them to play the game, (3) implementation of the algorithm and experimental evaluation on a set of games. The theoretical work is a direct extension of the theory behind the single-player case (i.e., POMDPs). In POMDPs, an optimal strategy in every step depends on the player’s belief over environment states and on the outcomes achievable in each state. In other words, we have a *value function* which takes a belief b and returns the optimal expected value that can be achieved under b (by following an optimal strategy in both the current decision point and those encountered afterwards).

Figure 1 visualizes the outline and key results provided in each section of the paper. After reviewing related work (Section 2) we state relevant technical background for POMDPs (Section 3). We then formally define one-sided POSG (Section 4) and restate some known results [44] regarding the characteristics of the value function (convexity) and show that the value function can be computed using a recursive formula (Section 5). We then observe that each strategy can be decomposed into the distribution that determines the very next action and the strategy for the remainder of the game and that this structure is mirrored on the level of value functions (Section 6). With these tools, we derive a Bellman equation for one-sided POSG and prove that the iterative application of the corresponding operator H is guaranteed to converge to the optimal value function V^* (Section 7). To get a baseline method of computing V^* , we show that the operator H can be computed using a linear program (Section 8). To get a method with better scaling properties, we design novel approximate algorithms that aim at approximating V^* (Section 9). Namely, we follow the heuristic search value iteration algorithm (HSVI) [39, 40] that uses two functions to approximate the value function, an upper bound function and a lower bound function. By decreasing the gap between these approximations, the algorithm approximates the optimal expected value for relevant belief points. We show that a similar approach can also work in one-sided POSG and that, while the overall idea remains, most of the technical parts of the algorithm have to be adapted for one-sided POSG. We identify and address these technical challenges in order to formally prove that our HSVI algorithm for one-sided POSG converges to optimal strategies. As defined, the HSVI algorithm primarily approximates optimal value for a given game. To extract strategies that reach computed values in expectation, we provide an additional online algorithm (based on ideas from online game-playing algorithms with imperfect information but finite horizon [32]) that generates actions from (approximate) optimal strategies according to the computed approximated value functions (Section 10). Finally, we experimentally evaluate the proposed algorithm on a set of different games, show scalability for these games, and provide deep insights into the performance for each specific part of the algorithm (Section 11). We demonstrate that our implementation of the algorithm is capable of solving non-trivial games with as much as 4 500 states and 120 000 transitions.

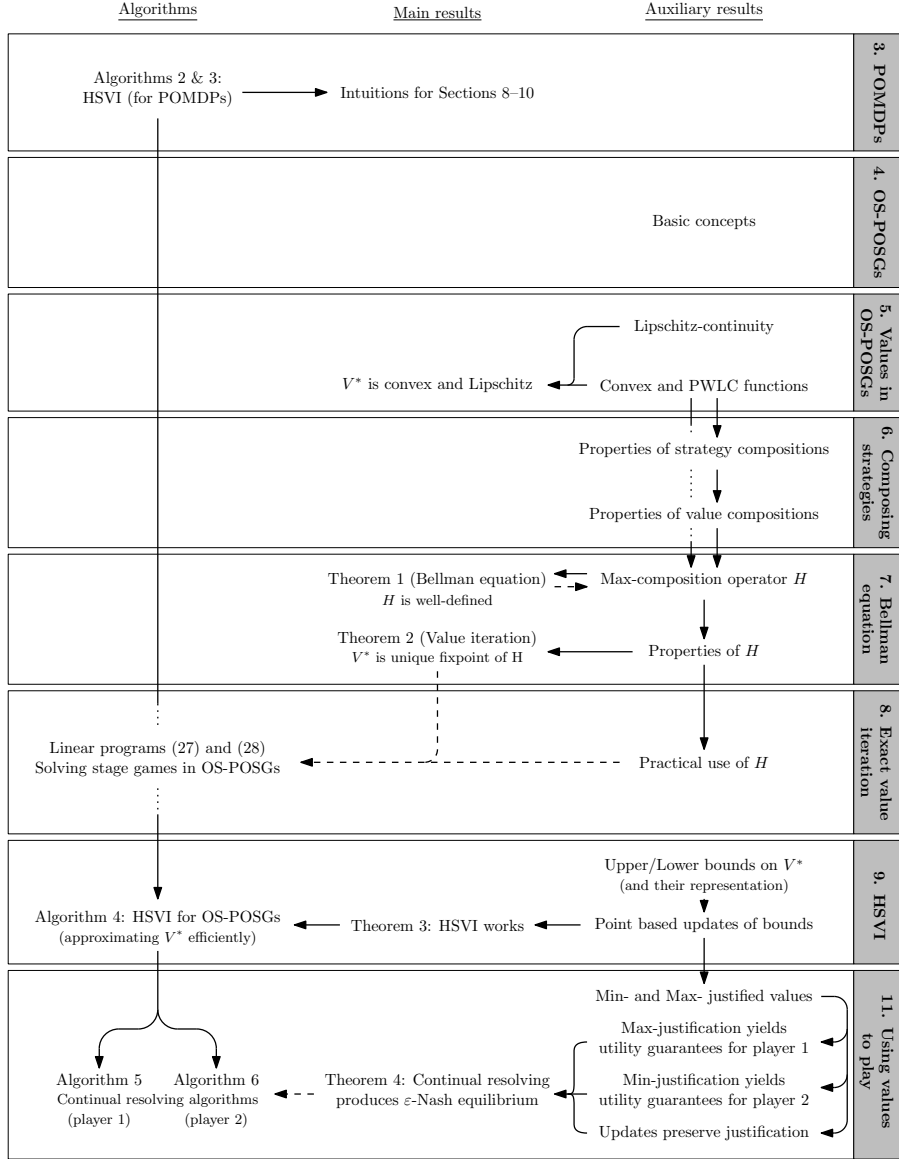


Figure 1: Outline of the theoretical results presented in the paper.

2. Related Work

General, domain-independent algorithms for solving³ (subclasses of) partially observable stochastic games with infinite horizon are not commonly studied. As argued in the introduction, the problem of nested beliefs is one of the reasons. One way of tackling this issue is by using history-dependent strategies. One of the few such approaches is the bottom-up dynamic programming for constructing relevant finite-horizon policy trees for individual players while pruning-out dominated strategies [18, 26]. However, while the history-dependent policies can cope with the necessity of considering the nested beliefs, the number of the strategies is doubly exponential in the horizon of the game (i.e., the number of turns in the game), which greatly limits the scalability and applicability of the algorithm.

We take another approach and restrict to subclasses of POSGs, where the problem of nested beliefs does not appear. Besides the works focused directly on one-sided POSGs, there are other works that consider specific subclasses of POSGs. For example, Ghosh et al. 2004 study zero-sum POSGs with public actions and observations. The authors show that the game has a well-defined value and present an algorithm that exploits the transformation of such a model into a game with complete information. In one-sided POSG, however, the actions are not publicly observable since the imperfectly-informed player lacks the information about their opponent’s action. Compared to existing works studying one-sided POSG [44, 13], our work is the first to provide a practical algorithm that can be directly used to solve games of non-trivial sizes.

Our algorithm focuses on the *offline problem* of (approximately) solving a given one-sided POSG. However, a part of our contribution is the extraction of the strategy that reaches the computed value. On the other hand, *online algorithms* focus on computing strategies that will be used while playing the game. For a long time, no online algorithms for dynamic imperfect-information games provided guarantees on the (near-)optimality of the resulting strategies. While several new algorithms with theoretical guarantees emerged [28, 32, 46] in recent years, they only considered limited-horizon games and produced history-dependent strategies. Using such online algorithms for POSGs is thus only possible with very limited lookahead or when using a heuristic evaluation function. Our approach is fully domain-independent and avoids considering complete histories and the use of evaluation functions while nevertheless being able to consider strategies with horizon of 100 turns or more. Finally, note that the recent work [47] has shown that online algorithms which seem to be consistent with some Nash equilibrium strategy might fail to be “sound” (i.e., there will be a way to exploit them). Fortunately, our algorithm is provably ϵ -sound in this sense, since (the proof of) Theorem 4 shows that it is always guaranteed to get at least the equilibrium value minus ϵ .

³Or even approximating an optimal solution to a given error.

3. Partially Observable MDPs

Partially observable Markov decision processes (POMDPs) [2, 43, 35, 39, 40, 45, 7, 41] are a standard tool for single-agent decision making in stochastic environment under uncertainty about the states. From the perspective of partially observable stochastic games, POMDPs can be seen as a variant of POSG that is only played by a single player.

Definition 3.1 (Partially observable Markov decision process). A *partially observable Markov decision process* is a tuple (S, A, O, T, R) where

- S is a finite set of states,
- A is a finite set of actions the agent can use,
- O is a finite set of observations the agent can observe,
- $T(o, s' \mid s, a)$ is a probability to transition to s' while generating observation o when the current state is s and agent uses action a ,
- $R(s, a)$ is the immediate reward of the agent when using action a in state s .

In POMDPs, the agent starts with a known belief $b^{\text{init}} \in \Delta(S)$ that characterizes the probability $b^{\text{init}}(s)$ that s is the initial state. The play proceeds similarly as in POSGs, except that there is only one decision-maker involved: The initial state $s^{(1)}$ is sampled from the distribution b^{init} . Then, in every stage t , the agent decides about the current action $a^{(t)}$ and receives reward $R(s^{(t)}, a^{(t)})$ based on the current state of the environment $s^{(t)}$. With probability $T(o^{(t)}, s^{(t+1)} \mid s^{(t)}, a^{(t)})$ the system transitions to $s^{(t+1)}$ and the agent receives observation $o^{(t)}$. The decision process is then repeated. Although many objectives have been studied in POMDPs, in this section we discuss only discounted POMDPs with infinite-horizon, i.e., the objective is to optimize $\sum_{t=1}^{\infty} \gamma^{t-1} r_t$ for a discount factor $\gamma \in (0, 1)$.

A strategy $\sigma : (A_1 O)^* \rightarrow A_1$ in POMDPs is traditionally called a *policy* and assigns a deterministic action to each observed history $\omega \in (A_1 O)^*$ of the agent.⁴ Since the agent is the only decision-maker within the environment, and the probabilistic characterization of the environment is known, the player is able to infer his belief $\mathbb{P}_{b^{\text{init}}}[s^{(t+1)} \mid (a^{(i)} o^{(i)})_{i=1}^t]$ (i.e., how likely it is to be in a particular state after a sequence of actions and observations $(a^{(i)} o^{(i)})_{i=1}^t$ has been used and observed). This belief can be defined recursively

$$\tau(b, a, o)(s') = \eta \sum_{s \in S} b(s) \cdot T(o, s' \mid s, a) \quad (1)$$

where η is a normalizing term, and $\tau(b, a, o) \in \Delta(S)$ is the updated belief of the agent when his current belief was b and he played and observed (a, o) . [42] has

⁴As usual, we take X^* to denote the set of all finite sequences over X . For a set Y of sequences, YZ denotes the set of sequences obtained by concatenating a single element of Z to some sequence from Y . (Combining this notation yields, e.g., $axbyc \in (AX)^*A$ for $a, b, c \in A$ and $x, y \in X$.)

shown that the belief of the agent is a sufficient statistic, and POMDPs can therefore be translated into *belief-space MDP*. In theory, standard methods for solving MDPs can be applied, and POMDPs can be solved, e.g., by iterating

$$V^{t+1}(b) = [HV^t](b) = \max_{a \in A} \left[\sum_{s \in S} b(s) \cdot R(s, a) + \gamma \sum_{o \in O} \mathbb{P}_b[o \mid a] \cdot V^t(\tau(b, a, o)) \right]. \quad (2)$$

Since H is a contraction, the repeated application of Equation (2) converges to a unique convex value function $V^* : \Delta(S) \rightarrow \mathbb{R}$ of the POMDP. However, since the number of beliefs is infinite, it is impossible to apply this formula to approximate V^* directly.

Exact value iteration. The value iteration can be, however, rewritten in terms of operations with so-called α -vectors [43]. An α -vector can be seen as a linear function $\alpha : \Delta(S) \rightarrow \mathbb{R}$ characterized by its values $\alpha(s)$ in the vertices $s \in S$ of the belief simplex $\Delta(S)$. We thus have $\alpha(b) = \sum_{s \in S} b(s) \cdot \alpha(s)$.

Assume that V^t is a piecewise-linear and convex function where $V^t(b) = \max_{\alpha \in \Gamma^t} \alpha(b)$ for a finite set of α -vectors Γ^t . We can then form a new (finite) set Γ^{t+1} of α -vectors to represent V^{t+1} from Equation (2) by considering all possible combinations of α -vectors from the set Γ^t :

$$\Gamma^{t+1} = \left\{ \alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha(s) = R(s, a) + \gamma \sum_{(o, s') \in O \times S} T(o, s' \mid s, a) \alpha^o(s') \right. \\ \left. \text{for some } a \in A \text{ and } \alpha^o \in \Gamma^t, o \in O \right\}. \quad (3)$$

As $|\Gamma^{t+1}| = |A| \cdot |\Gamma^t|^{|O|}$, this exact approach suffers from poor scalability. Several techniques have been proposed to reduce the size of sets Γ^t [29, 51], however, this still does not translate to an efficient algorithm.

In the remainder of this section, we present two scalable algorithms for solving POMDPs that are relevant to this thesis. First, we present RTDP-Bel that uses discretized value function and applies Equation (2) directly. Second, we present heuristic search value iteration (HSVI) [39, 40] that inspires our methods for solving POSGs.

RTDP-Bel. The RTDP-Bel algorithm [6] is based on RTDP [3] and has been originally framed in the context of Goal-POMDPs. Goal-POMDPs do not discount rewards (i.e., they set $\gamma = 1$ in Equation (2)). However, the agent is incentivized to reach the goal state g as his reward for every transition before reaching the goal is negative (i.e., it represents the cost). The RTDP-Bel also applies to discounted POMDPs as discounting can be modelled within the Goal-POMDP framework as a fixed probability $1 - \gamma$ of reaching the goal state during every transition [7].

RTDP-Bel adapts RTDP to partially observable domains by using a grid-based approximation of V^* and using a hash-table to store the values, where

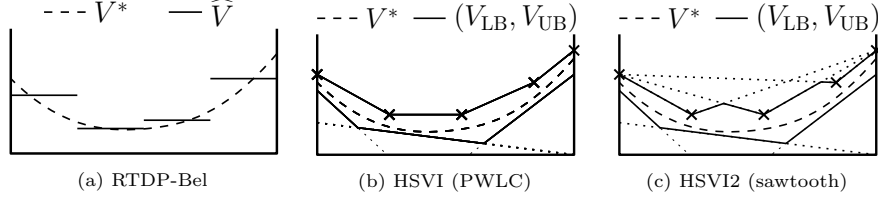


Figure 2: Comparison of value function approximation schemes

$V^*(b) \sim \hat{V}(\lfloor K \cdot b \rfloor)$ for some fixed parameter $K \in \mathbb{N}$. This approximation, however, loses the theoretical properties of RTDP. The algorithm need not converge as the values of the discretized value function may oscillate. Moreover, there is no guarantee that the values stored in the hash-table will provide a bound on the values of V^* [7, p. 3, last paragraph of Section 3]. Despite the lack of theoretical properties, RTDP-Bel has been shown to perform well in practice. The RTDP-Bel algorithm performs a sequence of trials (see Algorithm 1) that updates the discretized value function \hat{V} .

Algorithm 1: A single trial of the RTDP-Bel algorithm.

```

1  $b \leftarrow b^{\text{init}}; s \sim b$ 
2 while  $b(g) < 1$  do
3    $Q(b, a) \leftarrow \sum_{s \in S} b(s)R(s, a) + \sum_{o \in O} \mathbb{P}_b[o \mid a] \cdot \hat{V}(\lfloor K \cdot \tau(b, a, o) \rfloor)$ 
4    $a^* \leftarrow \arg \max_{a \in A} Q(b, a)$ 
5    $\hat{V}(\lfloor K \cdot b \rfloor) \leftarrow Q(b, a^*)$ 
6    $(o, s') \sim T(o, s' \mid s, a^*); b \leftarrow \tau(b, a^*, o); s \leftarrow s'$ 
```

Heuristic search value iteration (HSVI). Heuristic search value iteration [39, 40] is a representative of a class of point-based methods for solving POMDPs. Unlike RTDP-Bel, it approximates V^* using piecewise-linear functions. We illustrate the difference between a grid-based approximation used in RTDP-Bel and a piecewise-linear approximation in Figures 2a and 2b. Observe that unlike the grid-based approximation, a piecewise-linear approximation can yield a close approximation of V^* even in regions with a rapid change of value.

In the original version of the *heuristic-search value iteration* algorithm (HSVI) [39], the algorithm keeps two piecewise-linear and convex (PWLC) functions V_{LB}^Γ and V_{UB}^Υ to approximate V^* (see Figure 2b) and refines them over time. The lower bound on the value is represented in the vector-set representation using a finite set of α -vectors Γ , while the upper bound is formed as a lower convex hull of a set of points $\Upsilon = \{(b_i, y_i) \mid i = 1, \dots, m\}$ where $b_i \in \Delta(S)$ and

$y_i \in \mathbb{R}$. We then have

$$V_{\text{LB}}^\Gamma(b) = \max_{\alpha \in \Gamma} \sum_{s \in S} b(s) \cdot \alpha(s) \quad (4a)$$

$$V_{\text{UB}}^\Upsilon(b) = \min\{\sum_{i=1}^m \lambda_i y_i \mid \lambda \in \mathbb{R}_{\geq 0}^m : \sum_{i=1}^m \lambda_i b_i = b\} . \quad (4b)$$

Computing $V_{\text{UB}}^\Upsilon(b)$ according to Equation (4b) requires solving a linear program. In the second version of the algorithm (HSVI2, [40]), the PWLC representation of upper bound has been replaced by a sawtooth-shaped approximation [19] (see Figure 2c). While the sawtooth approximation is less tight with the same set of points, the computation of $V_{\text{UB}}^\Upsilon(b)$ does not rely on the use of linear programming and can be done in linear time in the size of Υ .

HSVI2 initializes the value function V_{LB}^Γ by considering policies ‘always play the action a ’ and construct one α -vector for each action $a \in A$ corresponding to the expected cost for playing such policy. For the initialization of the upper bound, the fast-informed bound is used [19].

The refinement of V_{LB}^Γ and V_{UB}^Υ is done by adding new elements to the sets Γ and Υ . Since the goal of each update is to improve the approximation quality in the selected belief b as much as possible, we refer to them as *point-based updates* (see Algorithm 2).

Algorithm 2: Point-based `update`(b) procedure of $(V_{\text{LB}}^\Gamma, V_{\text{UB}}^\Upsilon)$.

- 1 $\alpha^{a,o} \leftarrow \arg \max_{\alpha \in \Gamma} \sum_{s' \in S} \tau(b, a, o)(s') \cdot \alpha(s')$ for all $a \in A, o \in O$
 - 2 $\alpha^a(s) \leftarrow R(s, a) + \gamma \sum_{o, s'} T(o, s' \mid s, a) \cdot \alpha^{a,o}(s')$ for all $s \in S, a \in A$
 - 3 $\Gamma \leftarrow \Gamma \cup \{\arg \max_{\alpha^a} \sum_{s \in S} b(s) \cdot \alpha^a(s)\}$
 - 4 $\Upsilon \leftarrow \Upsilon \cup \{(b, \max_{a \in A} [\sum_{s \in S} b(s) R(s, a) + \gamma \sum_{o \in O} \mathbb{P}_b[o \mid a] \cdot V_{\text{UB}}^\Upsilon(\tau(b, a, o))])\}$
-

Algorithm 3: HSVI2 for discounted POMDPs. The pseudocode follows the ZMDP implementation and includes `update` on line 6.

- 1 Initialize V_{LB}^Γ and V_{UB}^Υ
 - 2 **while** $V_{\text{UB}}^\Upsilon(b^{\text{init}}) - V_{\text{LB}}^\Gamma(b^{\text{init}}) > \varepsilon$ **do** `explore`($b^{\text{init}}, \varepsilon, 0$)
 - 3 **procedure** `explore`(b, ε, t)
 - 4 **if** $V_{\text{UB}}^\Upsilon(b) - V_{\text{LB}}^\Gamma(b) \leq \varepsilon \gamma^{-t}$ **then return**
 - 5 $a^* \leftarrow \arg \max_{a \in A} [\sum_s b(s) \cdot R(s, a) + \gamma \sum_{o \in O} \mathbb{P}_b[o \mid a] V_{\text{UB}}^\Upsilon(\tau(b, a, o))]$
 - 6 `update`(b)
 - 7 $o^* \leftarrow \arg \max_{o \in O} \mathbb{P}_b[o \mid a] \cdot \text{excess}_{t+1}(\tau(b, a^*, o))$
 - 8 `explore`($\tau(b, a^*, o^*), \varepsilon, t + 1$)
 - 9 `update`(b)
-

Similarly to RTDP-Bel, HSVI2 selects beliefs where the update should be performed based on the simulated play (selecting actions according to V_{UB}^Υ). Unlike RTDP-Bel, however, observations are not selected randomly. Instead,

HSVI2 selects an observation with the highest *weighted excess gap*, i.e. the excess approximation error

$$\text{excess}_{t+1}(\tau(b, a^*, o)) = V_{\text{UB}}^{\Upsilon}(\tau(b, a^*, o)) - V_{\text{LB}}^{\Gamma}(\tau(b, a^*, o)) - \varepsilon\gamma^{-(t+1)} \quad (5)$$

in $\tau(b, a^*, o)$ weighted by the probability $\mathbb{P}_b[o \mid a^*]$. This heuristic choice attempts to target beliefs where the update will have the most significant impact on $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}})$.

The HSVI2 algorithm for discounted-sum POMDPs ($\gamma \in (0, 1)$) is shown in Algorithm 3. This algorithm provably converges to an ε -approximation of $V^*(b^{\text{init}})$ using values $V_{\text{LB}}^{\Gamma}(b^{\text{init}})$ and $V_{\text{UB}}^{\Upsilon}(b^{\text{init}})$, see [39].

4. Game Model: One-Sided Partially Observable Stochastic Games (OS-POSGs)

We now define the model of one-sided POSGs and describe strategies for this class of games.

Definition 4.1 (one-sided POSGs). A *one-sided POSG* (or OS-POSG) is a tuple $G = (S, A_1, A_2, O, T, R, \gamma)$ where

- S is a finite set of game *states*,
- A_1 and A_2 are finite sets of *actions* of player 1 and player 2, respectively,
- O is a finite set of *observations*
- for every $(s, a_1, a_2) \in S \times A_1 \times A_2$, $T(\cdot \mid s, a_1, a_2) \in \Delta(O \times S)$ represents probabilistic transition function,
- $R : S \times A_1 \times A_2 \rightarrow \mathbb{R}$ is a reward function of player 1,
- $\gamma \in (0, 1)$ is a discount factor.

The game starts by sampling the initial state $s^{(1)} \sim b^{\text{init}}$ from the *initial belief* b^{init} . Then the game proceeds for an infinite number of *stages* where the players choose their actions simultaneously and receive feedback from the environment. At the beginning of i -th stage, the current state $s^{(i)}$ is revealed to player 2, but not to player 1. Then player 1 selects action $a_1^{(i)} \in A_1$ and player 2 selects action $a_2^{(i)} \in A_2$. Based on the current state of the game $s^{(i)}$ and the actions $(a_1^{(i)}, a_2^{(i)})$ taken by the players, an unobservable reward $R(s^{(i)}, a_1^{(i)}, a_2^{(i)})$ is assigned⁵ to player 1, and the game transitions to a state $s^{(i+1)}$ while generating observation $o^{(i)}$ with probability $T(o^{(i)}, s^{(i+1)} \mid s^{(i)}, a_1^{(i)}, a_2^{(i)})$. After committing to action $a_2^{(i)}$, player 2 observes the entire outcome of the current stage, including the action $a_1^{(i)}$ taken by player 1 and the observation $o^{(i)}$. player 1, on the other hand, knows only his own action $a_1^{(i)}$ and the observation $o^{(i)}$, while the action

⁵Note that we consider a zero-sum setting, hence the reward of player 2 is $-R(s^{(i)}, a_1^{(i)}, a_2^{(i)})$. We do however consider that player 2 focuses on minimizing the reward of player 1 instead of reasoning about the rewards of player 2 directly.

$a_2^{(i)}$ of player 2 and both the past and new states of the system $s^{(i)}$ and $s^{(i+1)}$ remain unknown to him.

The information asymmetry in the game means that while player 2 can observe entire course of the game $(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^t s^{(t+1)} \in (SA_1A_2O)^*S$ up to the current decision point at time $t+1$, player 1 only knows his own actions and observations $(a_1^{(i)} o_{i=1}^{(i)})^t \in (A_1O)^*$.⁶ The players make decisions solely based on this information - formally, this is captured by the following definition:

Definition 4.2 (Behavioral strategy). Let G be a one-sided POSG. Mappings $\sigma_1 : (A_1O)^* \rightarrow \Delta(A_1)$ and $\sigma_2 : (SA_1A_2O)^*S \rightarrow \Delta(A_2)$ are *behavioral strategies* of imperfectly informed player 1 and perfectly informed player 2, respectively. The sets of all behavioral strategies of player 1 and player 2 are denoted Σ_1 and Σ_2 , respectively.

Plays in OS-POSGs. Players use their behavioral strategies (σ_1, σ_2) to play the game. A *play* is an infinite word $(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^\infty$, while finite prefixes of plays $w = (s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^T s^{(T+1)}$ are called *histories* of length T , and plays having w as a prefix are denoted $\text{Cone}(w)$. Formally, a *cone* of w is a set of all plays extending w ,

$$\text{Cone}(w) := \left\{ (s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^\infty \in (SA_1A_2O)^* \mid (s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^\infty \text{ extends } w \right\}. \quad (6)$$

At a decision point at time t , players extend a history $(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^t s^{(t+1)}$ of length t by sampling actions from their strategies $a_1^{(t+1)} \sim \sigma_1((a_1^{(i)} o^{(i)})_{i=1}^t)$ and $a_2^{(t+1)} \sim \sigma_2((s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^t s^{(t+1)})$. We consider a discounted-sum objective with discount factor $\gamma \in (0, 1)$. The payoff associated with a play $(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^\infty$ is thus $\text{Disc}^\gamma := \sum_{i=1}^\infty \gamma^{i-1} R(s^{(i)}, a_1^{(i)}, a_2^{(i)})$. Player 1 is aiming to maximize this quantity while player 2 is minimizing it.

Apart from reasoning about decision rules of the players for the entire game (i.e., their behavioural strategies σ_1 and σ_2), we also consider the strategies they use for a single decision point—or stage—of the game only (i.e., assuming that the course of the previous stages $(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^t$ is fixed and considered a parameter of the given stage).

Definition 4.3 (Stage strategy). Let G be a one-sided POSG. A *stage strategy* of player 1 is a distribution $\pi_1 \in \Delta(A_1)$ over the actions player 1 can use at the current stage. A *stage strategy* of player 2 is a mapping $\pi_2 : S \rightarrow \Delta(A_2)$ from the possible current states of the game (player 2 observes the true state at the beginning of the current stage) to a distribution over actions of player 2.

⁶Recall that we use the standard notation where $X^* :=$ all finite sequences over X (and, if Y is a set of sequences, YZ denotes the set of sequences obtained by appending a single element of Z at the end of some $y \in Y$).

The sets of all stage strategies of player 1 and player 2 are denoted Π_1 and Π_2 , respectively.

Note that a stage strategy of player 2 is essentially a conditional probability distribution given the current state of the game. For the reasons of notational convenience, we use notation $\pi_2(a_2 | s)$ instead of $\pi_2(s)(a_2)$ wherever applicable.

4.1. Subgames

Recall that both players know past actions of player 1 and all observations player 1 has received. The action-observation history is thus public knowledge. This allows us to define a notion of *subgames*. A subgame induced by an action-observation history ω (or ω -subgame) is formed by histories h such that the action-observation history $\omega(h)$ of player 1 in h is a suffix of ω , i.e., $\omega(h) \sqsupseteq \omega$.

Later in the text, we will specifically reason about subgames that follow directly after the first stage of the game—these correspond to (a_1, o) -subgames for some action a_1 and observation o . Observe that, once (a_1, o) is played and observed, both players know exactly which (a_1, o) -subgame they are currently in. Consequently, reasoning about (a_1, o) -subgame can be done without considering any other (a'_1, o') -subgame.

4.2. Probability measures

We now proceed by defining a probability measure on the space of infinite plays in one-sided POSGs. Assuming that $b \in \Delta(S)$ is the initial belief characterizing the distribution over possible initial states, and players use strategies (σ_1, σ_2) to play the game from the current situation, we can define the probability distribution over histories (i.e., prefixes of plays) recursively as follows.

$$\mathbb{P}_{b, \sigma_1, \sigma_2}[s^{(1)}] = b(s^{(1)}) \quad (7a)$$

$$\begin{aligned} \mathbb{P}_{b, \sigma_1, \sigma_2}[(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^t s^{(t+1)}] &= \mathbb{P}_{b, \sigma_1, \sigma_2}[(s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^{t-1} s^{(t)}] \cdot \\ &\cdot \sigma_1((a_1^{(i)} o^{(i)})_{i=1}^{t-1}, a_1^{(t)}) \cdot \sigma_2((s^{(i)} a_1^{(i)} a_2^{(i)} o^{(i)})_{i=1}^{t-1} s^{(t)}, a_2^{(t)}) \cdot \\ &\cdot T(o^{(t)}, s^{(t+1)} | s^{(t)}, a_1^{(t)}, a_2^{(t)}) \end{aligned} \quad (7b)$$

This probability distribution also coincide with a measure μ defined over the cones, i.e. plays having w as a prefix.

$$\mu(\text{Cone}(w)) = \mathbb{P}_{b, \sigma_1, \sigma_2}[w] \quad (8)$$

The measure μ uniquely extends to the probability measure $\mathbb{P}_{b, \sigma_1, \sigma_2}[\cdot]$ over infinite plays of the game, which allows us to define the expected utility $\mathbb{E}_{b, \sigma_1, \sigma_2}[\text{Disc}^\gamma]$ of the game when the initial belief of the game is b and strategies $\sigma_1 \in \Sigma_1$ and $\sigma_2 \in \Sigma_2$ are played by player 1 and player 2, respectively.

In a similar manner, we can define a probability measure $\mathbb{P}_{b, \pi_1, \pi_2}[s, a_1, a_2, o, s']$ that predicts events only one step into the future (for *stage* strategies $\pi_1 \in \Pi_1$, $\pi_2 \in \Pi_2$). For belief b and stage strategies π_1, π_2 , we consider the probability that a stage starts in state $s \in S$ (sampled from b), players select actions $a_1 \sim \pi_1$

and $a_2 \sim \pi_2$, and that this results into a transition to a new state $s' \in S$ while generating an observation $o \in O$:

$$\mathbb{P}_{b,\pi_1,\pi_2}[s, a_1, a_2, o, s'] = b(s)\pi_1(a_1)\pi_2(a_2 | s)T(o, s' | s, a_1, a_2) . \quad (9)$$

The probability distribution in Equation (9) can be marginalized to obtain, e.g., the probability that player 1 plays action $a_1 \in A_1$ and observes $o \in O$,

$$\begin{aligned} \mathbb{P}_{b,\pi_1,\pi_2}[a_1, o] &= \sum_{(s,a_2,s') \in S \times A_2 \times S} \mathbb{P}_{b,\pi_1,\pi_2}[s, a_1, a_2, o, s'] \\ &= \sum_{(s,a_2,s') \in S \times A_2 \times S} b(s)\pi_1(a_1)\pi_2(a_2 | s)T(o, s' | s, a_1, a_2) . \end{aligned} \quad (10)$$

At the beginning of each stage, the imperfectly informed player 1 selects their action based on their belief about the current state of the game. For a fixed current stage-strategy π_2 of player 2, player 1 can derive the distribution over possible states at the beginning of the next stage. If player 1 starts with a belief b , takes an action $a_1 \in A_1$, and observes $o \in O$, his updated belief $b' = \tau(b, a_1, \pi_2, o)$ over states $s' \in S$ is going to be $\tau(b, a_1, \pi_2, o)(s') =$

$$= \mathbb{P}_{b,\pi_1,\pi_2}[s' | a_1, o] = \sum_{(s,a_2) \in S \times A_2} \mathbb{P}_{b,\pi_1,\pi_2}[s, a_2, s' | a_1, o] \quad (11a)$$

$$= \frac{1}{\mathbb{P}_{b,\pi_1,\pi_2}[a_1, o]} \sum_{(s,a_2) \in S \times A_2} \mathbb{P}_{b,\pi_1,\pi_2}[s, a_1, a_2, o, s'] \quad (11b)$$

$$= \frac{1}{\mathbb{P}_{b,\pi_1,\pi_2}[a_1, o]} \sum_{(s,a_2) \in S \times A_2} b(s)\pi_1(a_1)\pi_2(a_2 | s)T(o, s' | s, a_1, a_2) . \quad (11c)$$

In Section 7, this expression will prove useful for describing the Bellman equation in one-sided POSGs.

5. Value of One-Sided POSGs

We now proceed by establishing the value function of one-sided POSGs. The value function represents the utility player 1 can achieve in each possible initial belief of the game. First, we define the value of a strategy $\sigma_1 \in \Sigma_1$ of player 1, which assigns a payoff player 1 is guaranteed to get by playing σ_1 in the game (parameterized by the initial belief of the game). Based on the value of strategies, we define the optimal value function of the game where player 1 chooses the best strategy for the given initial belief.

Definition 5.1 (Value of strategy). Let G be a one-sided POSG and $\sigma_1 \in \Sigma_1$ be a behavioral strategy of the imperfectly informed player 1. The *value of strategy* σ_1 , denoted val^{σ_1} , is a function mapping each belief $b \in \Delta(S)$ to the expected utility that σ_1 guarantees against a best-responding player 2 given that the initial belief is b :

$$\text{val}^{\sigma_1}(b) = \inf_{\sigma_2 \in \Sigma_2} \mathbb{E}_{b,\sigma_1,\sigma_2}[\text{Disc}^\gamma] . \quad (12)$$

When given an instance of a one-sided POSG with initial belief b , player 1 aims for a strategy that yields the best possible expected utility $\text{val}^{\sigma_1}(b)$. The value player 1 can guarantee in belief b is characterized by the optimal value function V^* of the game.

Definition 5.2 (Optimal value function). Let G be a one-sided POSG. The *optimal value function* $V^* : \Delta(S) \rightarrow \mathbb{R}$ of G represents the supinf value of player 1 for each of the beliefs, i.e.

$$V^*(b) = \sup_{\sigma_1 \in \Sigma_1} \text{val}^{\sigma_1}(b) . \quad (13)$$

Note that according to von Neumann's minimax theorem [49] (resp. its generalization [38]), every zero-sum POSG with discounted-sum objective Disc^γ is determined in the sense that the lower values (in the supinf sense) and the upper values (in the infsup sense) of the game coincide and represent the value of the game. Therefore, $V^*(b)$ also represents the value of the game when the initial belief of the game is $b \in \Delta(S)$.

Since the Disc^γ objective is considered (for $0 < \gamma < 1$), the infinite discounted sum of rewards of player 1 converge. As a result, the values of strategies $\text{val}^{\sigma_1}(b)$ and the value of the game $V^*(b)$ can be bounded.

Proposition 5.3. *Let G be a one-sided POSG. Then the payoff Disc^γ of an arbitrary play in G is bounded by values*

$$L = \min_{(s, a_1, a_2)} R(s, a_1, a_2)/(1 - \gamma) \quad U = \max_{(s, a_1, a_2)} R(s, a_1, a_2)/(1 - \gamma) . \quad (14)$$

It also follows that $L \leq V^(b) \leq U$ and $L \leq \text{val}^{\sigma_1}(b) \leq U$ holds for every belief $b \in \Delta(S)$ and strategy $\sigma_1 \in \Sigma_1$ of the imperfectly informed player 1.*

Since the values L and U are uniquely determined by the given one-sided POSG, we will use these symbols in the remainder of the text. We now focus on the discussion of structural properties of solutions of OS-POSGs. First, we show that the value of an arbitrary strategy $\sigma_1 \in \Sigma_1$ of player 1 is linear in $b \in \Delta(S)$ — that is, it can be represented as a convex combination of its values in the vertices of the simplex $\Delta(S)$.

In accordance with the notation used in the POMDP literature, we refer to linear functions defined over the $\Delta(S)$ simplex as α -vectors. For $s \in S$, we overload the notation as $\alpha(s) :=$ the value of α in the vertex corresponding to s . This allows us to write the following for every $b \in \Delta(S)$

$$\alpha(b) = \sum_{s \in S} \alpha(s) \cdot b(s) \quad \text{where } \alpha(s) = \alpha(\mathbb{1}_s), \quad \mathbb{1}_s(s') = \begin{cases} 1 & s = s' \\ 0 & \text{otherwise} \end{cases} . \quad (15)$$

The following lemma shows the result we promised earlier:

Lemma 5.4. *Let G be a one-sided POSG and $\sigma_1 \in \Sigma_1$ be an arbitrary behavioral strategy of player 1. Then the value val^{σ_1} of strategy σ_1 is a linear function in the belief space $\Delta(S)$.*

Proof. According to the Definition 5.1, the value val^{σ_1} of strategy σ_1 is defined as the expected utility of σ_1 against the best-response strategy σ_2 of player 2. However, before having to act, player 2 observes the true initial state $s \sim b$. Therefore, he will play a best-response strategy σ_2 against σ_1 (with expected utility $\text{val}^{\sigma_1}(s)$) given that the initial state is s . Since the probability that the initial state is s is $b(s)$, we have

$$\text{val}^{\sigma_1}(b) = \sum_{s \in S} b(s) \text{val}^{\sigma_1}(s) . \quad (16)$$

This shows that val^{σ_1} is a linear function in the belief $b \in \Delta(S)$. \square

Since a point-wise supremum of a set of linear functions is convex, Lemma 5.4 implies that the optimal value function V^* is convex:

Lemma 5.5. *Optimal value function V^* of a one-sided POSG is convex.*

Unless otherwise specified, we endow any space $\Delta(X)$ over a finite set X with the $\|\cdot\|_1$ metric. To prepare the ground for the later proof of correctness of our main algorithm (presented in Section 9), we now show that both the value of strategies and the optimal value function V^* are Lipschitz continuous. (Recall that for $k > 0$ a function $f : \Delta(X) \rightarrow \mathbb{R}$ is k -Lipschitz continuous if for every $p, q \in \Delta(X)$ it holds $|f(p) - f(q)| \leq k \cdot \|p - q\|_1$.)

Lemma 5.6. *Let X be a finite set and let $f : \Delta(X) \rightarrow [y_{\min}, y_{\max}]$ be a linear function. Then f is k -Lipschitz continuous for $k = (y_{\max} - y_{\min})/2$.*

Lemma 5.6 directly implies that both values val^{σ_1} of strategies σ_1 of the imperfectly informed player 1, as well as the optimal value function V^* are Lipschitz continuous.

Lemma 5.7. *Let $\sigma_1 \in \Sigma_1$ be an arbitrary strategy of the imperfectly informed player 1. Then val^{σ_1} is $(U - L)/2$ -Lipschitz continuous.⁷*

Proof. Value val^{σ_1} of strategy σ_1 is linear (Lemma 5.4) and its values are bounded by L and U (Proposition 5.3). Therefore, according to Lemma 5.6, the function val^{σ_1} is $(U - L)/2$ -Lipschitz. \square

For notational convenience, we denote this constant as $\delta := (U - L)/2$ in the remainder of the text.

Proposition 5.8. *Value function V^* of one-sided POSGs is δ -Lipschitz continuous.*

Remark. In the remainder of the text, we will use term *value function* to refer to an arbitrary function $V : \Delta(S) \rightarrow \mathbb{R}$ that assigns numbers $V(b)$ (estimates of the value achieved under optimal play) to beliefs $b \in \Delta(S)$ of player 1.

⁷Recall that L and U , introduced in Proposition 5.3, are the minimum and maximum possible utilities in the game.

5.1. Elementary Properties of Convex Functions

In 16 5.5, we have shown that the optimal value function V^* of one-sided POSGs is convex. In this section, we will explicitly state some of the important properties of convex functions that motivate our approach and are used throughout the rest of the text.

Proposition 5.9. *Let $f : \Delta(S) \rightarrow \mathbb{R}$ be a point-wise supremum of linear functions, i.e.,*

$$f(b) = \sup_{\alpha \in \Gamma} \alpha(b), \quad \Gamma \subseteq \{\alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha \text{ is linear}\}. \quad (17)$$

Then f is convex and continuous. Furthermore, if every $\alpha \in \Gamma$ is k -Lipschitz continuous, f is k -Lipschitz continuous as well.

Proof. Let $b, b' \in \Delta(S)$ and $\lambda \in [0, 1]$ be arbitrary. We have

$$\begin{aligned} \lambda f(b) + (1 - \lambda)f(b') &= \lambda \sup_{\alpha \in \Gamma} \alpha(b) + (1 - \lambda) \sup_{\alpha \in \Gamma} \alpha(b') \\ &= \sup_{\alpha \in \Gamma} \lambda \alpha(b) + \sup_{\alpha \in \Gamma} (1 - \lambda) \alpha(b') \\ &\geq \sup_{\alpha \in \Gamma} [\lambda \alpha(b) + (1 - \lambda) \alpha(b')] \\ &= \sup_{\alpha \in \Gamma} \alpha(\lambda b + (1 - \lambda)b') \\ &= f(\lambda b + (1 - \lambda)b'), \end{aligned}$$

which shows that f is convex.

We now prove the continuity of f . Since every convex function is continuous on the interior of its domain, it remains to show that f is continuous on the boundary of $\Delta(S)$. Assume to the contradiction that it is not continuous, i.e., there exists b_0 on the boundary such that for all b from its neighborhood $f(b_0) > f(b) + C$ for some $C > 0$. Since f is a pointwise supremum of linear functions, there exists $\alpha \in \Gamma$ such that $\alpha(b_0) > f(b_0) - C/2$. However, at the same time, we have $\alpha(b) \leq f(b) - C$. This is in contradiction with the fact that all $\alpha \in \Gamma$ are linear, and hence continuous.

Furthermore, suppose that every $\alpha \in \Gamma$ is k -Lipschitz continuous and let $b, b' \in \Delta(S)$. We have

$$\begin{aligned} f(b) &= \sup_{\alpha \in \Gamma} \alpha(b) \\ &\leq \sup_{\alpha \in \Gamma} [\alpha(b') + k\|b - b'\|_1] \quad (\text{since every } \alpha \in \Gamma \text{ is } k\text{-Lipschitz}) \\ &= \left[\sup_{\alpha \in \Gamma} \alpha(b') \right] + k\|b - b'\|_1 \\ &= f(b') + k\|b - b'\|_1. \end{aligned}$$

Since the identical argument proves the inequality $f(b') \leq f(b) + k\|b - b'\|_1$, this shows that f is k -Lipschitz continuous. \square

Recall that we aim to emulate the HSVI algorithm from POMDPs, where the optimal value function V^* is approximated by a series of piecewise linear and convex functions. One of the common ways to represent these functions is as a point-wise maximum of a finite set of linear functions (typically called α -vectors in the POMDP context):

Definition 5.10 (Piecewise linear and convex function on $\Delta(S)$). A function $f : \Delta(S) \rightarrow \mathbb{R}$ is said to be *piecewise linear and convex* (PWLC) if it is of the form $f(b) = \max_{\alpha \in \Gamma} \alpha(b)$ (for each $b \in \Delta(S)$) for some finite set $\Gamma \subset \{\alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha \text{ is linear}\}$.

We immediately see that the preceding Proposition 5.9 applies to any function of this type. The next result shows that PWLC functions remain unchanged if we replace the set Γ by its convex hull:

Proposition 5.11. *Let $\Gamma \subset \{\alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha \text{ is linear}\}$ be a set of linear functions. Then for every $b \in \Delta(S)$ we have*

$$\sup_{\alpha \in \Gamma} \alpha(b) = \sup_{\alpha \in \text{Conv}(\Gamma)} \alpha(b) . \quad (18)$$

In the opposite direction, every convex function can be represented as a supremum over some set of linear functions. The following proposition shows this using the largest possible set, i.e. $\{\alpha \leq f \mid \alpha \text{ linear}\}$:

Proposition 5.12. *Let $f : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function. Then there exists a set Γ of linear functions such that $\alpha \leq f$ for every $\alpha \in \Gamma$ and $f(b) = \sup_{\alpha \in \Gamma} \alpha(b)$ for every $b \in \Delta(S)$.*

6. Composing Strategies

Every behavioural strategy of the imperfectly informed player 1 can be split into the stage strategy π_1 player 1 uses in the first stage of the game, and behavioural strategies he uses in the rest of the game after he reaches an (a_1, o) -subgame. We can also use the inverse principle, called *strategy composition*, to form new strategies by choosing the stage strategy π_1 for the first stage and then selecting a separate behavioral strategy $\bar{\zeta} = (\zeta_{a_1, o})_{(a_1, o) \in A_1 \times O}$ for each subgame (see Figure 3 for illustration).

Definition 6.1 (Strategy composition). Let G be a one-sided POSG and $\pi_1 \in \Pi_1$ a stage strategy of player 1. Furthermore, let $\bar{\zeta} \in (\Sigma_1)^{A_1 \times O}$ be a vector representing behavioral strategies of player 1 for each (a_1, o) -subgame where $a_1 \in A_1$ and $o \in O$. The *strategy composition* $\text{comp}(\pi_1, \bar{\zeta})$ is a behavioral strategy of player 1 such that

$$\text{comp}(\pi_1, \bar{\zeta})(\omega) = \begin{cases} \pi_1 & \omega = \emptyset \\ \zeta_{a_1, o}(\omega') & \omega = a_1 o \omega' \end{cases} \quad \text{for each } \omega \in (A_1 O)^* . \quad (19)$$

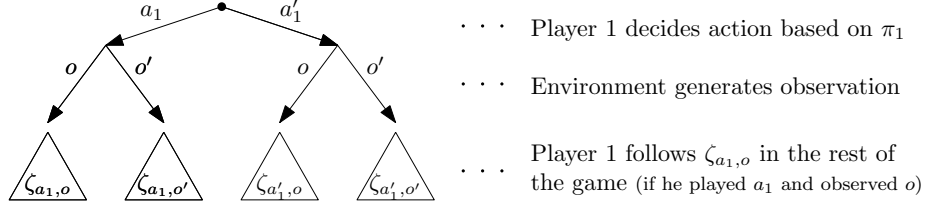


Figure 3: Composition of strategies ζ using a stage strategy π_1 .

By composing strategies $\bar{\zeta}$ using π_1 , we obtain a new strategy where the probability of playing a_1 in the first stage of the game is $\pi_1(a_1)$, and strategy $\zeta_{a_1,o}$ is used after playing action a_1 and receiving observation o in the first stage of the game. Importantly, the newly formed strategy $\text{comp}(\pi_1, \bar{\zeta}) \in \Sigma_1$ is also a behavioral strategy (of imperfectly informed player 1), and therefore the properties of strategies presented in Section 5 apply also to $\text{comp}(\pi_1, \bar{\zeta})$. As the next result shows, the opposite property also holds — for each strategy $\sigma_1 \in \Sigma_1$ of player 1, we can find the appropriate π_1 and $\bar{\zeta}$ such that $\sigma_1 = \text{comp}(\pi_1, \bar{\zeta})$:

Proposition 6.2. *Every behavioral strategy $\sigma_1 \in \Sigma_1$ of player 1 can be represented as a strategy composition of some stage strategy $\pi_1 \in \Pi_1$ and player 1 behavioral strategies $\zeta_{a_1,o}$.*

Importantly, we can obtain values $\text{val}^{\text{comp}(\pi_1, \bar{\zeta})}$ of composite strategies without considering the entire strategy $\text{comp}(\pi_1, \bar{\zeta})$. As the following lemma shows, it suffices to consider only the first stage of the game and the *values* of the strategies $\bar{\zeta} \in (\Sigma_1)^{A_1 \times O}$.

Lemma 6.3. *Let G be a one-sided POSG and $\text{comp}(\pi_1, \bar{\zeta})$ a composite strategy. Then the following holds:*

$$\begin{aligned} \text{val}^{\text{comp}(\pi_1, \bar{\zeta})}(s) &= \min_{a_2 \in A_2} \mathbb{E}_{a_1 \sim \pi_1, (o, s') \sim T(\cdot | s, a_1, a_2)} \left[R(s, a_1, a_2) + \gamma \text{val}^{\zeta_{a_1, o}}(s') \right] \\ &= \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1(a_1) \left[R(s, a_1, a_2) + \gamma \sum_{(o, s') \in O \times S} T(o, s' | s, a_1, a_2) \text{val}^{\zeta_{a_1, o}}(s') \right]. \end{aligned} \tag{20}$$

The proof relies on the fact that when player 1 takes the action a_1 , observes o , and ends up in s' , the strategy $\zeta_{a_1,o}$ guarantees the player gets at least $\text{val}^{\zeta_{a_1,o}}(s')$ utility (in expectation), no matter what player 2 does. Since the values in the rest of the game are known, it suffices to focus on the best-response strategy of player 2 in the first stage of the game.

6.1. Generalized Composition

Lemma 6.3 suggests that we can use composition of *values* of strategies $\text{val}^{\zeta_{a_1,o}}$ to form values of composite strategies $\text{val}^{\text{comp}(\pi_1, \bar{\zeta})}$. In this section,

still consider linear functions $\text{val}^{\zeta_{a_1,o}}$, but we relax the assumption that these functions represent values of some specific behavioural strategy. This allows us to derive a generalized principle of composition and approximate the value function V^* by a supremum of arbitrary linear functions (as opposed to functions val^{σ_1}). Throughout the text, we will use $\text{lin}_{\Delta(S)}$ to denote the set of linear functions on $\Delta(S)$ (i.e., α -vectors). We will also use the term ‘linear’ to refer to functions that satisfy $f(\lambda b + (1 - \lambda)b') = \lambda f(b) + (1 - \lambda)f(b')$ on $\Delta(S)$.

Definition 6.4 (Value composition). Let $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in (\text{lin}_{\Delta(S)})^{A_1 \times O}$. *Value composition* $\text{valcomp}(\pi_1, \bar{\alpha}) : \Delta(S) \rightarrow \mathbb{R}$ is a linear function defined by the values in vertices of the $\Delta(S)$ simplex as follows:

$$\begin{aligned} \text{valcomp}(\pi_1, \bar{\alpha})(s) = \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1(a_1) & \left[R(s, a_1, a_2) + \right. \\ & \left. \gamma \sum_{(o,s') \in O \times S} T(o, s' | s, a_1, a_2) \alpha_{a_1,o}(s') \right]. \end{aligned} \quad (21)$$

Observe that according to Lemma 6.3, $\text{valcomp}(\pi_1, \bar{\alpha}) = \text{val}^{\text{comp}(\pi_1, \bar{\zeta})}$ for $\alpha_{a_1,o} = \text{val}^{\zeta_{a_1,o}}$. The value composition $\text{valcomp}(\pi_1, \bar{\alpha})$, however, admits arbitrary linear function $\alpha_{a_1,o}$ and not only the value $\text{val}^{\zeta_{a_1,o}}$ of some strategy $\zeta_{a_1,o} \in \Sigma_1$. Moreover, as long as linear functions $\alpha_{a_1,o}$ serve as lower bounds for values of some strategies, so will the corresponding value composition serve as a lower bound for the corresponding composite strategy:

Lemma 6.5. *Let $\pi_1 \in \Pi_1$ be a stage strategy of player 1 and $\bar{\alpha} \in (\text{lin}_{\Delta(S)})^{A_1 \times O}$ a vector of linear functions s.t. for each $\alpha_{a_1,o}$ there exists a strategy $\zeta_{a_1,o} \in \Sigma_1$ with $\text{val}^{\zeta_{a_1,o}} \geq \alpha_{a_1,o}$. Then there exists a strategy $\sigma_1 \in \Sigma_1$ such that $\sigma_1(\emptyset) = \pi_1$ and $\text{val}^{\sigma_1} \geq \text{valcomp}(\pi_1, \bar{\alpha})$.*

In case of value of composite strategies, we know that $\text{val}^{\text{comp}(\pi_1, \zeta)}$ is a δ -Lipschitz continuous linear function (since $\text{comp}(\pi_1, \zeta) \in \Sigma_1$ is a behavioral strategy of player 1 and Lemma 5.7 applies). Additionally, we prove that as long as linear functions $\alpha_{a_1,o}$ are bounded by $L \leq \alpha_{a_1,o}(b) \leq U$ for every belief $b \in \Delta(S)$, and are therefore δ -Lipschitz continuous, the value composition $\text{valcomp}(\pi_1, \bar{\alpha})$ is also δ -Lipschitz.

Lemma 6.6. *Let $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in (\text{lin}_{\Delta(S)})^{A_1 \times O}$ such that $L \leq \alpha_{a_1,o}(b) \leq U$ for every $b \in \Delta(S)$. Then $L \leq \text{valcomp}(\pi_1, \bar{\alpha})(b) \leq U$ for every $b \in \Delta(S)$ and $\text{valcomp}(\pi_1, \bar{\alpha})$ is a δ -Lipschitz continuous function.*

7. Bellman Equation for One-Sided POSGs

In Section 5, we have defined the value function V^* as the supremum over the strategies player 1 can achieve in each of the beliefs (see Definition 5.2). However, while this correctly defines the value function, it does not provide a straightforward recipe to obtaining value $V^*(b)$ for the given belief $b \in \Delta(S)$.

Obtaining the value for the given belief according to Definition 5.2 is as hard as solving the game itself.

In this section, we provide an alternative characterization of the optimal value function V^* inspired by the value iteration methods, e.g., for Markov decision processes (MDPs) and their partially observable variant (POMDPs). The high-level idea behind these approaches is to start with a coarse approximation $V_0 : \Delta(S) \rightarrow \mathbb{R}$ of the value function V^* , and then iteratively improve the approximation by applying the Bellman's operator H , i.e., generate a sequence such that $V_{i+1} = HV_i$. In our case, the improvement is based on finding a new, previously unknown, strategy that achieves higher values for each of the beliefs by means of value composition principle (Definition 6.4). Throughout this section, we will consider value functions that are represented as a point-wise supremum over a (possibly infinite) set Γ of linear functions (called α -vectors), i.e.,

$$V(b) = \sup_{\alpha \in \Gamma} \alpha(b) \quad \text{for } \Gamma \subset \{\alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha \text{ is linear}\} . \quad (22)$$

By Proposition 5.11, we can always assume that the set Γ is convex (since this doesn't come at the loss of generality). For more details on this representation of value functions see Section 5.1.

Definition 7.1 (Max-composition). Let $V : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function and let Γ be a convex set of linear functions such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$. The *max-composition* operator H is defined as

$$[HV](b) = \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) . \quad (23)$$

We will now prove several fundamental properties of the max-composition operator H from Definition 7.1. First, we will show that this operator preserves continuity and convexity, allowing us to apply the operator iteratively. Second, we introduce equivalent formulations of the operator H , which represent the solution of $[HV](b)$ in a more traditional form of finding a Nash equilibrium of a stage-game. These formulations also allow us to show that the behaviour of H is not sensitive to the choice of the set Γ used to represent the value function V . Finally, we conclude by showing that the operator H can indeed be used to approximate the optimal value function V^* . Namely, we show that H is a contraction mapping (and thus iterated application converges to a unique fixpoint) and that its fixpoint is the optimal value function V^* .

Proposition 7.2. *Proposition Let $V : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function and let Γ be a convex set of linear functions such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$. Then HV is also convex and continuous. Furthermore, if V is δ -Lipschitz continuous, the function HV is δ -Lipschitz continuous as well.*

The proof of this result goes by rewriting HV as a supremum over all value-compositions and using our earlier observations about convexity and Lipschitz continuity of such suprema.

We will now prove that the max-composition operator H can be alternatively characterized using max-min and min-max optimization. Recall that $\tau(b, a_1, \pi_2, o)$ denotes the Bayesian update of belief b given that player 1 played a_1 and observed o , and player 2 is assumed to follow stage strategy π_2 in the current round (see Equation (11)).

Theorem 1. *Let $V : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function and let Γ be a convex set of linear functions on $\Delta(S)$ such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$ for every belief $b \in \Delta(S)$. Then the following definitions of operator H are equivalent:*

$$\begin{aligned} [HV](b) &= \\ &= \max_{\pi_1 \in \Delta(S)} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) \end{aligned} \quad (24a)$$

$$= \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} \left[\mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)) \right] \quad (24b)$$

$$= \min_{\pi_2 \in \Pi_2} \max_{\pi_1 \in \Pi_1} \left[\mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)) \right]. \quad (24c)$$

The proof consists of verifying the assumptions of von Neumann's minimax theorem, which shows the equivalence of (24b) and (24c). The equivalence of (24b) and (24a) can be then shown by reformulating each stage game as a separate zero-sum game and verifying that it satisfies the assumptions of a Sion's generalization of the minimax theorem [38].

Corollary 7.3. *Bellman's operator H does not depend on the convex set Γ of linear functions used to represent the convex value function V .*

Since the maximin and minimax values of the game (from equations (24b) and (24c)) coincide, the value $[HV](b)$ corresponds to the Nash equilibrium in the stage game. We define the stage game formally.

Definition 7.4 (Stage game). A *stage game* with respect to a convex continuous value function $V : \Delta(S) \rightarrow \mathbb{R}$ and belief $b \in \Delta(S)$ is a two-player zero sum game with strategy spaces Π_1 for the maximizing player 1 and Π_2 for the minimizing player 2, and payoff function

$$u^{V,b}(\pi_1, \pi_2) = \mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)). \quad (25)$$

With a slight abuse of notation, we use $[HV](b)$ to refer both to the max-composition operator (Definition 7.1) as well as to this stage game.

We will now show that the Bellman's operator H is a contraction mapping. Recall that the mapping H is a contraction, if there exists $0 \leq k < 1$ such that $\|HV_1 - HV_2\| \leq k\|V_1 - V_2\|$. We consider the metric $\|V_1 - V_2\|_\infty =$

$\max_{b \in \Delta(S)} |V_1(b) - V_2(b)|$ corresponding to the l_∞ . First, we focus on a single belief point and identify a criterion which ensures that $|HV_1(b) - HV_2(b)| \leq k|V_1(b) - V_2(b)|$. While somewhat technical, this criterion will enable us to demonstrate the contractivity of H . Moreover, it will also be useful in Section 9.3 to prove the correctness of the HSVI algorithm proposed therein.

Lemma 7.5. *Let $V, W : \Delta(S) \rightarrow \mathbb{R}$ be two convex continuous value functions and $b \in \Delta(S)$ a belief such that $[HV](b) \leq [HW](b)$. Let (π_1^V, π_2^V) and (π_1^W, π_2^W) be Nash equilibrium strategy profiles in stage games $[HV](b)$ and $[HW](b)$, respectively, and $C \geq 0$. If $W(\tau(b, a_1, o, \pi_2^V)) - V(\tau(b, a_1, o, \pi_2^V)) \leq C$ for every action $a_1 \in \text{Supp}(\pi_1^W)$ of player 1 and every observation $o \in O$ such that $\mathbb{P}_{b, \pi_1^W, \pi_2^V}[o | a_1] > 0$, then $[HW](b) - [HV](b) \leq \gamma C$.*

Lemma 7.6. *Operator H is a contraction on the space of convex continuous functions $V : \Delta(S) \rightarrow \mathbb{R}$ (under the supremum norm), with contraction-factor γ .*

Proof. Let $V, W : \Delta(S) \rightarrow \mathbb{R}$ be convex functions such that $\|V - W\|_\infty = \max_{b \in \Delta(S)} |V(b) - W(b)| \leq C$. To prove the contractivity of H , it suffices to show that $\|HV - HW\|_\infty \leq \gamma C$, i.e., $|[HV](b) - [HW](b)| \leq \gamma C$ for every belief $b \in \Delta(S)$. Since $|V(b) - W(b)| \leq C$ holds for every belief b , Lemma 7.5 yields both $HV(b) - HW(b) \leq \gamma C$ and $HW(b) - HV(b) \leq \gamma C$. \square

Next, we show that the optimal value function from Definition 5.2 is the fixpoint of the Bellman's operator H . Intuitively, this holds because V^* can be represented as a supremum over all possible value functions val^{σ^1} , which remains unchanged as we apply the operator H (resp. the value-compositions it consists of).

Lemma 7.7. *Lemma The optimal value function V^* satisfies $V^* = HV^*$.*

Together, the two results ensure that H can be applied iteratively to obtain V^* :

Theorem 2. *V^* is a unique fixpoint of H . Moreover, for any convex function V_0 , the sequence $\{V_i\}_{i=0}^\infty$ such that $V_i = HV_{i-1}$ converges to V^* .*

Proof. By Lemma 7.7, V^* is a fixpoint of H . By Lemma 7.6, H is a contraction mapping on the space of convex value functions. Banach's fixed point theorem [16] then implies the uniqueness and the “moreover” part. \square

8. Exact Value Iteration

In Section 7, we have shown that the optimal value function can be approximated by means of composing strategies in the sense of max-composition introduced in Definition 7.1. In this section, we provide a linear programming formulation to perform such optimal composition for value functions that are piecewise linear and convex, i.e., can be represented as a point-wise maximum of a finite set Γ of linear functions. Furthermore, we show that as long as the value function V is piecewise linear and convex, HV is also piecewise linear and convex. This allows for using the same linear program (LP) iteratively to approximate the optimal value function V^* by means of constructing a sequence of piecewise linear and convex value functions $\{V_i\}_{i=1}^\infty$ such that $V_i = HV_{i-1}$.

8.1. Computing Max-Compositions

In order to compute HV given a piecewise linear and convex (PWLC) value function V , it is essential to solve Equation (23). Every PWLC value function can be represented as a point-wise maximum over a finite set of linear functions $\{\alpha_1, \dots, \alpha_k\}$ (see Definition 5.10). Without loss of generality, we consider that the set Γ used to represent the value function V is the convex hull of the aforementioned set:

$$\Gamma := \text{Conv}(\{\alpha_1, \dots, \alpha_k\}) = \left\{ \sum_{i=1}^k \lambda_i \alpha_i \mid \lambda \in \mathbb{R}_{\geq 0}^k, \|\lambda\|_1 = 1 \right\}. \quad (26)$$

Recall that forming a convex hull of the set of linear functions used to represent V does not affect the values V attains (by Proposition 5.11). We will now show that when the set Γ is represented as in Equation (26), linear programming can be used to compute $HV(b)$:

Lemma 8.1. *Let $\Gamma = \text{Conv}(\{\alpha_1, \dots, \alpha_k\})$ be a convex hull of a finite set of α -vectors. Then $[HV](b)$ coincides with the solution of the following linear program:*

$$\max_{\pi_1, \lambda, \hat{\alpha}, V} \sum_{s \in S} b(s) \cdot V(s) \quad (27a)$$

$$\begin{aligned} s.t. \quad V(s) \leq & \sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2) + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} T(o, s' \mid s, a_1, a_2) \hat{\alpha}^{a_1, o}(s') \\ & \forall (s, a_2) \in S \times A_2 \end{aligned} \quad (27b)$$

$$\hat{\alpha}^{a_1, o}(s') = \sum_{i=1}^k \hat{\lambda}_i^{a_1, o} \cdot \alpha_i(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (27c)$$

$$\sum_{i=1}^k \hat{\lambda}_i^{a_1, o} = \pi_1(a_1) \quad \forall (a_1, o) \in A_1 \times O \quad (27d)$$

$$\sum_{a_1 \in A_1} \pi_1(a_1) = 1 \quad (27e)$$

$$\pi_1(a_1) \geq 0 \quad \forall a_1 \in A_1 \quad (27f)$$

$$\hat{\lambda}_i^{a_1, o} \geq 0 \quad \forall (a_1, o) \in A_1 \times O, 1 \leq i \leq k \quad (27g)$$

In the latter text, we also use the following dual formulation of the linear

program (27) (with some minor modifications to improve readability):

$$\min_{V, \pi_2, \hat{\tau}} V \quad (28a)$$

$$\text{s.t. } V \geq \sum_{(s, a_2) \in S \times A_2} \pi_2(s \wedge a_2) R(s, a_1, a_2) + \gamma \sum_{o \in O} \hat{V}(a_1, o) \quad \forall a_1 \quad (28b)$$

$$\hat{V}(a_1, o) \geq \sum_{s' \in S} \hat{\tau}(b, a_1, o, \pi_2)(s') \cdot \alpha_i(s') \quad \forall (a_1, o), 1 \leq i \leq k \quad (28c)$$

$$\hat{\tau}(b, a_1, \pi_2, o)(s') = \sum_{(s, a_2) \in S \times A_2} T(o, s' | s, a_1, a_2) \pi_2(s \wedge a_2) \quad \forall (a_1, o, s') \quad (28d)$$

$$\sum_{a_2 \in A_2} \pi_2(s \wedge a_2) = b(s) \quad \forall s \quad (28e)$$

$$\pi_2(s \wedge a_2) \geq 0 \quad \forall (s, a_2) \quad (28f)$$

Here, the stage strategy of player 2 is represented as a joint probability $\pi_2(s \wedge a_2)$ of playing action $a_2 \in A_2$ while being in state $s \in S$ (i.e., $\pi_2(a_2 | s) = \pi_2(s \wedge a_2)/b(s)$ where applicable). Player 1 then seeks the best response $a_1 \in A_1$ (constraint (28b)) that maximizes the sum of the expected immediate reward and the γ -discounted utility in the (a_1, o) -subgames. The beliefs $\tau(b, a_1, \pi_2, o)$ in the subgames are multiplied by the probability of reaching the (a_1, o) -subgame (i.e., there is no division by $\mathbb{P}_{b, a_1, \pi_2}[a_1, o]$ in Equation (28d)), hence also the values of subgames $V(a_1, o)$ need not be multiplied by $\mathbb{P}_{b, a_1, \pi_2}[a_1, o]$. The value of an (a_1, o) -subgame $V(a_1, o)$ is expressed as a maximum $\max_{\alpha \in \Gamma} \alpha(\tau(b, a_1, \pi_2, o))$ expressed by constraints (28c).

8.2. Value Iteration

To run a value iteration algorithm that would apply the linear program (27) repeatedly, we require that every V_i in the sequence $\{V_i\}_{i=0}^{\infty}$, starting from an arbitrary PWLC value function V_0 , is also piecewise linear and convex. By Lemma 8.3 this is always the case.

Lemma 8.2. *Let Q be the set of vertices of the polytope defined by constraints (27b)-(27g), and let $(\pi_1^q, \hat{\alpha}^q)$ be the assignment of the variables π_1 and $\hat{\alpha}$ corresponding to the vertex $q \in Q$. Then⁸*

$$[HV](b) = \max_{q \in Q} \text{valcomp}(\pi_1^q, \bar{\alpha}^q) \quad \text{for } \bar{\alpha}^q(a_1, o) = \hat{\alpha}^q(a_1, o)/\pi_1^q(a_1) . \quad (29)$$

⁸Note that $\bar{\alpha}^q(a_1, o)$ for a_1 with $\pi_1^q(a_1) = 0$ do not contribute to $\text{valcomp}(\pi_1^q, \bar{\alpha}^q)$. In parts of the game that are not reached by player 1, we can thus define $\bar{\alpha}^q$ arbitrarily.

Proof. Consider the LP (27) which computes the optimal value composition $\text{valcomp}(\pi_1, \bar{\alpha})$ in $[HV](b)$ (see Lemma 8.1). The polytope of feasible solutions of the LP defined by the constraints (27b)–(27g) is independent of the belief b (which only appears in the objective (27a)). Therefore, the set Q of vertices of this polytope is also independent of belief $b \in \Delta(S)$. The optimal solution of a linear programming problem (27) representing $[HV](b)$ can be found within the vertices Q of the polytope of feasible solutions [48]. There is a finite number of vertices $q \in Q$, and each vertex $q \in Q$ corresponds to some assignment of variables defining the value composition $\text{valcomp}(\pi_1^q, \bar{\alpha}^q)$. Since the set Q of the vertices of the polytope is independent of the belief b , we get the desired result. \square

Lemma 8.3. *If V is a piecewise linear and convex function, then so is HV .*

Proof. This lemma is a direct consequence of 25 8.2. Since the number of vertices of the polytope of LP (27) is finite, the pointwise maximization in (29) defines a PWLC function. \square

We can use the above-stated results to iteratively construct a sequence of value functions $\{V_i\}_{i=0}^\infty$ such that V_0 is an arbitrary PWLC function and $V_i = HV_{i-1}$. Namely, we construct V_i by enumerating the vertices of the polytope defined by the linear program (27) and constructing appropriate linear functions $\text{valcomp}(\pi_1^q, \bar{\alpha}^q)$. By 25 8.2, these linear functions form the set of α -vectors needed to represent a PWLC (Lemma 8.3) function V_i . According to Theorem 2 this sequence converges to V^* :

Corollary 8.4. *Starting from an arbitrary PWLC value function V_0 , a repeated application of the LP (27), as described in 25 8.2, converges to V^* .*

A more efficient algorithm can be devised based on, e.g., the linear support algorithm for POMDPs [14]. Here, the set Γ' of linear functions defining HV is constructed incrementally, terminating once it is provably large enough to represent the value function HV . Exact value iteration algorithms to solve POMDPs are, however, generally considered to only be capable of solving very small problems. We cannot, therefore, expect a decent performance of such approaches when solving one-sided POSGs that are more general than POMDPs. The next section remedies this issue by providing a point-based approach for solving one-sided POSGs

9. Heuristic Search Value Iteration for OS-POSGs

In this section, we provide a scalable algorithm for solving one-sided POSGs, inspired by the *heuristic search value iteration* (HSVI) algorithm [39, 40] for approximating value function of POMDPs (summarized in Section 3). Our algorithm approximates the convex optimal value function V^* using a pair of piecewise linear and convex value functions V_{LB}^Γ (lower bound on V^*) and V_{UB}^Γ (upper bound on V^*). These bounds are refined over time and, given the

initial belief b^{init} and the desired precision $\varepsilon > 0$, the algorithm is guaranteed to approximate the value $V^*(b^{\text{init}})$ within ε . In Section 10, we show that this process also generates value functions that allow us to extract ε -Nash equilibrium strategies of the game.

We first show the approximation schemes used to represent V_{LB}^Γ and V_{UB}^Υ , and the methods to initialize these bounds (Section 9.1). We then discuss the so-called “point-based updates” which are used to refine the bounds induced by V_{LB}^Γ and V_{UB}^Υ (Section 9.2). Finally, in Section 9.3, we describe the algorithm and prove its correctness.

9.1. Value Function Representations

Following the results on POMDPs and the original HSVI algorithm [19, 39], we use two distinct methods to represent upper and lower PWLC bounds on V^* .

Lower bound V_{LB}^Γ . Similarly as in the previous sections, the lower bound $V_{\text{LB}}^\Gamma : \Delta(S) \rightarrow \mathbb{R}$ is represented as a point-wise maximum over a finite set Γ of linear functions called α -vectors, i.e., $V_{\text{LB}}^\Gamma(b) = \max_{\alpha \in \Gamma} \alpha(b)$. Each $\alpha \in \Gamma$ is a linear function $\alpha : \Delta(S) \rightarrow \mathbb{R}$ represented by its values $\alpha(s)$ in the vertices of the $\Delta(S)$ simplex, i.e., $\alpha(b) = \sum_{s \in S} b(s) \cdot \alpha(s)$.

Upper bound V_{UB}^Υ . Upper bound $V_{\text{UB}}^\Upsilon : \Delta(S) \rightarrow \mathbb{R}$ is represented as a lower convex hull of a set of points $\Upsilon = \{(b_i, y_i) \mid 1 \leq i \leq k\}$. Each point $(b_i, y_i) \in \Upsilon$ provides an upper bound y_i on the value $V^*(b_i)$ in belief b_i , i.e., $y_i \geq V^*(b_i)$. Since the value function V^* is convex, it holds that

$$\left(\forall \lambda \in \mathbb{R}_{\geq 0}^k \text{ s.t. } \sum_{i=1}^k \lambda_i = 1 \right) : V^* \left(\sum_{i=1}^k \lambda_i b_i \right) \leq \sum_{i=1}^k \lambda_i \cdot V^*(b_i) \leq \sum_{i=1}^k \lambda_i \cdot y_i. \quad (30)$$

This fact is used in the first variant of the HSVI algorithm (HSV1 [39]) to obtain the value of the upper bound $V_{\text{HSV1}}^\Upsilon(b)$ for belief b : A linear program can be used to find coefficients $\lambda \in \mathbb{R}_{\geq 0}^k$ such that $b = \sum_{i=1}^k \lambda_i \cdot b_i$ holds and $\sum_{i=1}^k \lambda_i \cdot y_i$ is minimal:

$$V_{\text{HSV1}}^\Upsilon(b) = \min \left\{ \sum_{i=1}^k \lambda_i y_i \mid \lambda \in \mathbb{R}_{\geq 0}^k : \sum_{i=1}^k \lambda_i = 1 \wedge \sum_{i=1}^k \lambda_i b_i = b \right\}, \quad (31)$$

In the latter proof of the Theorem 3 showing the correctness of the algorithm, we require the bounds V_{LB}^Γ and V_{UB}^Υ to be δ -Lipschitz continuous. Since this needs not hold for V_{HSV1}^Υ , we define V_{UB}^Υ as a lower δ -Lipschitz envelope of V_{HSV1}^Υ :

$$V_{\text{UB}}^\Upsilon(b) = \min_{b' \in \Delta(S)} [V_{\text{HSV1}}^\Upsilon(b') + \delta \|b - b'\|_1]. \quad (32)$$

This computation can be expressed as a linear programming problem

$$V_{\text{UB}}^{\Upsilon}(b) = \min_{\lambda, \Delta, b'} \sum_{i=1}^k \lambda_i y_i + \delta \sum_{s \in S} \Delta_s \quad (33a)$$

$$\text{s.t. } \sum_{i=1}^k \lambda_i b_i(s) = b'(s) \quad \forall s \in S \quad (33b)$$

$$\Delta_s \geq b'(s) - b(s) \quad \forall s \in S \quad (33c)$$

$$\Delta_s \geq b(s) - b'(s) \quad \forall s \in S \quad (33d)$$

$$\sum_{i=1}^k \lambda_i = 1 \quad (33e)$$

$$\lambda_i \geq 0 \quad \forall 1 \leq i \leq k \quad (33f)$$

Here, we have $\Delta_s = |b'(s) - b(s)|$ (and hence $\sum_{s \in S} \Delta_s = \|b - b'\|_1$). Using the definitions of V_{UB}^{Υ} and $V_{\text{HSV11}}^{\Upsilon}$ together with the fact that V^* is δ -Lipschitz continuous and convex, we can prove that the function V_{UB}^{Υ} represents an upper bound on V^* :

Lemma 9.1. *Let $\Upsilon = \{(b_i, y_i) \mid 1 \leq i \leq k\}$ such that $y_i \geq V^*(b_i)$ for every $1 \leq i \leq k$. Then the value function V_{UB}^{Υ} is δ -Lipschitz continuous and satisfies*

$$V^* \leq V_{\text{UB}}^{\Upsilon} \leq V_{\text{HSV11}}^{\Upsilon}.$$

The dichotomy in representation of value functions V_{LB}^{Γ} and V_{UB}^{Υ} allows for easy refinement of the bounds. By adding new elements to the set Γ , the value $V_{\text{LB}}^{\Gamma}(b) = \max_{\alpha \in \Gamma} \alpha(b)$ can only increase—and hence the lower bound V_{LB}^{Γ} gets tighter. Similarly, by adding new elements to the set of points Υ , the solution of linear program (33) can only decrease and hence the upper bound V_{UB}^{Υ} tightens.

9.1.1. Initial Bounds

We now describe our approach to obtaining the initial bounds V_{LB}^{Γ} and V_{UB}^{Υ} on the optimal value function V^* of the game:

Lower bound V_{LB}^{Γ} . We initially set the lower bound to the value $\text{val}^{\sigma_1^{\text{unif}}}$ of the uniform strategy $\sigma_1^{\text{unif}} \in \Sigma_1$ of player 1 (i.e., the strategy that plays every action with probability $1/|A_1|$ in all stages of the game). Recall that the value $\text{val}^{\sigma_1^{\text{unif}}}$ of the strategy σ_1^{unif} is a linear function (see Lemma 5.4), and hence the initial lower bound V_{LB}^{Γ} is a piecewise linear and convex function represented as a pointwise maximum of the set $\Gamma = \{\text{val}^{\sigma_1^{\text{unif}}}\}$.

Upper bound V_{UB}^{Υ} . We use the solution of a perfect information variant of the game (i.e., where player 1 is assumed to know the entire history of the game, unlike in the original game). We form a modified game G' which is identical to the OS-POSG G (i.e., has the same states S , actions A_1 and A_2 , dynamics

T and rewards R), except that all information is revealed to player 1 in each step. G' is a perfect information stochastic game, and we can apply the value iteration algorithm to solve G' [9]. The additional information player 1 in G' (compared to G) can only increase the utility he can achieve. Hence V_s^* of the state s of game G' forms an upper bound on the utility player 1 can achieve in G if he knew that the initial state of the game is s (i.e., his belief is b_s where $b_s(s) = 1$). We initially define Υ as the set that contains one point for each state $s \in S$ of the game (i.e., for each vertex of the $\Delta(S)$ simplex),

$$\Upsilon = \{(b_s, V_s^*) \mid s \in S\} \quad b_s(s') = \begin{cases} 1 & s = s' \\ 0 & \text{otherwise} \end{cases} . \quad (34)$$

9.2. Point-based Updates

Unlike the exact value iteration algorithm (Section 8) which constructs all α -vectors needed to represent HV in each iteration, the HSVI algorithm focuses on a single belief at a time. Performing a *point-based* update in belief $b \in \Delta(S)$ corresponds to solving the stage-games $[HV_{\text{LB}}^\Gamma](b)$ and $[HV_{\text{UB}}^\Upsilon](b)$ where the values of subsequent stages are represented using value functions V_{LB}^Γ and V_{UB}^Υ , respectively.

Update of lower bound V_{LB}^Γ . First, the LP (27) is used to compute the optimal value composition $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ in $[HV_{\text{LB}}^\Gamma](b)$, i.e.,

$$(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}}) = \arg \max_{\substack{\pi_1 \in \Pi_1 \\ \bar{\alpha} \in \text{Conv}(\Gamma)^{\mathcal{A}_1 \times \mathcal{O}}}} \text{valcomp}(\pi_1, \bar{\alpha})(b) . \quad (35)$$

The $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ function is a linear function corresponding to a new α -vector that forms a lower bound on V^* . This new α -vector is used to refine the bound by setting $\Gamma := \Gamma \cup \{\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})\}$. As the following lemma shows, refining the lower bound V_{LB}^Γ via a point-based update preserves its desirable properties:

Lemma 9.2. *The lower bound V_{LB}^Γ initially satisfies the following conditions, which are subsequently preserved during point-based updates:*

- (1) V_{LB}^Γ is δ -Lipschitz continuous.
- (2) V_{LB}^Γ is lower bound on V^* .

Update of upper bound V_{UB}^Υ . Similarly to the case of the point-based update of the lower bound V_{LB}^Γ , the update of upper bound is performed by solving the stage game $[HV_{\text{UB}}^\Upsilon](b)$. Since V_{UB}^Υ is represented by a set of points Υ , it is not necessary to compute the optimal value composition. Instead, we form a refined upper bound $V_{\text{UB}}^{\Upsilon'}$ (which corresponds to V_{UB}^Υ after the point-based update is made) by adding a new point $(b, [HV_{\text{UB}}^\Upsilon](b))$ to the set Υ' representing $V_{\text{UB}}^{\Upsilon'}$, i.e., $\Upsilon' = \Upsilon \cup \{(b, [HV_{\text{UB}}^\Upsilon](b))\}$. We now show that the upper bound V_{UB}^Υ has the desired properties, and these properties are retained when applying the point-based update—and hence we can perform point-based updates of V_{UB}^Υ repeatedly.

Lemma 9.3. *The upper bound V_{UB}^Υ initially satisfies the following conditions, which are subsequently preserved during point-based updates:*

- (1) V_{UB}^Υ is δ -Lipschitz continuous.
- (2) V_{UB}^Υ is an upper bound on V^* .

The LPs (27) and (28) solve the stage game $[HV](b)$ when the value function V is represented as a maximum over a set of linear functions (i.e., the way lower bound V_{LB}^Γ is). It is, however, possible to adapt constraints in (28) to solve the $[HV_{\text{UB}}^\Upsilon](b)$ problem. We replace constraint (28c) by constraints inspired by the LP (33) used to solve $V_{\text{UB}}^\Upsilon(b)$.

$$\hat{V}(a_1, o) = \sum_{i=1}^{|\Upsilon|} \lambda_i^{a_1, o} y_i + \delta \sum_{s' \in S} \Delta_{a_1, o}^{s'} \quad \forall (a_1, o) \in A_1 \times O \quad (36a)$$

$$\sum_{i=1}^{|\Upsilon|} \lambda_{a_1, o}^i b_i(s') = b'_{a_1, o}(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (36b)$$

$$\Delta_{a_1, o}^{s'} \geq b'_{a_1, o}(s') - \hat{\tau}(b, a_1, \pi_2, o)(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (36c)$$

$$\Delta_{a_1, o}^{s'} \geq \hat{\tau}(b, a_1, \pi_2, o)(s') - b'_{a_1, o}(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (36d)$$

$$\sum_{i=1}^{|\Upsilon|} \lambda_i^{a_1, o} = \sum_{s' \in S} \hat{\tau}(b, a_1, \pi_2, o)(s') \quad \forall (a_1, o) \in A_1 \times O \quad (36e)$$

$$\lambda_{a_1, o}^i \geq 0 \quad \forall (a_1, o) \in A_1 \times O, 1 \leq i \leq |\Upsilon| \quad (36f)$$

9.3. The Algorithm

We are now ready to present the heuristic search value iteration (HSVI) algorithm for one-sided POSGs (Algorithm 4) and prove its correctness. The algorithm is similar to the HSVI algorithm for POMDPs [39, 40]. First, the bounds V_{LB}^Γ and V_{UB}^Υ on the optimal value function V^* are initialized (as described in Section 9.1) on line 1. Then, until the desired precision $V_{\text{UB}}^\Upsilon(b^{\text{init}}) - V_{\text{LB}}^\Gamma(b^{\text{init}}) \leq \varepsilon$ is reached, the algorithm performs a sequence of trials using the **Explore** procedure, starting from the initial belief b^{init} (lines 2–3).

The recursive procedure **Explore** generates a sequence of beliefs $\{b_i\}_{i=0}^k$ (for some $k \geq 0$) where $b_0 = b^{\text{init}}$ and each belief b_t reached at the recursion depth t satisfied $\text{excess}_t(b_t) > 0$ on line 2 or 10. The algorithm tries to ensure that values of beliefs b_t reached at t -th level of recursion (i.e., t -th stage of the game) are approximated with sufficient accuracy and the gap between $V_{\text{UB}}^\Upsilon(b)$ and $V_{\text{LB}}^\Gamma(b)$ is at most $\rho(t)$, where $\rho(t)$ is defined by

$$\rho(0) = \varepsilon \quad \rho(t+1) = [\rho(t) - 2\delta D]/\gamma. \quad (37)$$

To ensure that the sequence ρ is monotonically increasing and unbounded, we need to select the parameter D from the interval $(0, (1 - \gamma)\varepsilon/2\delta)$. When the approximation quality $V_{\text{UB}}^\Upsilon(b_t) - V_{\text{LB}}^\Gamma(b_t)$ of the value of a belief b_t reached at the t -th recursion level of **Explore** (i.e., at the $(t+1)$ -th stage of the game) exceeds

Algorithm 4: HSVI algorithm for one-sided POSGs

Data: Game G , initial belief b^{init} , discount factor $\gamma \in (0, 1)$, desired precision $\varepsilon > 0$, neighborhood parameter D

Result: Approximate value functions V_{LB}^Γ and V_{UB}^Υ satisfying $V_{\text{UB}}^\Upsilon(b) - V_{\text{LB}}^\Gamma(b) \leq \varepsilon$, sets Γ and Υ constructed by point-based updates that represent V_{LB}^Γ and V_{UB}^Υ

- 1 Initialize V_{LB}^Γ and V_{UB}^Υ (see Section 9.1)
- 2 **while** $\text{excess}_0(b^{\text{init}}) > 0$ **do**
- 3 **Explore**($b^{\text{init}}, 0$)
- 4 **return** V_{LB}^Γ and V_{UB}^Υ , sets Γ and Υ that represent V_{LB}^Γ and V_{UB}^Υ
- 5 **procedure** **Explore**(b_t, t)
- 6 $(\pi_1^{\text{LB}}, \pi_2^{\text{LB}}) \leftarrow$ equilibrium strategy profile in $[HV_{\text{LB}}^\Gamma](b_t)$
- 7 $(\pi_1^{\text{UB}}, \pi_2^{\text{UB}}) \leftarrow$ equilibrium strategy profile in $[HV_{\text{UB}}^\Upsilon](b_t)$
- 8 Perform point-based updates of V_{LB}^Γ and V_{UB}^Υ at belief b_t (see Section 9.2)
- 9 $(a_1^*, o^*) \leftarrow$ select according to forward exploration heuristic
- 10 **if** $\mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}}[a_1^*, o^*] \cdot \text{excess}_{t+1}(\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*)) > 0$ **then**
- 11 **Explore**($\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*), t + 1$)
- 12 Perform point-based updates of V_{LB}^Γ and V_{UB}^Υ at belief b_t (see Section 9.2)

the desired approximation quality $\rho(t)$, it is said to have a positive *excess gap* $\text{excess}_t(b_t)$,

$$\text{excess}_t(b_t) = V_{\text{UB}}^\Upsilon(b_t) - V_{\text{LB}}^\Gamma(b_t) - \rho(t) . \quad (38)$$

Note that our definition of excess gap is more strict compared to the original HSVI algorithm for POMDPs, where the $-2\delta D$ term from Equation (37) is absent (see Equation (5)). Unlike in POMDPs, which are single-agent, the belief transitions $\tau(b, a_1, \pi_2, o)$ in one-sided POSGs depend on player 2 as well (resp., on her strategy π_2). The tighter bounds on the approximation quality allow us to prove the correctness of the proposed algorithm in Theorem 3.

Forward exploration heuristic. The algorithm uses a heuristic approach to select which belief $\tau(b, a_1, \pi_2^{\text{LB}}, o)$ will be considered in the next recursion level of the **Explore** procedure, i.e., what action-observation pair $(a_1, o) \in A_1 \times O$ will be chosen by player 1, on line 9. This selection is motivated by Lemma 7.5—in order to ensure that $\text{excess}_t(b_t) \leq 0$ (or more precisely $\text{excess}_t(b_t) \leq -2\delta D$) at the currently considered belief b_t in t -th recursion level, all beliefs $\tau(b_t, a_1, \pi_2^{\text{LB}}, o)$ reached with positive probability when playing π_1^{UB} have to satisfy $\text{excess}_{t+1}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) \leq 0$. Specifically, we focus on a belief that has the highest *weighted excess gap*. Inspired by the original HSVI algorithm for POMDPs [39, 40]), we define the weighted excess gap as the excess gap $\text{excess}_{t+1}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o))$ multiplied by the probability that the action-

observation pair (a_1, o) that leads to the belief $\tau(b_t, a_1, \pi_2^{\text{LB}}, o)$ occurs. As a result, the next action-observation pair (a_1^*, o^*) for further exploration is selected according to the formula

$$(a_1^*, o^*) = \arg \max_{(a_1, o) \in A_1 \times O} \mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}} [a_1, o] \cdot \text{excess}_{t+1}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) . \quad (39)$$

We now show formally that if the weighted excess gap of the optimal (a_1^*, o^*) satisfies $\mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}} [a_1^*, o^*] \cdot \text{excess}_{t+1}(\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*)) \leq 0$, performing the point based update at b_t ensures that $\text{excess}_t(b_t) \leq -2\delta D$.

Lemma 9.4. *Let b_t be a belief encountered at t -th recursion level of **Explore** procedure and assume that the corresponding action-observation pair (a_1^*, o^*) (from line 9 of Algorithm 4) satisfies*

$$\mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}} [a_1^*, o^*] \cdot \text{excess}_{t+1}(\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*)) \leq 0 . \quad (40)$$

Then $\text{excess}_t(b_t) \leq -2\delta D$ after performing a point-based update at b_t . Furthermore, all beliefs $b'_t \in \Delta(S)$ such that $\|b_t - b'_t\|_1 \leq D$ satisfy $\text{excess}_t(b'_t) \leq 0$.

The proof goes by verifying the assumptions of Lemma 7.5 (“a criterion for contractivity”), which allows us to bound the difference between V_{UB}^Υ and V_{LB}^Γ by $\rho(t+1)$. The “furthermore” part then follows from δ -Lipschitz continuity of the bounds.

We now use Lemma 9.4 (especially its second part) to prove that Algorithm 4 terminates with $V_{\text{UB}}^\Upsilon(b^{\text{init}}) - V_{\text{LB}}^\Gamma(b^{\text{init}}) \leq \varepsilon$. As we mentioned earlier, we can also use value functions V_{LB}^Γ and V_{UB}^Υ to play the game and obtain ε -Nash equilibrium of the game (see Section 10).

Theorem 3. *For any $\varepsilon > 0$ and $0 < D < (1 - \gamma)\varepsilon/2\delta$, Algorithm 4 terminates with $V_{\text{UB}}^\Upsilon(b^{\text{init}}) - V_{\text{LB}}^\Gamma(b^{\text{init}}) \leq \varepsilon$.*

Proof. By the choice of parameter D , the sequence $\rho(t)$ (for $\rho(0) = \varepsilon$) is monotonically increasing and unbounded, and the difference between value functions V_{LB}^Γ and V_{UB}^Υ is bounded by $U - L$ (since $L \leq V_{\text{LB}}^\Gamma(b) \leq V_{\text{UB}}^\Upsilon(b) \leq U$ for every belief $b \in \Delta(S)$). Therefore, there exists T_{max} such that $\rho(T_{\text{max}}) \geq U - L \geq V_{\text{UB}}^\Upsilon(b) - V_{\text{LB}}^\Gamma(b)$ for every $b \in \Delta(S)$, so the recursive procedure **Explore** always terminates.

To prove that the whole algorithm terminates, we reason about sets $\Psi_t \subset \Delta(S)$ of belief points where the trials performed by the **Explore** terminated. Initially, $\Psi_t = \emptyset$ for every $0 \leq t < T_{\text{max}}$. Whenever the **Explore** recursion terminates at recursion level t (i.e., the condition on line 10 does not hold), the belief b_t (which was the last belief considered during the trial) is added into set Ψ_t ($\Psi_t := \Psi_t \cup \{b_t\}$). Recall that since $\Delta(S)$ is compact, it is, in particular, totally bounded (that is, if any two distinct elements b, b' of a set $\Psi_t \subset \Delta(S)$ satisfy $\|b - b'\|_1 > D$, the set Ψ_t must be finite). Since the number of possible termination depths is finite ($0 \leq t \leq T_{\text{max}}$), the algorithm has to terminate unless some Ψ_t is infinite. To show that the algorithm terminates, it thus remains that every two distinct points $b, b' \in \Psi_t$ are at least D apart.

Assume to the contradiction that two trials terminated at recursion level t with the last beliefs considered $b_t^{(1)}$ (for the earlier trial) and $b_t^{(2)}$ (for the trial that occurred at a later time), and that these beliefs satisfy $\|b_t^{(1)} - b_t^{(2)}\|_1 \leq D$. When the former trial has been terminated in belief $b_t^{(1)}$, all reachable beliefs from $b_t^{(1)}$ had a negative excess gap (otherwise the trial would have continued as the condition on line 10 would have been satisfied). According to Lemma 9.4, after the point-based update is performed in $b_t^{(1)}$, the excess gap of all beliefs b'_t with $\|b_t^{(1)} - b'_t\|_1 \leq D$ have negative excess gap $\text{excess}_t(b'_t) \leq 0$. When $b_t^{(2)}$ has been selected for exploration in $(t-1)$ -th level of recursion, the condition on line (10) was met and $b_t^{(2)}$ must have had positive excess gap $\text{excess}_t(b_t^{(2)}) > 0$. This, however, contradicts the assumption that all beliefs b'_t with $\|b_t^{(1)} - b'_t\|_1 \leq D$ (i.e., including $b_t^{(2)}$) already have negative excess gap.

Now that we know that Algorithm 4 always terminates, note that at least one trial must have terminated in the first level of recursion (unless the Algorithm 4 has terminated on line 2 with $\text{excess}_0(b^{\text{init}}) \leq 0$ beforehand). By Lemma 9.4, the update in b^{init} then renders $\text{excess}_0(b^{\text{init}}) \leq -2\delta D \leq 0$. We then have that $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}}) \leq \rho(0) = \varepsilon$ which completes the proof. \square

10. Using Value Function to Play

In the previous section, we have presented an algorithm that can approximate the value $V^*(b^{\text{init}})$ of the game within an arbitrary given precision $\varepsilon > 0$ starting from an arbitrary initial belief b^{init} . However, in many games, knowing only the game's value is not enough. Indeed, to solve the game, we also need access to strategies that achieve the desired near-optimal performance. In this section, we show that using the value functions V_{LB}^{Γ} and V_{UB}^{Υ} computed by the HSVI algorithm (Algorithm 4) enables us to obtain ε -Nash equilibrium strategies for both players.

The Bellman's equation from Theorem 1 may suggest that the near-optimal strategies can be extracted by employing the lookahead decision rule (similarly to POMDPs) and obtaining strategies to play in the current stage by computing the Nash equilibrium of stage games $[HV_{\text{LB}}^{\Gamma}](b)$ and $[HV_{\text{UB}}^{\Upsilon}](b)$, respectively. However, unlike in POMDPs and Markov games of imperfect information, this approach does *not* work in one-sided POSGs because the belief of player 1 does not constitute a sufficient statistic for playing the game. The reasons for this are similar to the usage of unsafe resolving [12, 36] in the realm of extensive-form games. We use the following example to demonstrate the insufficiency of the belief to play the game.

Example 10.1. Consider a *matching pennies* game shown in Figure 4a. This game can be formalized as a one-sided POSG that is shown in Figure 4b. The game starts in the state s_0 (i.e., the initial belief is $b^{\text{init}}(s_0) = 1$) and player 2 chooses her action H or T . Next, after transitioning to s_H or s_T (based on the decision of player 2), player 1 is *unaware* of the true state of the game (i.e., the past decision of player 1) and chooses his action H or T . Based on the

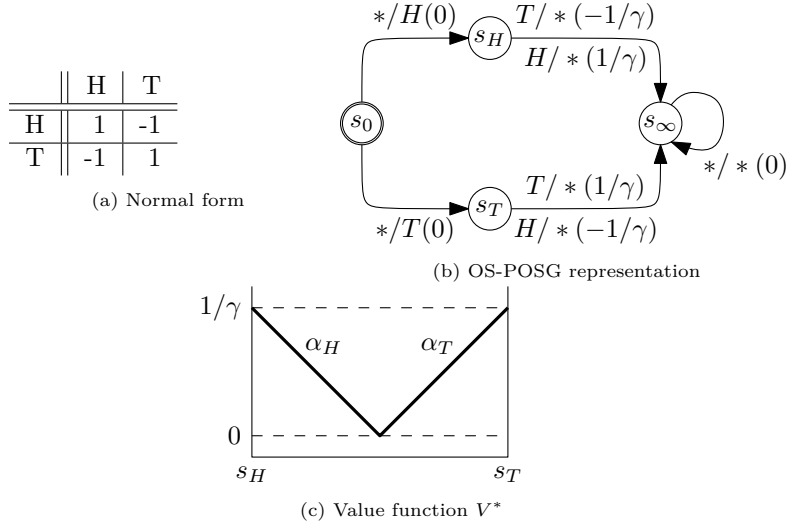


Figure 4: A game where belief is not a sufficient statistic for the imperfectly informed player.

combination of decisions taken by the players, player 1 gets either $1/\gamma$ or $-1/\gamma$ and the game transitions to the state s_∞ where it stays forever with zero future rewards.

To understand the caveats of using belief $b \in \Delta(S)$ to derive the stage strategy to play, let us consider the optimal value function V^* of the OS-POSG representation (Figure 4b) of the matching pennies game. Figure 4c shows the values of V^* over simplex $\Delta(\{s_H, s_T\})$. If it is more likely that the player 2 played H in the first stage of the game (i.e., the current state is s_H), it is optimal for player 1 to play strategy prescribing him to play H in the current stage (with value α_H). Conversely, if it is more likely that the current state is s_T , player 1 is better off with playing T (with value α_T). The value function V^* is then a point-wise maximum over these two linear functions.

Now, since the uniform mixture between H and T is the Nash equilibrium strategy for both players in the matching pennies game, player 1 will find himself in a situation when he assumes that the current belief is $\{s_H : 0.5, s_T : 0.5\}$. In this belief, any decision of player 1 yields expected reward 0—hence based purely on the belief, player 1 may opt to play, e.g., “always T ”. However, such strategy is not in equilibrium and player 2 is able to exploit it by playing “always H ”. This example illustrates that the belief alone does not provide sufficient information to choose the right strategy π_1 for the current stage based on the Equation (24b).

10.1. Justified Value Functions

First of all, we define conditions under which it makes sense to use value function to play a one-sided POSG. The conditions are similar to *uniform improvability* in, e.g., POMDPs. Our definitions, however, reflect the fact that

we deal with a two-player problem (and we thus introduce the condition for each player separately). Moreover, we use a stricter condition for player 1 who does *not* have perfect information about the belief—and thus defining the condition based solely on the beliefs is not sufficient.

Definition 10.2 (Min-justified value function). Convex continuous value function V is said to be *min-justified* (or, justified from the perspective of the minimizing player 2) if for every belief $b \in \Delta(S)$ it holds that $[HV](b) \leq V(b)$.

Definition 10.3 (Max-justified value function). Let Γ be a compact set of linear functions, and V be a value function such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$ for every b . V is said to be *max-justified* by Γ (or, justified from the perspective of the maximizing player 1) if for every $\alpha \in \Gamma$ there exists $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \Gamma^{A_1 \times O}$ such that $\text{valcomp}(\pi_1, \bar{\alpha}) \geq \alpha$.

While the reason for the terminology is not apparent just yet, we will show in Sections 10.2 that the “max-justifying” set Γ can be used to construct a strategy σ_1 of player 1 such that $\text{val}^{\sigma_1}(b) \geq V(b)$ for every b . Similarly, we will show in Section 10.3 that if the value function V is min-justified, we can construct a strategy σ_2 of player 2 that *justifies* the value $V(b)$ for every belief $b \in \Delta(S)$, i.e., we have $\mathbb{E}_{b, \sigma_1, \sigma_2}[\text{Disc}^\gamma] \leq V(b)$ against every strategy σ_1 of player 1.

As preparation for more substantial proofs that follow, the remainder of this subsection presents several basic properties of min- and max-justified functions.

Recall that no matter how well things go for the maximizing player, the corresponding utility will never get above U . Similarly, the minimizing player cannot push the utility below L . Lemma 10.5 and Lemma 10.4 prove that max- and min-justified functions obey the same restrictions. This is in agreement with our intuition that max-justification should guarantee utility of *at least* some value (which therefore cannot be higher than U) and min-justification should guarantee utility of *no more than* some value (which therefore cannot be lower than L).

Lemma 10.4. *Let V be a value function that is min-justified. Then $V(b) \geq L$.*

Lemma 10.5. *Let V be a value function that is max-justified by a set of α -vectors Γ . Then for every $\alpha \in \Gamma$ we have $\alpha \leq U$.*

To prepare for showing that the value function V_{LB}^Γ resulting produced by Algorithm 4 is max-justified by $\text{Conv}(\Gamma)$, we state the following technical lemma:

Lemma 10.6. *Let Γ be a set of linear functions, and V a value function that is max-justified by Γ . Then V is also max-justified by $\text{Conv}(\Gamma)$.*

10.2. Strategy of Player 1

In this section, we will show that when the value function V is max-justified by a set of α -vectors Γ , we can implicitly form a strategy σ_1 of player 1 that achieves utility of at least $V(b^{\text{init}})$ for any given initial belief b^{init} . We provide

an online game-playing algorithm (Algorithm 5) which implicitly constructs the desired strategy. This algorithm is inspired by the ideas of continual resolving for extensive-form games [32].

While playing the game, Algorithm 5 maintains a lower bound ρ on the values the reconstructed strategy has to achieve. Inspired by the terminology of continual resolving for extensive-form games, we call this lower-bounding linear function a *gadget*. The goal of the $\text{Act}(b, \rho)$ method is to reconstruct a strategy σ_1 of player such that its value satisfies $\text{val}^{\sigma_1} \geq \rho$. We will now show that the Act method achieves precisely this. The reasoning about the current gadget allows us to obtain guarantees on the quality of the reconstructed strategy, even when player 1 does not have an accurate belief because he does not have access to the stage strategies used by the adversary.

Algorithm 5: Continual resolving algorithm for one-sided POSGs

input : one-sided POSG G
a finite set Γ of linear functions representing convex value
function V

- 1 $b \leftarrow b^{\text{init}}$
- 2 $\rho^{\text{init}} \leftarrow \arg \max_{\alpha \in \Gamma} \alpha(b^{\text{init}})$
- 3 $\text{Act}(b^{\text{init}}, \rho^{\text{init}})$
- 4 **procedure** $\text{Act}(b, \rho)$
- 5 $(\pi_1^*, \bar{\alpha}^*) \leftarrow \arg \max_{\pi_1, \bar{\alpha}} \{ \text{valcomp}(\pi_1, \bar{\alpha})(b) \mid \pi_1 \in \Pi_1, \bar{\alpha} \in \text{Conv}(\Gamma)^{A_1 \times O} \text{ s.t. } \text{valcomp}(\pi_1, \bar{\alpha}) \geq \rho \}$
- 6 $\pi_2 \leftarrow \text{solve } [HV](b) \text{ to obtain assumed stage strategy of the adversary}$
- 7 sample and play $a_1 \sim \pi_1^*$
- 8 $o \leftarrow \text{observed observation}$
- 9 $b' \leftarrow \tau(b, a_1, \pi_2, o)$
- 10 $\text{Act}(b', \alpha_{a_1, o}^*)$

Proposition 10.7. *Let V be a value function that is max-justified by a set of α -vectors Γ . Let $b^{\text{init}} \in \Delta(S)$ and $\rho^{\text{init}} \in \Gamma$. By playing according to $\text{Act}(b^{\text{init}}, \rho^{\text{init}})$, player 1 implicitly forms a strategy σ_1 for which $\text{val}^{\sigma_1} \geq \rho^{\text{init}}$.*

This proposition is proven by constructing a sequence of strategies under which player 1 follows Algorithm 5 for K steps (for $K = 0, 1, \dots$). We provide a lower bound on the value each of these strategies, and show that the limit of these lower bounds coincides with ρ^{init} , as well as with the lower bound on the value guaranteed by following Algorithm 5 for *infinite* period of time.

Corollary 10.8. *Let V be a value function that is max-justified by a compact set Γ and let b^{init} be the initial belief of the game. The Algorithm 5 implicitly constructs a strategy σ_1 which guarantees that the utility to player 1 will be at least $V(b^{\text{init}})$.*

Proof. ρ^{init} from line 2 of Algorithm 5 has value $\rho^{\text{init}}(b^{\text{init}}) = V(b^{\text{init}})$ in the initial belief b^{init} . By Proposition 10.7, we can construct a strategy σ_1 with value $\text{val}^{\sigma_1} \geq \rho^{\text{init}}$. Hence $\text{val}^{\sigma_1}(b^{\text{init}}) \geq \rho^{\text{init}}(b^{\text{init}}) = V(b^{\text{init}})$. \square

10.3. Strategy of Player 2

We will now present an analogous algorithm to obtain a strategy for player 2 when the value function V is min-justified. Recall that the stage strategies π_2 of player 2 influence the belief of player 1 (Equation 11). Unlike player 1, player 2 knows which stage strategies π_2 have been used in the past, and he is thus able to infer the current belief of player 1. As a result, the **Act** method of Algorithm 6 depends on the current belief of player 1, but not on the gadget ρ as it did in Algorithm 5.

Algorithm 6: Strategy of player 2

input : one-sided POSG G
convex value function V

- 1 **Act**(b^{init})
- 2 **procedure** **Act**(b)
- 3 $\pi_2^* \leftarrow$ optimal strategy of player 2 in the stage game $[HV](b)$
- 4 $s \leftarrow$ currently observed state
- 5 sample and play $a_2 \sim \pi_2^*(\cdot | s)$
- 6 $(a_1, o) \leftarrow$ action of the adversary and the corresponding observation
- 7 **Act**($\tau(b, a_1, \pi_2^*, o)$)

We will now show that if the value function V is min-justified, playing according to Algorithm 6 guarantees that the utility will be at most⁹ $V(b^{\text{init}})$.

Proposition 10.9. *Let V be a min-justified value function and let b^{init} be the initial belief of the game. The Algorithm 6 implicitly constructs a strategy σ_2 which guarantees that the utility to player 1 will be at most $V(b^{\text{init}})$.*

The proof of Proposition 10.9 is similar to the proof of Proposition 10.7. We derive an upper bound on the utility player 1 can achieve against player 2 who follows Algorithm 6 for K steps (for $K = 0, 1, \dots$). We show that the limit of these upper bounds coincides with $V(b^{\text{init}})$ and with the upper bound on the utility player 1 can achieve when player 2 follows 6 for an infinite number of iterations.

10.4. Using Value Functions V_{LB}^Γ and V_{UB}^Υ to Play the Game

In Sections 10.2 and 10.3, we have shown that we can obtain strategies to play the game when the value functions are max-justified or min-justified, respectively. In this section, we will show that the heuristic search value iteration

⁹In other words, this is a performance guarantee for the (minimizing) player 2.

algorithm for solving one-sided POSGs (Section 9) generates value functions with these properties. Namely, at any time, the lower bound V_{LB}^Γ is max-justified value function by the set of α -vectors $\text{Conv}(\Gamma)$, and the upper bound V_{UB}^Υ is min-justified.

This allows us to derive two important properties of the algorithm. First, since Theorem 3 guarantees that the algorithm terminates with $V_{\text{UB}}^\Upsilon(b^{\text{init}}) - V_{\text{LB}}^\Gamma(b^{\text{init}}) \leq \varepsilon$, we can use the resulting value functions V_{LB}^Γ (represented by Γ) and V_{UB}^Υ to obtain ε -Nash equilibrium strategies for both players. Next, we can also run the algorithm in anytime fashion and, since the bounds V_{LB}^Γ and V_{UB}^Υ satisfy the properties at any point of time, use these bounds to extract strategies with performance guarantees.

We will first prove that at any point of time in the execution of Algorithm 4, the lower bound V_{LB}^Γ is max-justified by the set $\text{Conv}(\Gamma)$, and the upper bound V_{UB}^Υ is a min-justified value function. To prove this, it suffices to show initial lower-bound value function V_{LB}^Γ is max-justified by $\text{Conv}(\Gamma)$ and the initial upper-bound value function V_{UB}^Υ is min-justified, and that this property is preserved after any sequence of point-based updates performed on V_{LB}^Γ and V_{UB}^Υ . With the help of Lemma 10.6, we can prove that this is true for V_{LB}^Γ :

Lemma 10.10. *Let Γ be the set of α -vectors that have been generated at any time during the execution of the HSVI algorithm for one-sided POSGs (Algorithm 4). Then the lower bound V_{LB}^Γ is max-justified by the set $\text{Conv}(\Gamma)$.*

Even though the proof is more complicated, the analogous result holds for V_{UB}^Υ as well:

Lemma 10.11. *Let V_{UB}^Υ be the upper bound considered at any time of the execution of the HSVI algorithm for one-sided POSGs (Algorithm 4). Then V_{UB}^Υ is min-justified.*

Proof. Upper bound V_{UB}^Υ is only modified by means of point-based update on lines 8 and 12 of Algorithm 4. Therefore, it suffices to show that (1) the initial upper bound is min-justified and that (2) the upper bound $V_{\text{UB}}^{\Upsilon'}$ resulting from applying a point-based update on a min-justified upper bound V_{UB}^Υ is min-justified as well.

First, let us prove that the initial value function V_{UB}^Υ is min-justified. Initially, $V_{\text{UB}}^\Upsilon(b)$ is set to the value of a *perfect information* version of the game, where the imperfectly informed player 1 gets to know the initial state of the game. By removing this information from player 1, the utility player 1 can achieve can only decrease. It follows that $[HV_{\text{UB}}^\Upsilon](b) \leq V_{\text{UB}}^\Upsilon(b)$, so the initial value function $V_{\text{UB}}^\Upsilon(b)$ is min-justified.

Now, let us consider an upper bound V_{UB}^Υ represented by a set $\Upsilon = \{(b_i, y_i) \mid 1 \leq i \leq k\}$ that is considered by the Algorithm 4 and let us assume that V_{UB}^Υ is min-justified. Consider that a point-based update in b_{k+1} is to be performed. We show that the function $V_{\text{UB}}^{\Upsilon'}$ resulting from the point-based update in b_{k+1} is min-justified as well. Recall that $\Upsilon' = \Upsilon \cup \{(b_{k+1}, y_{k+1})\}$ and $y_{k+1} = [HV_{\text{UB}}^\Upsilon](b_{k+1})$. Clearly, since $\Upsilon \subset \Upsilon'$, it holds $V_{\text{UB}}^{\Upsilon'}(b) \leq V_{\text{UB}}^\Upsilon(b)$ and $[HV_{\text{UB}}^{\Upsilon'}](b) \leq [HV_{\text{UB}}^\Upsilon](b)$

for every $b \in \Delta(S)$. Due to this and since V_{UB}^{Υ} is assumed to be min-justified, we have $y_i \geq [HV_{\text{UB}}^{\Upsilon'}](b)$ for every $1 \leq i \leq k+1$. We will now prove that $V_{\text{UB}}^{\Upsilon'}$ is min-justified by showing that $[HV_{\text{UB}}^{\Upsilon'}](b) \leq V_{\text{UB}}^{\Upsilon'}(b)$ holds for arbitrary belief $b \in \Delta$. Let λ_i and b' correspond to the optimal solution of the linear program (33) for solving $V_{\text{UB}}^{\Upsilon'}(b)$. We have

$$\begin{aligned}
V_{\text{UB}}^{\Upsilon'}(b) &= \sum_{i=1}^{k+1} \lambda_i y_i + \delta \|b - b'\|_1 \\
&\quad \lambda_i \text{ and } b' \text{ represent an optimal solution of } V_{\text{UB}}^{\Upsilon'}(b) \\
&\geq \sum_{i=1}^{|\Upsilon|} \lambda_i \cdot [HV_{\text{UB}}^{\Upsilon'}](b_i) + \delta \|b - b'\|_1 \\
&\geq [HV_{\text{UB}}^{\Upsilon'}](b') + \delta \|b - b'\|_1 \\
&\quad HV_{\text{UB}}^{\Upsilon'} \text{ is convex, see Proposition 7.2} \\
&\geq [HV_{\text{UB}}^{\Upsilon'}](b) \quad V_{\text{UB}}^{\Upsilon'} \text{ is } \delta\text{-Lipschitz continuous, and hence,} \\
&\quad \text{by Proposition 7.2, } HV_{\text{UB}}^{\Upsilon'} \text{ is as well.}
\end{aligned}$$

This shows that any point-based update results in a min-justified value function $V_{\text{UB}}^{\Upsilon'}$. As a result, Algorithm 4 only considers upper bounds V_{UB}^{Υ} that are min-justified. \square

We are now in a position to show that Algorithm 4 produces ε -Nash equilibrium strategies.

Theorem 4. *In any OS-POSG, applying Algorithms 5 and 6 to the output of Algorithm 4 yields an ε -Nash equilibrium.*

Proof. According to Theorem 3, Algorithm 4 terminates and the value functions V_{LB}^{Γ} and V_{UB}^{Υ} that result from the execution of the algorithm satisfy $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}}) \leq \varepsilon$. Furthermore, we know that lower bound V_{LB}^{Γ} is max-justified by the set Γ resulting from the execution of Algorithm 4 (38 10.10), and the upper bound V_{UB}^{Υ} is min-justified (Lemma 10.11). We can therefore use Algorithm 5 to obtain a strategy for player 1 that achieves utility of at least $V_{\text{LB}}^{\Gamma}(b^{\text{init}})$ for player 1 (Corollary 10.8). Similarly, we can use Algorithm 6 to obtain a strategy for player 2 that ensures that the utility of player 1 will be at most $V_{\text{UB}}^{\Upsilon}(b^{\text{init}})$ (Proposition 10.9). It follows that if either player were to deviate from the strategy prescribed by the algorithm, they would not be able to improve their utility by more than $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}})$. Since $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}}) \leq \varepsilon$, these strategies must form a ε -Nash equilibrium of the game. \square

11. Experimental evaluation

In this section, we focus on the experimental evaluation of the heuristic search value iteration algorithm for solving one-sided partially observable stochastic

games from Section 9. We demonstrate the scalability of the algorithm in three security domains. Rewards in all of the domains have been scaled to the interval $[0, 100]$ or $[-100, 0]$, respectively, and we report the runtime required to reach $V_{\text{UB}}^{\Upsilon}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}}) \leq 1$. We first outline the details of our experimental setup.

11.1. Algorithm Settings

Compared to the version of the HSVI algorithm presented in Section 9, we adopt several modifications to improve the scalability of the algorithm. In this section, we describe these modifications and show that the theoretical guarantees of the algorithm still hold.

Pruning the Sets Γ and Υ . Each time a point-based update is performed, the size of the sets Γ and Υ used to represent value functions V_{LB}^{Γ} and V_{UB}^{Υ} increases. As new elements are generated, some of the elements in these sets may become unnecessary to accurately represent the bounds V_{LB}^{Γ} and V_{UB}^{Υ} . Since the sizes of sets V_{LB}^{Γ} and V_{UB}^{Υ} have a direct impact on the sizes of linear programs used throughout the algorithm, removing unnecessary elements from V_{LB}^{Γ} and V_{UB}^{Υ} improves the performance. Whenever a new α -vector $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ is generated according to Equation (35), all dominated elements in the set Γ get removed and only those elements of $\alpha \in \Gamma$ that dominate $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ in at least one state remain, i.e.,

$$\Gamma := \left\{ \alpha' \in \Gamma \mid \exists s \in S : \alpha'(s) > \text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})(s) \right\} \cup \left\{ \text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}}) \right\}. \quad (41)$$

For the set Υ used to represent the upper bound V_{UB}^{Υ} , we use a batch approach instead of removing dominated elements immediately. We remove dominated elements every time the size of the set Υ increases by 10% compared to the size after the last pruning was performed (this is analogous to the pruning technique proposed in [39]). Algorithm 7 inspects each point $(b_i, y_i) \in \Upsilon$ and checks whether it is needed to represent value function V_{UB}^{Υ} —and if it is not needed, the point gets removed.

Algorithm 7: Pruning set Υ representing the upper bound V_{UB}^{Υ}

input : Set Υ used to represent V_{UB}^{Υ}
1 for $(b_i, y_i) \in \Upsilon$ **do**
2 **if** $y_i > V_{\text{UB}}^{\Upsilon}(b_i)$ **then** $\Upsilon := \Upsilon \setminus \{(b_i, y_i)\}$

Removing elements from sets Γ and Υ does not violate the theoretical properties of the algorithm. First of all, only elements that are not necessary to represent currently considered bounds are removed—hence the values of value functions V_{LB}^{Γ} and V_{UB}^{Υ} considered at each step of the algorithm remain unchanged, and the convergence property is hence retained. Furthermore, we can still use pruned value functions to extract strategies with guaranteed

performance. Since the resulting upper bound value function V_{UB}^{Γ} is identical to the one obtained without pruning, it is still min-justified. It can thus be used to obtain a strategy of the minimizing player 2 with guaranteed utility at most $V_{\text{UB}}^{\Gamma}(b^{\text{init}})$ (Section 10.3). Similarly, V_{LB}^{Γ} can be used to obtain a strategy of player 1 (Section 10.2). Despite the fact that the resulting set Γ of α -vectors is different from the set constructed by Algorithm 4 when no pruning is used, we can see that for every missing element α' there has to exist an element α such that $\alpha \geq \alpha'$ (see Equation (41)). Therefore, we can always replace missing α -vectors in value compositions (i.e., linear functions $\alpha^{a_1, o}$) without decreasing the values of the resulting value composition—and hence V_{LB}^{Γ} remains max-justified by the set of α -vectors $\text{Conv}(\Gamma)$.

Partitioning States and Value Functions. In many games, even the imperfectly informed player 1 has access to some information about the game. For example, in the pursuit-evasion games we discuss below, the pursuer *knows* his position—and representing his uncertainty about his position within the belief is unnecessary. To reduce the dimension of the beliefs, we allow for partitioning states into disjoint sets such that the imperfectly informed player 1 always *knows* in which set he is currently. Formally, let $S = \bigcup_{i=1}^K S_i$ such that $S_i \cap S_j = \emptyset$ for every $i \neq j$. Player 1 has to know the initial partition, i.e., $\text{Supp}(b^{\text{init}}) \subseteq S_i$ for some $1 \leq i \leq K$. Furthermore, he has to be able to infer which partition he is in at any time, i.e., for every belief b over a partition S_i (i.e., $\text{Supp}(b) \subseteq S_i$), every achievable action-observation pair (a_1, o) and every stage strategy $\pi_2 \in \Pi_2$ of player 2, we have $\text{Supp}(\tau(b, a_1, \pi_2, o)) \subseteq S_j$ for some $1 \leq j \leq K$. We use $T(S_i, a_1, o)$ to denote such S_j .

This partitioning allows for reducing the size of LP (27) used to compute stage game solutions. Namely, the quantification over $s \in S$ can be replaced by $s \in S_i$, where S_i is the current partition. Furthermore, since also the partition of the next stage has to be known, we can also replace $(a_1, o, s') \in A_1 \times O \times S$ by $(a_1, o, s') \in A_1 \times O \times T(S_i, a_1, o)$.

Parameters and Hardware. We use value iteration for stochastic games, or MDPs, respectively, to initialize the upper and lower bounds. The upper bound is initialized by solving a perfect-information variant of the game (see Section 9.1). The lower bound is computed by fixing the uniform strategy σ_1^{unif} for player 1 and solving the resulting Markov decision process from the perspective of player 2. We terminate the algorithms when either change in valuations between iterations of value iteration is lower than 0.025, or 20 minutes time limit has expired. The initialization time is included in the computation times of the HSVI algorithm.

We use $\varepsilon = 1$. However, similarly to [39], we adjust ε in each iteration, and we get ε_{imm} that is about to be used in the current iteration using formula $\varepsilon_{\text{imm}} = 0.25 + \eta(V_{\text{UB}}^{\Gamma}(b^{\text{init}}) - V_{\text{LB}}^{\Gamma}(b^{\text{init}}) - 0.25)$ with $\eta = 0.9$. We set the parameter D to the largest value such that $\rho(t) \geq 0.25^{-t}$ holds for every $t \geq 0$.

Each experiment has been run on a single core of Intel Xeon Platinum 8160. We have used CPLEX 12.9 to solve the linear programs.

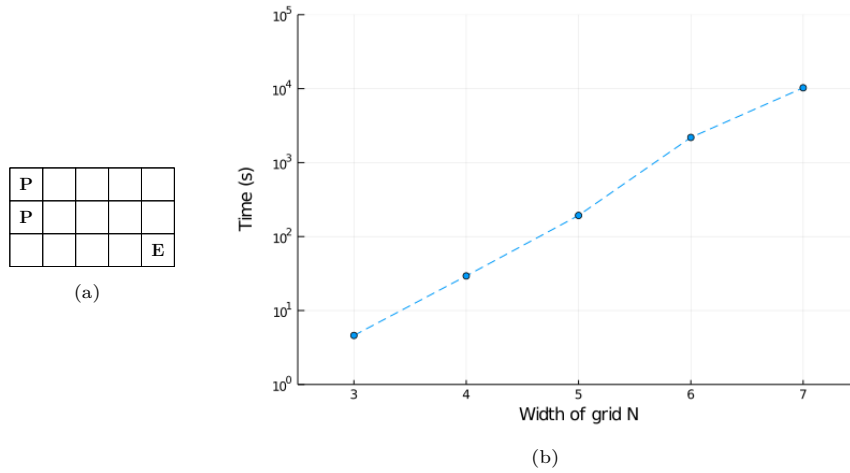


Figure 5: Pursuit evasion games: (a) Pursuit evasion game $5 \times N$. The **P**s denote the initial positions of the pursuers, the **E** denotes the initial position of the evader. (b) Time to reach gap 1 for different grid-widths (N).

11.2. Experimental Results

We now turn our attention to the discussion of experimental results. We introduce the domains used in our experiments and comment on the scalability of the proposed algorithm.

Pursuit-Evasion Games (inspired by [15, 24]). In pursuit-evasion games, a team of K centrally controlled pursuers (we consider a team of $K = 2$) is trying to locate and capture the evader—who is trying to avoid getting captured. The game is played on a grid (dimensions $3 \times N$), with the pursuers starting in the top-left corner and the evader in the bottom-right corner – see Figure 5a. In each step, the units move to one of their adjacent locations (i.e., the actions of the evader are $A_2 = \{\text{left, right, up, down}\}$, while the actions available to the team of pursuers are joint actions for all units in the team, $A_1 = (A_2)^K$). The game ends when one of the units from the team of pursuers enters the same cell as the evader—and the team of pursuers (player 1) then receives a reward of +100. The reward for all other transitions in the game is zero. The pursuer knows the location of their units, but the current location of the evader is not known.

The game with $N = 3$ was solved in 4.5 s on average, while the game with $N = 7$ took 10 267 s to be solved to the gap $\varepsilon = 1$ – full results can be found in Figure 5b. The game $8 \times N$ has not been solved successfully within 10 hours time limit, and the gap of $V_{UB}^\Upsilon(b_{\text{init}}) - V_{LB}^\Gamma(b_{\text{init}}) = 1.245$ has been reached after 10 hours. Sizes of the games range from 143 states and 2671 transitions (for $3 \times N$ game) to 3171 states and 92531 transitions (for $8 \times N$ game).

Search Games (inspired by [8]). In search games that model intrusion, the defender patrols checkpoint zones (see Figure 6a, the zones are marked with

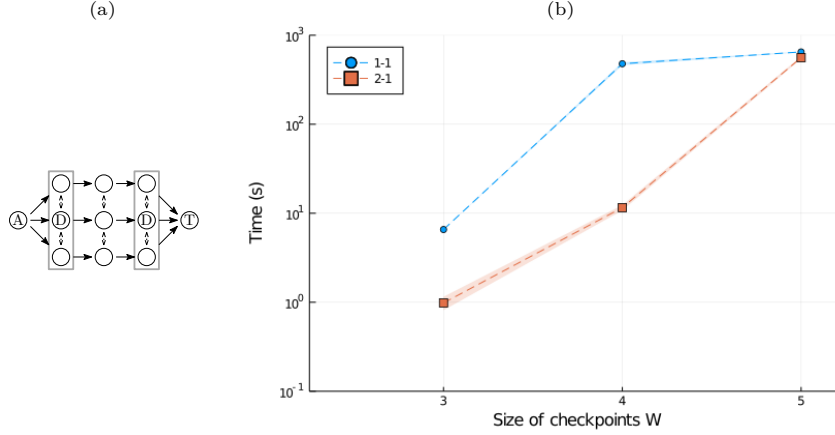


Figure 6: Intrusion search games: (a) Intrusion-search game with width $W = 3$ in configuration 1-1: A denotes the initial position of the attacker and D the positions of the defender’s units. T is the attacker’s target. (b) Time to reach $V_{UB}^T(b_{init}) - V_{LB}^T(b_{init}) \leq 1$.

box). The attacker aims to cross the graph while not being captured by the defender. She can either wait for one move to conceal her presence (and clean up the trace), or move further. Each unit of the defender can move to adjacent nodes within its assigned zone. The goal of the attacker is to cross the graph to reach node marked by T without encountering any unit of the defender. If she manages to do so, the defender receives a reward of -100 .

We consider games with two checkpoint zones with a varying number of nodes in a zone W (i.e. the width of the graph). We use two configurations of the defending forces: (1) one defender in each checkpoint and (2) two defenders in the first checkpoint and one defender in the second checkpoint. We denote these settings as 1-1 and 2-1.

The results are shown in Figure 6b (with five runs for each parameterization, the confidence intervals mark the standard error in our graphs). The largest game ($W = 5$ and two defenders in the first zone) has 4 656 states and 121 239 transitions and can be solved within 560 s. This case highlights that our algorithm can solve even large games. However, a much smaller game with $W = 5$ and configuration 1-1 (964 states and 9 633 transitions) is more challenging, since the coordination problem with just one defender in the first zone is harder, and despite its smaller size it is solved within 640 s.

Patrolling Games (inspired by [4, 50]). In a patrolling game, a patroller (player 1) aims to protect a set of targets V . The targets are represented by vertices of a graph, and the possible movements of the patroller are represented by the edges of the graph. The attacker observes the movement of the patroller and decides which target $v \in V$ he will attack, or whether he will postpone the decision. Once the attacker decides to attack a target v , the defender has t_x steps to reach the attacked vertex. If he fails to do so, he receives a negative reward $-C(v)$

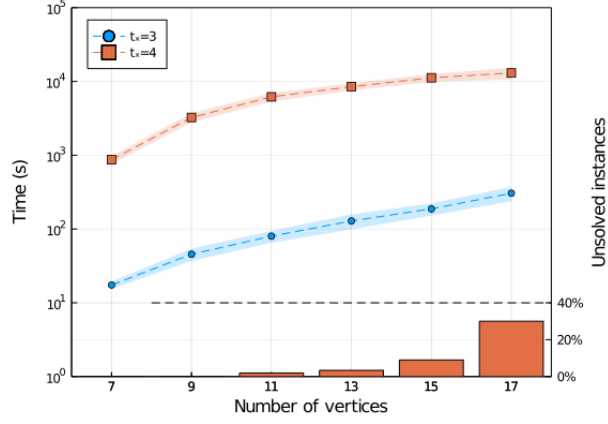


Figure 7: Time to reach $V_{UB}^{\Upsilon}(b_{init}) - V_{LB}^{\Gamma}(b_{init}) \leq 1$ for patrolling games with attack times $t_{\times} = 3$ and $t_{\times} = 4$. Bars indicate percentage of unsolved instances for $t_{\times} = 4$.

associated to the target v —otherwise, he successfully protects the target, and the reward is zero. The patroller does not know whether and where the attack has already started. The costs $C(v)$ are scaled so the $\max_{v \in V} C(v) = 100/\gamma^{t_{\times}}$, i.e., the minimum possible payoff for the defender is -100 .

Following the setting in [50], we focus on graphs generated from Erdos-Renyi model [33] with parameter $p = 0.25$ (denoted $ER(0.25)$) with attack times $t_{\times} \in \{3, 4\}$ and number of vertices $|\mathcal{V}|$ ranging from 7 to 15. The time to solve even the largest instances ($V = 17$) with $t_{\times} = 3$ was 305.5 s. For attack time $t_{\times} = 4$, however, some number of instances failed to reach the precision $V_{UB}^{\Upsilon}(b^{init}) - V_{LB}^{\Gamma}(b^{init}) \leq 1$ within the time limit of 10 hours. For the most difficult setting, $|\mathcal{V}| = 17$ and $t_{\times} = 4$, the algorithm reached desired precision in 70% of instances. For unsolved instances in this setting, mean $V_{UB}^{\Upsilon}(b^{init}) - V_{LB}^{\Gamma}(b^{init})$ after the cutoff after 10 hours is however reasonably small at 3.77 ± 0.54 . The results include games with up to 856 states and 6409 transitions. See Figure 7 for more details.

11.3. Impact of Initialization on Solution Time

Recall that we use value iteration algorithms for solving perfect information stochastic games and Markov decision processes, respectively, to initialize upper and lower bounds on V^* . In our experiments, we terminate the algorithms whenever the change in valuation between iterations of value iteration is smaller than $\beta = 0.025$. In Figure 8, we analyze the impact of the choice of β on the running time of the algorithm when applied to pursuit evasion games. Observe that the tighter initial bounds are used, the faster the convergence of the algorithm. In fact, the difference between $\beta = 1$ and $\beta = 0.025$ is approximately an order of magnitude in run time. Recall that the bounds V_{UB}^{Υ} and V_{LB}^{Γ} not only serve as bounds on V^* , but they are also used to obtain strategies that are considered during the forward exploration phase of the algorithm (see lines 6

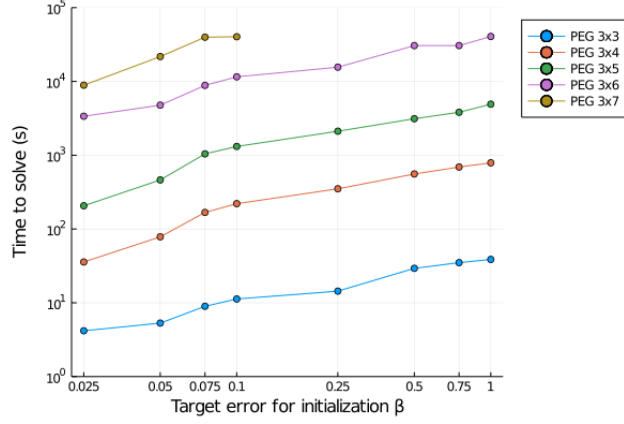


Figure 8: Effect of initialization on runtime. The target error is measured as Bellman residual $\|TV - V\|_\infty$ of the value iteration algorithms used to obtain initial bounds.

and 7 of Algorithm 4). We believe that these results indicate that the use of, e.g., domain-dependent initialization of the bounds can greatly improve the run time of the algorithm in complex domains.

11.4. Performance Analysis

Based on the the algorithm’s runtime data, we observed that most of the computation time is split between solving the linear programs used to compute HV_{UB}^Υ and HV_{LB}^Γ and pruning the representations of these bounds. Together, these three tasks typically took around 85% of the total runtime (and always at least 70%), with the remaining time being spent on computation of initial bounds, construction of the linear programs, and other smaller tasks. More specifically, solving HV_{UB}^Υ took 30-50% of the runtime in typical games while reaching as far as 60% in large pursuit evasion games (e.g., 60.5% in the 3×7 pursuit evasion game). Solving HV_{LB}^Γ was faster — in most games, it took between 10 and 20% of the total runtime. Finally, time required to perform pruning of the bounds V_{UB}^Υ and V_{LB}^Γ also took 10-20% of the runtime, with the exception of the patrolling games with attack time $t_\times = 4$, where it required over 40% of the total runtime.

12. Conclusions

We cover two-player zero-sum partially observable stochastic games (POSGs) with discounted rewards and one-sided observability — that is, those where the second player has perfect information about the game. We describe the theoretical properties of the value function in these games and show that algorithms based on value-iteration converge to an optimal value of the game. We also propose the first approximate algorithm that generalizes the ideas behind point-based algorithms

designed for partially observable Markov decision processes (POMDPs) and transfers these techniques to POSGs.

The presented work shows that it is possible to translate selected results from the single-agent setting to zero-sum games. Moreover, in future work, this work could be further extended in several ways: First, as already demonstrated by existing follow-up works [23], the scalability of the algorithm can be substantially improved for specific security games. Second, many heuristics and methods proven useful in the POMDP setting can be translated and evaluated in the game-theoretic setting, further improving the scalability. Third, generalization beyond the strictly adversarial setting (e.g., by computing a Stackelberg equilibrium) is another key direction supporting the applicability of these game-theoretic models to security.

Acknowledgements

This research was supported by the Czech Science Foundation (no. 19-24384Y), by the OP VVV MEYS funded project CZ.02.1.01/0.0/0.0/16 019/0000765 “Research Center for Informatics”, and by the U.S. Army Combat Capabilities Development Command Army Research Laboratory and was accomplished under Cooperative Agreement Number W911NF-13-2-0045 (ARL Cyber Security CRA). The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Combat Capabilities Development Command Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

References

- [1] F. Amigoni and N. Basilico. A game theoretical approach to finding optimal strategies for pursuit evasion in grid environments. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 2155–2162, 2012.
- [2] K. J. Astrom. Optimal control of Markov processes with incomplete state information. *Journal of mathematical analysis and applications*, 10(1): 174–205, 1965.
- [3] A. G. Barto, S. J. Bradtke, and S. P. Singh. Learning to act using real-time dynamic programming. *Artificial intelligence*, 72(1-2):81–138, 1995.
- [4] N. Basilico, N. Gatti, and F. Amigoni. Leader-follower strategies for robotic patrolling in environments with arbitrary topologies. In *8th International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 57–64, 2009. URL <http://portal.acm.org/citation.cfm?id=1558020>.

- [5] N. Basilico, G. D. Nittis, and N. Gatti. A Security Game Combining Patrolling and Alarm-Triggered Responses Under Spatial and Detection Uncertainties. In *30th AAAI Conference on Artificial Intelligence*, pages 397–403, 2016.
- [6] B. Bonet. Solving large POMDPs using real time dynamic programming. In *AAAI Fall Symposium on POMDPs*, 1998.
- [7] B. Bonet and H. Geffner. Solving POMDPs: RTDP-Bel vs. Point-based Algorithms. In *19th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1641–1646, 2009.
- [8] B. Bošanský, C. Kiekintveld, V. Lisý, and M. Pěchouček. An Exact Double-Oracle Algorithm for Zero-Sum Extensive-Form Games with Imperfect Information. *Journal of Artificial Intelligence Research*, 51:829–866, 2014.
- [9] M. Bowling and M. Veloso. An analysis of stochastic game theory for multi-agent reinforcement learning. Technical report, Carnegie Mellon University, 2000.
- [10] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [11] N. Brown and T. Sandholm. Superhuman ai for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.
- [12] N. Burch. *Time and space: Why imperfect information games are hard*. PhD thesis, University of Alberta, 2018.
- [13] K. Chatterjee and T. A. Henzinger. Semiperfect-information games. In *International Conference on Foundations of Software Technology and Theoretical Computer Science*, pages 1–18. Springer, 2005.
- [14] H.-T. Cheng. *Algorithms for partially observable Markov decision processes*. PhD thesis, University of British Columbia, 1988.
- [15] T. H. Chung, G. A. Hollinger, and V. Isler. Search and pursuit-evasion in mobile robotics. *Autonomous robots*, 31(4):299–316, 2011.
- [16] K. Ciesielski et al. On Stefan Banach and some of his results. *Banach Journal of Mathematical Analysis*, 1(1):1–10, 2007.
- [17] M. K. Ghosh, D. McDonald, and S. Sinha. Zero-Sum Stochastic Games with Partial Information. *Journal of Optimization Theory and Applications*, 121(1):99–118, 2004. ISSN 1573-2878.
- [18] E. A. Hansen, D. S. Bernstein, and S. Zilberstein. Dynamic Programming for Partially Observable Stochastic Games. In *National Conference on Artificial Intelligence (AAAI)*, pages 709–715, 2004.

- [19] M. Hauskrecht. Value-function approximations for partially observable Markov decision processes. *Journal of artificial intelligence research*, 13: 33–94, 2000.
- [20] K. Horák, B. Bošanský, and M. Pěchouček. Heuristic Search Value Iteration for One-Sided Partially Observable Stochastic Games. In *31st AAAI Conference on Artificial Intelligence*, pages 558–564, 2017.
- [21] K. Horák, B. Bošanský, and K. Chatterjee. Goal-HSVI: Heuristic Search Value Iteration for Goal POMDPs. In *27th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 4764–4770, 2018.
- [22] K. Horák, B. Bošanský, C. Kiekintveld, and C. Kamhoua. Compact Representation of Value Function in Partially Observable Stochastic Games. In *28th International Joint Conference on Artificial Intelligence (IJCAI)*, pages 350–356, 2019.
- [23] K. Horák, B. Bošanský, P. Tomášek, C. Kiekintveld, and C. Kamhoua. Optimizing honeypot strategies against dynamic lateral movement using partially observable stochastic games. *Computers & Security*, 87:101579, 2019. ISSN 0167-4048.
- [24] V. Isler and N. Karnad. The role of information in the cop-robber game. *Theoretical Computer Science*, 399(3):179–190, 2008.
- [25] V. Isler, S. Kannan, and S. Khanna. Randomized pursuit-evasion in a polygonal environment. *IEEE Transactions on Robotics*, 21(5):875–884, 2005.
- [26] A. Kumar and S. Zilberstein. Dynamic programming approximations for partially observable stochastic games. In *22nd International FLAIRS Conference*, pages 547–552, 2009.
- [27] H. Kurniawati, D. Hsu, and W. S. Lee. Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. In *Robotics: Science and systems*, volume 2008. Zurich, Switzerland., 2008.
- [28] V. Lisý, M. Lanctot, and M. Bowling. Online monte carlo counterfactual regret minimization for search in imperfect information games. In *Proceedings of the 2015 international conference on autonomous agents and multiagent systems*, pages 27–36, 2015.
- [29] M. L. Littman. *Algorithms for sequential decision making*. PhD thesis, Brown University, 1996.
- [30] C. L. MacDermed. *Value methods for efficiently solving stochastic games of complete and incomplete information*. PhD thesis, Georgia Institute of Technology, 2013.

- [31] O. Madani, S. Hanks, and A. Condon. On the undecidability of probabilistic planning and infinite-horizon partially observable Markov decision problems. In *16th National Conference on Artificial Intelligence (AAAI)*, pages 541–548, 1999.
- [32] M. Moravčík, M. Schmid, N. Burch, V. Lisý, D. Morrill, N. Bard, T. Davis, K. Waugh, M. Johanson, and M. Bowling. DeepStack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, pages 508–513, 2017. ISSN 0036-8075.
- [33] M. Newman. *Networks: an introduction*. Oxford university press, 2010.
- [34] H. Nikaido. On a minimax theorem and its applications to functional analysis. *Journal of the Mathematical Society of Japan*, 5(1):86–94, 1953.
- [35] J. Pineau, G. Gordon, S. Thrun, et al. Point-based value iteration: An anytime algorithm for POMDPs. In *3rd International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1025–1032, 2003.
- [36] D. Seitz, V. Kovařík, V. Lisý, J. Rudolf, S. Sun, and K. Ha. Value functions for depth-limited solving in imperfect-information games beyond poker. *arXiv preprint arXiv:1906.06412*, 2019.
- [37] D. Silver and J. Veness. Monte-carlo planning in large pomdps. In *Advances in neural information processing systems*, pages 2164–2172, 2010.
- [38] M. Sion. On general minimax theorems. *Pacific Journal of mathematics*, 8(1):171–176, 1958.
- [39] T. Smith and R. Simmons. Heuristic search value iteration for POMDPs. In *20th Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 520–527, 2004.
- [40] T. Smith and R. Simmons. Point-based POMDP algorithms: improved analysis and implementation. In *21st Conference on Uncertainty in Artificial Intelligence (UAI)*, pages 542–549, 2005.
- [41] A. Somani, N. Ye, D. Hsu, and W. S. Lee. DESPOT: Online POMDP planning with regularization. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1772–1780, 2013.
- [42] E. J. Sondik. The optimal control of partially observable Markov processes. Technical report, Stanford University, 1971.
- [43] E. J. Sondik. The optimal control of partially observable Markov processes over the infinite horizon: Discounted costs. *Operations research*, 26(2): 282–304, 1978.
- [44] S. Sorin. Stochastic games with incomplete information. In *Stochastic Games and applications*, pages 375–395. Springer, 2003.

- [45] M. T. Spaan and N. Vlassis. Perseus: Randomized point-based value iteration for POMDPs. *Journal of artificial intelligence research*, 24:195–220, 2005.
- [46] M. Šustr, V. Kovařík, and V. Lisý. Monte carlo continual resolving for online strategy computation in imperfect information games. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pages 224–232, 2019.
- [47] M. Sustr, M. Schmid, M. Moravcik, N. Burch, and M. Bowling. Sound search in imperfect information games. *arXiv preprint arXiv:2006.08740*, 2020.
- [48] R. J. Vanderbei. *Linear programming*. Springer, 2015.
- [49] J. von Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320, 1928.
- [50] Y. Vorobeychik, B. An, M. Tambe, and S. P. Singh. Computing Solutions in Infinite-Horizon Discounted Adversarial Patrolling Games. In *24th International Conference on Automated Planning and Scheduling (ICAPS)*, pages 314–322, 2014.
- [51] N. L. Zhang and W. Zhang. Speeding up the convergence of value iteration in partially observable Markov decision processes. *Journal of Artificial Intelligence Research*, 14:29–51, 2001.

Appendix A. Proofs

Proposition 5.3. *Let G be a one-sided POSG. Then the payoff Disc^γ of an arbitrary play in G is bounded by values*

$$L = \min_{(s, a_1, a_2)} R(s, a_1, a_2)/(1 - \gamma) \quad U = \max_{(s, a_1, a_2)} R(s, a_1, a_2)/(1 - \gamma) . \quad (14)$$

It also follows that $L \leq V^(b) \leq U$ and $L \leq \text{val}^{\sigma_1}(b) \leq U$ holds for every belief $b \in \Delta(S)$ and strategy $\sigma_1 \in \Sigma_1$ of the imperfectly informed player 1.*

Proof. The smallest payoff player 1 can hypothetically achieve in any play consists of getting $\underline{r} = \min_{(s, a_1, a_2)} R(s, a_1, a_2)$ in every timestep. The infinite discounted sum $\sum_{t=1}^{\infty} \gamma^{t-1} \underline{r}$ converges to $\underline{r}/(1 - \gamma) = L$. Conversely, the maximum payoff can be achieved if player 1 obtains $\bar{r} = \max_{(s, a_1, a_2)} R(s, a_1, a_2)$ in every timestep. Expected values of strategies (and therefore also the value of the game) are expectation over the payoffs of individual plays—hence are bounded by L and U as well. \square

Lemma 5.5. *Optimal value function V^* of a one-sided POSG is convex.*

Proof. Definition 5.2 defines V^* as the point-wise supremum over linear functions val^{σ_1} (over all strategies $\sigma_1 \in \Sigma_1$ of player 1). This implies the convexity of V^* [10, p.81]. \square

Lemma 5.6. *Let X be a finite set and let $f : \Delta(X) \rightarrow [y_{\min}, y_{\max}]$ be a linear function. Then f is k -Lipschitz continuous for $k = (y_{\max} - y_{\min})/2$.*

Proof. Let $p, q \in \Delta(X)$ be arbitrary two points in the probability simplex over the finite set X . Since f is a linear function, it can be represented as a convex combination of values $\alpha(x)$ in the vertices of the simplex corresponding to the elements $u \in X$,

$$f(p) = \sum_{u \in X} \alpha(u) \cdot p(u) \quad \text{where} \quad \alpha(u) = f(\mathbb{1}_u), \quad \mathbb{1}_u(v) = \begin{cases} 1 & v = u \\ 0 & \text{otherwise} \end{cases}. \quad (\text{A.1a})$$

Without loss of generality, let us assume $f(p) \geq f(q)$. Now, the difference $|f(p) - f(q)|$ satisfies

$$|f(p) - f(q)| = f(p) - f(q) = \sum_{u \in X} \alpha(u) \cdot [p(u) - q(u)]. \quad (\text{A.1b})$$

Denote $X^+ = \{u \in X \mid p(u) - q(u) \geq 0\}$ and $X^- = \{u \in X \mid p(u) - q(u) < 0\}$. We can now bound the difference from Equation (A.1b) by $|f(p) - f(q)| =$

$$= \sum_{u \in X^+} \alpha(u) \cdot [p(u) - q(u)] + \sum_{u \in X^-} \alpha(u) \cdot [p(u) - q(u)] \quad (\text{A.1c})$$

$$\leq \sum_{u \in X^+} y_{\max} \cdot [p(u) - q(u)] + \sum_{u \in X^-} y_{\min} \cdot [p(u) - q(u)] \quad (\text{A.1d})$$

$$= y_{\max} \sum_{u \in X^+} [p(u) - q(u)] + y_{\min} \sum_{u \in X^-} [p(u) - q(u)]. \quad (\text{A.1e})$$

Since both p and q belong to $\Delta(X)$, we have $\|p\|_1 = \|q\|_1 = 1$. Since $p(u), q(u) \geq 0$ are non-negative, we have

$$\|p\|_1 = \|q\|_1 + \sum_{u \in X^+} [p(u) - q(u)] - \sum_{u \in X^-} [q(u) - p(u)].$$

It follows that

$$\sum_{u \in X^+} [p(u) - q(u)] = - \sum_{u \in X^-} [p(u) - q(u)]. \quad (\text{A.1f})$$

From equation (A.1f), we further see that both terms in (A.1f) are equal to $\|p - q\|_1/2$. This implies that

$$\begin{aligned} |f(p) - f(q)| &\leq y_{\max} \|p - q\|_1/2 + y_{\min} (-\|p - q\|_1/2) \\ &= (y_{\max} - y_{\min})/2 \cdot \|p - q\|_1 \end{aligned}$$

and completes the proof. \square

Proposition 5.8. *Value function V^* of one-sided POSGs is δ -Lipschitz continuous.*

Proof. V^* is defined as a supremum over δ -Lipschitz continuous values val^{σ_1} of strategies $\sigma_1 \in \Sigma_1$ of the imperfectly informed player 1. Therefore for arbitrary $b, b' \in \Delta(S)$, we have the following

$$V^*(b) = \sup_{\sigma_1 \in \Sigma_1} \text{val}^{\sigma_1}(b) \leq \sup_{\sigma_1 \in \Sigma_1} [\text{val}^{\sigma_1}(b') + \delta \|b - b'\|_1] = V^*(b') + \delta \|b - b'\|_1 . \quad (\text{A.2})$$

□

Proposition 5.11. *Let $\Gamma \subset \{\alpha : \Delta(S) \rightarrow \mathbb{R} \mid \alpha \text{ is linear}\}$ be a set of linear functions. Then for every $b \in \Delta(S)$ we have*

$$\sup_{\alpha \in \Gamma} \alpha(b) = \sup_{\alpha \in \text{Conv}(\Gamma)} \alpha(b) . \quad (18)$$

Proof. Clearly, it suffices to prove the inequality \geq . Let $b \in \Delta(S)$ and let $\sum_{i=1}^k \lambda_i \alpha_i$ be an arbitrary¹⁰ convex combination of linear functions from Γ (i.e., we have $\alpha_i \in \Gamma$). We need to show that $\alpha(b) \geq \sum_{i=1}^k \lambda_i \alpha_i(b)$ holds for some $\alpha \in \Gamma$. This is straightforward, as can be witnessed by the function $\alpha_{i^*} \in \Gamma$, $i^* := \arg \max_i \alpha_i(b)$:

$$\sum_{i=1}^k \lambda_i \alpha_i(b) \leq \sum_{i=1}^k \lambda_i \max_{1 \leq i \leq k} \alpha_i(b) = \max_{1 \leq i \leq k} \alpha_i(b) = \alpha_{i^*}(b) .$$

□

Proposition 5.12. *Let $f : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function. Then there exists a set Γ of linear functions such that $\alpha \leq f$ for every $\alpha \in \Gamma$ and $f(b) = \sup_{\alpha \in \Gamma} \alpha(b)$ for every $b \in \Delta(S)$.*

Proof. Let $\Gamma := \{\alpha : \Delta(S) \rightarrow \mathbb{R} \text{ linear} \mid \alpha \leq f\}$. Clearly, the pointwise supremum of Γ is no greater than f . It remains to show that $\sup_{\alpha \in \Gamma} \alpha(b_0) \geq f(b_0)$ for each b_0 . Let b_0 be an interior point of $\Delta(S)$. By the standard convex-analysis result, there exists a subdifferential of f at b_0 , that is, a vector v such that $f(b) \geq f(b_0) + v \cdot (b - b_0)$ holds for each $b \in \Delta(S)$. The function $\alpha(b) := f(b_0) + v \cdot (b - b_0)$ therefore belongs to Γ and witnesses that $\sup_{\alpha \in \Gamma} \alpha(b_0) \geq f(b_0)$.

Suppose that b_0 lies at the boundary of $\Delta(S)$ and let η , $\|\eta\|_1 = 1$, be a direction in which every nearby point $b_\delta := b_0 - \delta\eta$, $\delta \in (0, \Delta]$, lies in the interior of $\Delta(S)$ (for some $\Delta > 0$). Since f is convex, the directional derivatives $f'_\eta(b_\delta) = \lim_{g \rightarrow 0^+} \frac{f(b_\delta + g\eta) - f(b_\delta)}{g}$ are non-decreasing as the points b_δ get closer to b_0 . In particular, the linear functions α_δ found for b_δ in the previous step satisfy

$$\alpha_\delta(b_0) \geq f(b_\delta) + f'_\eta(b_\delta)\delta \geq f(b_\delta) + f'_\eta(b_\Delta)\delta .$$

The right-hand side converges to $f(b_0) + f'_\eta(b_\Delta) \cdot 0 = f(b_0)$, which shows that the supremum of $\alpha_\delta(b_0)$ is at least $f(b_0)$. Since $\alpha_\delta \in \Gamma$, this proves the remaining part of the proposition. □

¹⁰Recall that according to the Carathéodory's theorem, it suffices to consider finite convex combinations.

Proposition 6.2. *Every behavioral strategy $\sigma_1 \in \Sigma_1$ of player 1 can be represented as a strategy composition of some stage strategy $\pi_1 \in \Pi_1$ and player 1 behavioral strategies $\zeta_{a_1,o}$.*

Proof. Let $\sigma_1 \in \Sigma_1$ be an arbitrary behavioral strategy of player 1, and let $\pi_1 = \sigma_1(\emptyset)$ and $\zeta_{a_1,o}(\omega') = \sigma_1(\omega')$ for every $(a_1, o) \in A_1 \times O$ and $\omega' \in (A_1 O)^*$. It can be easily verified that $\text{comp}(\pi_1, \bar{\zeta})$ defined in Definition 6.1 satisfies $\text{comp}(\pi_1, \bar{\zeta}) = \sigma_1$. \square

Lemma 6.3. *Let G be a one-sided POSG and $\text{comp}(\pi_1, \bar{\zeta})$ a composite strategy. Then the following holds:*

$$\begin{aligned} \text{val}^{\text{comp}(\pi_1, \bar{\zeta})}(s) &= \min_{a_2 \in A_2} \mathbb{E}_{a_1 \sim \pi_1, (o, s') \sim T(\cdot | s, a_1, a_2)} \left[R(s, a_1, a_2) + \gamma \text{val}^{\zeta_{a_1,o}}(s') \right] \\ &= \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1(a_1) \left[R(s, a_1, a_2) + \gamma \sum_{(o, s') \in O \times S} T(o, s' | s, a_1, a_2) \text{val}^{\zeta_{a_1,o}}(s') \right]. \end{aligned} \quad (20)$$

Proof. Let us evaluate the payoff if player 2 uses a_2 in the first stage of the game given that the initial state of the game is s . The expected reward of playing action a_2 against $\text{comp}(\pi_1, \bar{\zeta})$ in the first stage is $\sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2)$, i.e., the expectation over the actions player 1 can take. Now, at the beginning of the next stage, player 2 knows everything about the past stage—including action a_1 taken by player 1, observation o he received, and the new state of the game s' . Therefore, player 2 knows the strategy $\zeta_{a_1,o}$ player 1 is about to use in the rest of the game. By definition of $\text{val}^{\zeta_{a_1,o}}$ (Definition 5.1), the best payoff player 2 can achieve in (a_1, o) -subgame is $\text{val}^{\zeta_{a_1,o}}(s')$. After reaching the subgame, however, one stage has already passed and the rewards originally received at time t are now received at time $t + 1$. As a result, the reward $\text{val}^{\zeta_{a_1,o}}(s')$ gets discounted by γ . The probability that the (a_1, o) -subgame is reached is $\sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1(a_1) T(o, s' | s, a_1, a_2)$, and the expectation over $\gamma \text{val}^{\zeta_{a_1,o}}(s')$ is thus computed. Player 2 chooses an action which achieves the minimum payoff which completes the proof. \square

Lemma 6.5. *Let $\pi_1 \in \Pi_1$ be a stage strategy of player 1 and $\bar{\alpha} \in (\text{lin}_{\Delta(S)})^{A_1 \times O}$ a vector of linear functions s.t. for each $\alpha_{a_1,o}$ there exists a strategy $\zeta_{a_1,o} \in \Sigma_1$ with $\text{val}^{\zeta_{a_1,o}} \geq \alpha_{a_1,o}$. Then there exists a strategy $\sigma_1 \in \Sigma_1$ such that $\sigma_1(\emptyset) = \pi_1$ and $\text{val}^{\sigma_1} \geq \text{valcomp}(\pi_1, \bar{\alpha})$.*

Proof. Let $\bar{\zeta} \in (\Sigma_1)^{A_1 \times O}$ be as in the lemma, and let $\bar{\alpha}^\zeta$ be such that $\alpha_{a_1,o}^\zeta = \text{val}^{\zeta_{a_1,o}}$. According to the assumption we have $\alpha_{a_1,o}^\zeta \geq \alpha_{a_1,o}$. Replacing $\alpha_{a_1,o}$ by $\alpha_{a_1,o}^\zeta$ in Equation (21) can only increase the objective value, hence

$$\text{valcomp}(\pi_1, \bar{\alpha})(s) \leq \text{valcomp}(\pi_1, \bar{\alpha}^\zeta)(s) = \text{val}^{\text{comp}(\pi_1, \bar{\zeta})}(s). \quad (\text{A.3})$$

Composite strategies are behavioral strategies of player 1, hence $\sigma_1 = \text{comp}(\pi_1, \bar{\zeta})$. \square

Lemma 6.6. *Let $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in (\text{lin}_{\Delta(S)})^{A_1 \times O}$ such that $L \leq \alpha^{a_1, o}(b) \leq U$ for every $b \in \Delta(S)$. Then $L \leq \text{valcomp}(\pi_1, \bar{\alpha})(b) \leq U$ for every $b \in \Delta(S)$ and $\text{valcomp}(\pi_1, \bar{\alpha})$ is a δ -Lipschitz continuous function.*

Proof. Since $\text{valcomp}(\pi_1, \bar{\alpha})(b)$ is calculated as a convex combination of the values $\text{valcomp}(\pi_1, \bar{\alpha})(s)$ in the vertices of the $\Delta(S)$ simplex, it suffices to show that

$$(\forall s \in S) : L \leq \text{valcomp}(\pi_1, \bar{\alpha})(s) \leq U.$$

Let $a_2^* \in A_2$ be the minimizing action of player 2 in Equation (21). It holds $\underline{r} \leq R(s, a_1, a_2^*) \leq \bar{r}$, where \underline{r} and \bar{r} are minimum and maximum rewards in the game. Hence $\underline{r} \leq \sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2^*) \leq \bar{r}$. Similarly, from the assumption of the lemma, we have $L \leq \alpha_{a_1, o}(s') \leq U$ and hence $L \leq \sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1(a_1) T(o, s' | s, a_1, a_2^*) \alpha_{a_1, o}(s') \leq U$. We will now prove that $\text{valcomp}(\pi_1, \bar{\alpha})(s) \leq U$ (the proof of $\text{valcomp}(\pi_1, \bar{\alpha})(s) \geq L$ is analogous):

$$\begin{aligned} \text{valcomp}(\pi_1, \bar{\alpha})(s) &= \\ &= \sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2^*) + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1(a_1) T(o, s' | s, a_1, a_2^*) \alpha_{a_1, o}(s') \\ &\leq \bar{r} + \gamma U = \bar{r} + \gamma \frac{\bar{r}}{1 - \gamma} = U. \end{aligned}$$

The δ -Lipschitz continuity of $\text{valcomp}(\pi_1, \bar{\alpha})$ then follows directly from Lemma 5.6. \square

Proposition 7.2. *Proposition Let $V : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function and let Γ be a convex set of linear functions such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$. Then HV is also convex and continuous. Furthermore, if V is δ -Lipschitz continuous, the function HV is δ -Lipschitz continuous as well.*

Proof. According to Definition 7.1, operator H can be rewritten as a supremum over all possible value compositions:

$$[HV](b) = \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) = \sup_{(\pi_1, \bar{\alpha}) \in \Pi_1 \times \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b), \text{ and} \quad (\text{A.4a})$$

$$[HV](b) = \sup_{\alpha \in \Gamma'} \alpha(b) \quad \Gamma' = \{\text{valcomp}(\pi_1, \bar{\alpha}) \mid \pi_1 \in \Pi_1, \bar{\alpha} \in \Gamma^{A_1 \times O}\}. \quad (\text{A.4b})$$

In Equation (A.4b), HV is represented as a point-wise supremum from a set Γ' of linear functions $\text{valcomp}(\pi_1, \bar{\alpha})$, which is a convex continuous function (see Proposition 5.9).

Moreover, in case V is δ -Lipschitz continuous, the set Γ representing V can be assumed to contain only δ -Lipschitz continuous linear functions. According to Lemma 6.6, $\text{valcomp}(\pi_1, \bar{\alpha})$ is δ -Lipschitz continuous for every $\pi_1 \in \Pi_1$ and $\alpha^{a_1, o} \in \Gamma$. Hence, Γ' contains δ -Lipschitz continuous linear functions only and the point-wise maximum HV over Γ' is δ -Lipschitz continuous. \square

Theorem 1. Let $V : \Delta(S) \rightarrow \mathbb{R}$ be a convex continuous function and let Γ be a convex set of linear functions on $\Delta(S)$ such that $V(b) = \sup_{\alpha \in \Gamma} \alpha(b)$ for every belief $b \in \Delta(S)$. Then the following definitions of operator H are equivalent:

$$[HV](b) = \max_{\pi_1 \in \Delta(S)} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) \quad (24a)$$

$$= \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} \left[\mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)) \right] \quad (24b)$$

$$= \min_{\pi_2 \in \Pi_2} \max_{\pi_1 \in \Pi_1} \left[\mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)) \right]. \quad (24c)$$

Proof. We first prove the equality of (24b) and (24c). Let us define a payoff function $u : \Pi_1 \times \Pi_2 \rightarrow \mathbb{R}$ to be the objective of the maximin and minimax optimization in (24b) and (24c).

$$u(\pi_1, \pi_2) = \mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot V(\tau(b, a_1, \pi_2, o)) \quad (\text{A.5a})$$

After expanding the expectation $\mathbb{E}_{b, \pi_1, \pi_2} [R(s, a_1, a_2)]$ and expressing V as a supremum over linear functions $\alpha \in \Gamma$, we get

$$\begin{aligned} u(\pi_1, \pi_2) &= \sum_{s, a_1, a_2} b(s) \pi_1(a_1) \pi_2(a_2 | s) R(s, a_1, a_2) + \\ &\quad + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2} [a_1, o] \cdot \sup_{\alpha \in \Gamma} \sum_{s'} \tau(b, a_1, \pi_2, o)(s') \cdot \alpha(s') \end{aligned} \quad (\text{A.5b})$$

$$\begin{aligned} &= \sum_{s, a_1, a_2} b(s) \pi_1(a_1) \pi_2(a_2 | s) R(s, a_1, a_2) + \\ &\quad + \gamma \sum_{a_1, o} \pi_1(a_1) \cdot \sup_{\alpha \in \Gamma} \sum_{s, a_2, s'} b(s) \pi_2(a_2 | s) T(o, s' | s, a_1, a_2) \alpha(s'). \end{aligned} \quad (\text{A.5c})$$

Note that the term $\mathbb{P}_{b, \pi_1, \pi_2} [a_1, o]$ cancels out after expanding $\tau(b, a_1, \pi_2, o)$ in Equation (A.5c).

We now show that the von Neumann's minimax theorem [49, 34] applies to the game with utility function u and strategy spaces Π_1 and Π_2 for player 1 and player 2, respectively. The von Neumann's minimax theorem requires that the strategy spaces Π_1 and Π_2 are convex compact sets (which is clearly the case), and that the utility function u (as characterized by Equation (A.5c)) is continuous, convex in Π_2 and concave in Π_1 . We will now prove the latter and show that u is a convex-concave utility function. Clearly, for every $\pi_2 \in \Pi_2$, the function $u(\cdot, \pi_2) : \Pi_1 \rightarrow \mathbb{R}$ (where π_2 is considered constant) is linear in π_1 , and hence also concave. The convexity of $u(\pi_1, \cdot) : \Pi_2 \rightarrow \mathbb{R}$ (after fixing

arbitrary $\pi_1 \in \Pi_1$) is more involved. As weighted sum of convex functions with positive coefficients $\pi_1(a_1) \geq 0$ is also convex, it is sufficient to show that $f(\pi_2) = \sup_{\alpha \in \Gamma} \sum_{s,a_2,s'} b(s)\pi_2(a_2|s)T(o,s'|s,a_1,a_2)\alpha(s')$ is convex. Observe that for every $\alpha \in \Gamma$, the expression $\sum_{s,a_2,s'} b(s)\pi_2(a_2|s)T(o,s'|s,a_1,a_2)\alpha(s')$ is linear in π_2 and, as a result, the supremum over such linear expressions in π_2 is convex in π_2 (see Proposition 5.9). According to von Neumann's minimax theorem $\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} u(\pi_1, \pi_2) = \min_{\pi_2 \in \Pi_2} \max_{\pi_1 \in \Pi_1} u(\pi_1, \pi_2)$ which concludes the proof of equality of (24b) and (24c).

We now proceed by showing the equality of (24a) and (24b). By further rearranging Equation (A.5c), we get

$$u(\pi_1, \pi_2) = \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \left[\sum_{s,a_1,a_2} b(s)\pi_1(a_1)\pi_2(a_2|s)R(s,a_1,a_2) + \right. \quad (\text{A.6}) \\ \left. + \gamma \sum_{a_1,o} \pi_1(a_1) \sum_{s,a_2,s'} b(s)\pi_2(a_2|s)T(o,s'|s,a_1,a_2)\alpha_{a_1,o}(s') \right].$$

Let us define a game with strategy spaces Γ and Π_2 and payoff function $u'_{\pi_1} : \Gamma \times \Pi_2 \rightarrow \mathbb{R}$ where u'_{π_1} is the objective of the supremum in Equation (A.6) (Equation (A.7b) is an algebraic simplification of Equation (A.7a)).

$$u'_{\pi_1}(\bar{\alpha}, \pi_2) = \sum_{s,a_1,a_2} b(s)\pi_1(a_1)\pi_2(a_2|s)R(s,a_1,a_2) + \quad (\text{A.7a})$$

$$+ \gamma \sum_{a_1,o} \pi_1(a_1) \sum_{s,a_2,s'} b(s)\pi_2(a_2|s)T(o,s'|s,a_1,a_2)\alpha_{a_1,o}(s') \\ = \sum_s b(s) \sum_{a_2} \pi_2(a_2|s) \sum_{a_1} \pi_1(a_1) \left[R(s,a_1,a_2) + \right. \quad (\text{A.7b}) \\ \left. + \gamma \sum_{o,s'} T(o,s'|s,a_1,a_2)\alpha_{a_1,o}(s') \right].$$

Plugging (A.7b) into (A.6), we can write

$$\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} u(\pi_1, \pi_2) = \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} u'_{\pi_1}(\pi_2, \bar{\alpha}). \quad (\text{A.8})$$

To prove the equivalence of (24a) and (24b), we need to show that the minimum and supremum can be swapped. Since u'_{π_1} is linear in both π_2 and $\bar{\alpha}$, Π_2 is a compact convex set and Γ (and therefore also the set of mappings $\bar{\alpha} \in \Gamma^{A_1 \times O}$) is convex, it is possible to apply Sion's minimax theorem [38] to get

$$\max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} u'_{\pi_1}(\pi_2, \bar{\alpha}) = \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \min_{\pi_2 \in \Pi_2} u'_{\pi_1}(\pi_2, \bar{\alpha}). \quad (\text{A.9})$$

As u'_{π_1} is linear in π_2 (for fixed π_1 and $\bar{\alpha}$), the minimum over π_2 is attained in pure strategies. Denote $\hat{\pi}_2 : S \rightarrow A_2$ a pure strategy of player 2 assigning action $\hat{\pi}_2(s)$ to be played in state s , and $\hat{\Pi}_2$ the set of all pure strategies of player 2. We now rewrite u'_{π_1} to use pure strategies $\hat{\Pi}_2$ instead of randomized

stage strategies Π_2 . First, in Equation (A.10a), we replace the maximization over Π_2 by maximization over the pure strategies $\hat{\Pi}_2$ and replace expectation over actions of player 2 by using the deterministic action $\hat{\pi}_2(s)$ where appropriate. Then, in Equation (A.10b), we leverage the fact that, unlike player 1, player 2 knows the state before having to act, and hence he can optimize his actions $\hat{\pi}_2(s)$ independently. And, finally, in Equation (A.10c), we use Definition 6.4.

$$\begin{aligned} \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} u(\pi_1, \pi_2) &= \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \min_{\pi_2 \in \Pi_2} u'_{\pi_1}(\pi_2, \bar{\alpha}) = \\ &= \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \min_{\hat{\pi}_2 \in \hat{\Pi}_2} \sum_s b(s) \sum_{a_1} \pi_1(a_1) \left[R(s, a_1, \hat{\pi}_2(s)) + \right. \end{aligned} \quad (\text{A.10a})$$

$$\begin{aligned} &\quad \left. + \gamma \sum_{o, s'} T(o, s' \mid s, a_1, \hat{\pi}_2(s)) \alpha_{a_1, o}(s') \right] \\ &= \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \sum_s b(s) \min_{\hat{\pi}_2(s) \in A_2} \sum_{a_1} \pi_1(a_1) \left[R(s, a_1, \hat{\pi}_2(s)) + \right. \end{aligned} \quad (\text{A.10b})$$

$$\begin{aligned} &\quad \left. + \gamma \sum_{o, s'} T(o, s' \mid s, a_1, \hat{\pi}_2(s)) \alpha_{a_1, o}(s') \right] \\ &= \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \sum_s b(s) \cdot \text{valcomp}(\pi_1, \bar{\alpha})(s) \\ &= \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) . \end{aligned} \quad (\text{A.10c})$$

This concludes the proof of the equality of Equations (24a) and (24b). \square

Lemma 7.5. *Let $V, W : \Delta(S) \rightarrow \mathbb{R}$ be two convex continuous value functions and $b \in \Delta(S)$ a belief such that $[HV](b) \leq [HW](b)$. Let (π_1^V, π_2^V) and (π_1^W, π_2^W) be Nash equilibrium strategy profiles in stage games $[HV](b)$ and $[HW](b)$, respectively, and $C \geq 0$. If $W(\tau(b, a_1, o, \pi_2^V)) - V(\tau(b, a_1, o, \pi_2^V)) \leq C$ for every action $a_1 \in \text{Supp}(\pi_1^W)$ of player 1 and every observation $o \in O$ such that $\mathbb{P}_{b, \pi_1^W, \pi_2^V}[o \mid a_1] > 0$, then $[HW](b) - [HV](b) \leq \gamma C$.*

Proof. By deviating from the equilibrium strategy profiles in stage games $[HV](b)$ and $[HW](b)$, the players can only worsen their payoffs. Therefore, we have

$$\begin{aligned} u^{V,b}(\pi_1^W, \pi_2^V) &\leq u^{V,b}(\pi_1^V, \pi_2^V) = [HV](b) \leq \\ &\leq [HW](b) = u^{W,b}(\pi_1^W, \pi_2^W) \leq u^{W,b}(\pi_1^W, \pi_2^V) . \end{aligned} \quad (\text{A.11})$$

We can thus bound the difference $[HW](b) - [HV](b)$ by $u^{W,b}(\pi_1^W, \pi_2^V) - u^{V,b}(\pi_1^W, \pi_2^V)$ where, according to Definition 7.4,

$$\begin{aligned} u^{W,b}(\pi_1^W, \pi_2^V) - u^{V,b}(\pi_1^W, \pi_2^V) &= \\ &= \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1^W, \pi_2^V}[a_1, o] \cdot [W(\tau(b, a_1, \pi_2^V, o)) - V(\tau(b, a_1, \pi_2^V, o))] . \end{aligned} \quad (\text{A.12})$$

Since every $W(\tau(b, a_1, o, \pi_2^V)) - V(\tau(b, a_1, o, \pi_2^V))$ considered in Equation (A.12) with non-zero probability $\mathbb{P}_{b, \pi_1^W, \pi_2^V}[a_1, o]$ is assumed to be bounded by C , the

expectation over such $W(\tau(b, a_1, o, \pi_2^V)) - V(\tau(b, a_1, o, \pi_2^V))$ is likewise bounded by C . It follows that $u^{W,b}(\pi_1^W, \pi_2^V) - u^{V,b}(\pi_1^W, \pi_2^V) \leq \gamma C$, and hence we also have $[HW](b) - [HV](b) \leq \gamma C$. \square

Lemma 7.7. *Lemma The optimal value function V^* satisfies $V^* = HV^*$.*

Proof. According to Corollary 7.3, the Bellman's operator does not depend on the set Γ used to represent the value function V^* . To this end, we will assume that the set Γ used to represent V^* is

$$\Gamma = \text{Conv}\{\text{val}^{\sigma_1} \mid \sigma_1 \in \Sigma_1\} . \quad (\text{A.13a})$$

To prove the equivalence of value functions V^* and HV^* we consider that these functions are represented as follows:

$$V^*(b) = \sup_{\alpha \in \Gamma_{V^*}} \alpha(b) \quad \Gamma_{V^*} = \{\text{val}^{\sigma_1} \mid \sigma_1 \in \Sigma_1\} \quad (\text{A.13b})$$

$$[HV^*](b) = \sup_{\alpha \in \Gamma_{HV^*}} \alpha(b) \quad \Gamma_{HV^*} = \{\text{valcomp}(\pi_1, \bar{\alpha}) \mid \pi_1 \in \Pi_1, \bar{\alpha} \in \Gamma^{A_1 \times O}\} . \quad (\text{A.13c})$$

To prove the equivalence of V^* and HV^* , it suffices to show that for every $\alpha \in \Gamma_{V^*}$ there exists $\alpha' \in \Gamma_{HV^*}$ such that $\alpha' \geq \alpha$, and vice versa.

First, from Proposition 6.2, Lemma 6.3 and Definition 6.4, it follows that every strategy $\sigma_1 \in \Sigma_1$ can be represented as a value composition $\text{valcomp}(\pi_1, \bar{\zeta})$, and we have

$$\text{val}^{\sigma_1} = \text{val}^{\text{comp}(\pi_1, \bar{\zeta})} = \text{valcomp}(\pi_1, \bar{\alpha}^{\bar{\zeta}}) \quad (\text{A.13d})$$

where $\alpha_{a_1, o}^{\bar{\zeta}} = \text{val}^{\bar{\zeta}_{a_1, o}} \in \Gamma$. Hence $\text{val}^{\sigma_1} = \text{valcomp}(\pi_1, \bar{\zeta}) \in \Gamma_{HV^*}$.

The opposite direction of the proof, i.e., that for every $\alpha \in \Gamma_{HV^*}$ there exists $\alpha' \in \Gamma_{V^*}$ such that $\alpha' \geq \alpha$, is more involved. Let $\alpha = \text{valcomp}(\pi_1, \bar{\alpha}) \in \Gamma_{HV^*}$ be arbitrary. From (A.13c), each $\alpha_{a_1, o}$ can be written as a convex combination of finitely many elements of $\{\text{val}^{\sigma_1} \mid \sigma_1 \in \Sigma_1\}$.

$$\alpha_{a_1, o} = \sum_{i=1}^K \lambda_i^{a_1, o} \text{val}^{\sigma_1^{a_1, o, i}} \quad (\text{A.13e})$$

Let us form a vector of strategies $\bar{\zeta} \in (\Sigma_1)^{A_1 \times O}$ such that each $\zeta_{a_1, o}$ is a convex combination of strategies $\sigma_1^{a_1, o, i}$ using coefficients from Equation (A.13e),

$$\zeta_{a_1, o} = \sum_{i=1}^K \lambda_i^{a_1, o} \sigma_1^{a_1, o, i} . \quad (\text{A.13f})$$

We can interpret strategy $\zeta_{a_1, o}$ as player 1 first randomly choosing among strategies $\sigma_1^{a_1, o, i}$, and then following the chosen strategy in the rest of the game. If the player 2 knew which strategy $\sigma_1^{a_1, o, i}$ has been chosen, he is able to achieve

utility $\text{val}^{\sigma_1^{a_1, o, i}}$. However, he has no access to this information, and hence $\text{val}^{\zeta^{a_1, o}} \geq \sum_{i=1}^K \lambda_i^{a_1, o} \text{val}^{\sigma_1^{a_1, o, i}} = \alpha^{a_1, o}$. Now, we have

$$\alpha' = \text{val}^{\text{comp}(\pi_1, \bar{\zeta})} \geq \text{valcomp}(\pi_1, \bar{\alpha}) = \alpha \quad (\text{A.13g})$$

which concludes the proof. \square

Lemma 8.1. *Let $\Gamma = \text{Conv}(\{\alpha_1, \dots, \alpha_k\})$ be a convex hull of a finite set of α -vectors. Then $[HV](b)$ coincides with the solution of the following linear program:*

$$\max_{\pi_1, \lambda, \bar{\alpha}, V} \sum_{s \in S} b(s) \cdot V(s) \quad (27a)$$

$$\begin{aligned} s.t. \quad V(s) \leq & \sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2) + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} T(o, s' | s, a_1, a_2) \hat{\alpha}^{a_1, o}(s') \\ & \forall (s, a_2) \in S \times A_2 \end{aligned} \quad (27b)$$

$$\hat{\alpha}^{a_1, o}(s') = \sum_{i=1}^k \hat{\lambda}_i^{a_1, o} \cdot \alpha_i(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (27c)$$

$$\sum_{i=1}^k \hat{\lambda}_i^{a_1, o} = \pi_1(a_1) \quad \forall (a_1, o) \in A_1 \times O \quad (27d)$$

$$\sum_{a_1 \in A_1} \pi_1(a_1) = 1 \quad (27e)$$

$$\pi_1(a_1) \geq 0 \quad \forall a_1 \in A_1 \quad (27f)$$

$$\hat{\lambda}_i^{a_1, o} \geq 0 \quad \forall (a_1, o) \in A_1 \times O, 1 \leq i \leq k \quad (27g)$$

Proof. Since the set Γ is convex and compact, the dynamic programming operator H can be used:

$$[HV](b) = \max_{\pi_1 \in \Pi_1} \sup_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) \quad (\text{A.14a})$$

$$= \max_{\pi_1 \in \Pi_1} \max_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \text{valcomp}(\pi_1, \bar{\alpha})(b) \quad (\text{A.14b})$$

$$= \max_{\pi_1 \in \Pi_1} \max_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \sum_{s \in S} b(s) \cdot \text{valcomp}(\pi_1, \bar{\alpha})(s) \quad (\text{A.14c})$$

$$\begin{aligned} = \max_{\pi_1 \in \Pi_1} \max_{\bar{\alpha} \in \Gamma^{A_1 \times O}} \sum_{s \in S} b(s) \cdot \min_{a_2} \left[\sum_{a_1} \pi_1(a_1) R(s, a_1, a_2) + \right. & (\text{A.14d}) \\ & \left. + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} T(o, s' | s, a_1, a_2) \pi_1(a_1) \alpha^{a_1, o}(s') \right]. \end{aligned}$$

Equation (A.14b) follows from the fact that $\text{valcomp}(\pi_1, \bar{\alpha})$ is continuous in $\bar{\alpha}$, and Γ is a compact set (and hence also $\Gamma^{A_1 \times O}$ is). The Equation (A.14c) represents value of the linear function $\text{valcomp}(\pi_1, \bar{\alpha})$ as the convex combination of its values in the vertices of the $\Delta(S)$ simplex, and, finally, Equation (A.14d) rewrites $\text{valcomp}(\pi_1, \bar{\alpha})(s)$ using Definition 6.4.

Equation (A.14d) can be directly formalized as a mathematical program (A.15) whose solution is $[HV](b)$. Indeed, the minimization over $a_2 \in A_2$ can be rewritten as a set of constraints for each value of state $V(s)$ (one for each action $a_2 \in A_2$ of player 2) in Equation (A.15b). The convex hull of set $\{\alpha_1, \dots, \alpha_k\}$ is represented by (A.15c) where variables $\lambda_i^{a_1, o}$ represent coefficients of the convex combination. The stage strategy π_1 is characterized by (A.15e) and (A.15f).

$$\max_{\pi_1, \lambda, \bar{\alpha}, V} \sum_{s \in S} b(s) \cdot V(s) \quad (\text{A.15a})$$

$$\text{s.t. } V(s) \leq \sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2) + \quad \forall (s, a_2) \in S \times A_2 \quad (\text{A.15b})$$

$$+ \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} T(o, s' | s, a_1, a_2) \pi_1(a_1) \alpha^{a_1, o}(s') \\ \alpha^{a_1, o}(s') = \sum_{i=1}^k \lambda_i^{a_1, o} \cdot \alpha_i(s') \quad \forall (a_1, o, s') \in A_1 \times O \times S \quad (\text{A.15c})$$

$$\sum_{i=1}^k \lambda_i^{a_1, o} = 1 \quad \forall (a_1, o) \in A_1 \times O \quad (\text{A.15d})$$

$$\sum_{a_1 \in A_1} \pi_1(a_1) = 1 \quad (\text{A.15e})$$

$$\pi_1(a_1) \geq 0 \quad \forall a_1 \in A_1 \quad (\text{A.15f})$$

$$\lambda_i^{a_1, o} \geq 0 \quad \forall (a_1, o) \in A_1 \times O, 1 \leq i \leq k \quad (\text{A.15g})$$

This mathematical program is not linear since it contains a product of variables $\pi_1(a) \cdot \alpha^{a_1, o}(s')$. It can, however, be linearized by introducing substitution $\hat{\alpha}^{a_1, o}(s') = \pi_1(a_1) \alpha^{a_1, o}(s')$ and $\hat{\lambda}_i^{a_1, o} = \pi_1(a_1) \lambda_i^{a_1, o}$ to obtain (27). \square

Lemma 9.1. *Let $\Upsilon = \{(b_i, y_i) | 1 \leq i \leq k\}$ such that $y_i \geq V^*(b_i)$ for every $1 \leq i \leq k$. Then the value function V_{UB}^Υ is δ -Lipschitz continuous and satisfies*

$$V^* \leq V_{\text{UB}}^\Upsilon \leq V_{\text{HSV11}}^\Upsilon.$$

Proof. The inequality $V_{\text{UB}}^\Upsilon \leq V_{\text{HSV11}}^\Upsilon$ follows trivially from eq. (32) (with $b' := b$). Proving $V^*(b) \leq V_{\text{UB}}^\Upsilon(b)$ is more involved. Suppose that b' is the minimizer from

the definition of V_{UB}^{Υ} , i.e., that $V_{\text{UB}}^{\Upsilon}(b) = V_{\text{HSV11}}^{\Upsilon}(b') + \delta \|b - b'\|_1$. By definition of $V_{\text{HSV11}}^{\Upsilon}$, this b' can be represented as a convex combination $\sum_i \lambda_i b_i = b'$ for which $\sum_i \lambda_i y_i = V_{\text{HSV11}}^{\Upsilon}(b')$. We thus have

$$V_{\text{UB}}^{\Upsilon}(b) = \sum_{i=1}^k \lambda_i y_i + \delta \|b - b'\|_1. \quad (\text{A.16})$$

Our assumptions imply that every pair (b_i, y_i) satisfies $V^*(b_i) \leq y_i$. Combining this observations with the fact that V^* is convex and δ -Lipschitz continuous (Lemma 5.5 and Proposition 5.8), we have

$$\begin{aligned} V^*(b) &\leq V^*(b') + \delta \|b - b'\|_1 = V^*\left(\sum_i \lambda_i b_i\right) + \delta \|b - b'\|_1 \leq \\ &\leq \sum_{i=1}^k \lambda_i V^*(b_i) + \delta \|b - b'\|_1 \leq \sum_{i=1}^k \lambda_i y_i + \delta \|b - b'\|_1 = V_{\text{UB}}^{\Upsilon}(b). \end{aligned}$$

Finally, let us prove that V_{UB}^{Υ} is δ -Lipschitz continuous. Let us consider beliefs $b_1, b_2 \in \Delta(S)$. Without loss of generality, assume that $V_{\text{UB}}^{\Upsilon}(b_1) \geq V_{\text{UB}}^{\Upsilon}(b_2)$. Let $b_{\arg \min}$ be the minimizer of $V_{\text{UB}}^{\Upsilon}(b_2)$, i.e.,

$$b_{\arg \min} = \arg \min_{b'} [V_{\text{HSV11}}^{\Upsilon}(b') + \delta \|b_2 - b'\|_1]. \quad (\text{A.17})$$

By triangle inequality, we have

$$\begin{aligned} V_{\text{UB}}^{\Upsilon}(b_1) &= \\ &= \min_{b' \in \Delta(S)} [V_{\text{HSV11}}^{\Upsilon}(b') + \delta \|b_1 - b'\|_1] \leq V_{\text{HSV11}}^{\Upsilon}(b_{\arg \min}) + \delta \|b_1 - b_{\arg \min}\|_1 \leq \\ &\leq [V_{\text{HSV11}}^{\Upsilon}(b_{\arg \min}) + \delta \|b_2 - b_{\arg \min}\|_1] + \delta \|b_1 - b_2\|_1 = V_{\text{UB}}^{\Upsilon}(b_2) + \delta \|b_1 - b_2\|_1 \end{aligned}$$

which completes the proof. \square

Lemma 9.2. *The lower bound V_{LB}^{Γ} initially satisfies the following conditions, which are subsequently preserved during point-based updates:*

- (1) V_{LB}^{Γ} is δ -Lipschitz continuous.
- (2) V_{LB}^{Γ} is lower bound on V^* .

Proof. Initially, value function V_{LB}^{Γ} satisfies both conditions. Indeed, the set Γ contains only the value $\text{val}^{\sigma_1^{\text{unif}}}$ of the uniform strategy σ_1^{unif} , i.e., $V_{\text{LB}}^{\Gamma}(b) = \text{val}^{\sigma_1^{\text{unif}}}(b)$ for every belief $b \in \Delta(S)$. Value $\text{val}^{\sigma_1^{\text{unif}}}$ is the value for a valid strategy σ_1^{unif} of player 1—hence it is δ -Lipschitz continuous (Lemma 5.7) and lower bounds V^* .

Assume that every α -vector in the set Γ is δ -Lipschitz continuous, and that for each $\alpha \in \Gamma$ there exists strategy $\sigma_1 \in \Sigma_1$ with $\text{val}^{\sigma_1} \geq \alpha$ (which holds also for the initial V_{LB}^{Γ}). Let $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ be the value composition from Equation (35) obtained when performing the point-based update of V_{LB}^{Γ} by

solving $[HV_{\text{LB}}^\Gamma](b)$. We will now show that the refined function $V_{\text{LB}}^{\Gamma'}$ represented by the set $\Gamma' = \Gamma \cup \{\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})\}$ satisfies both properties, and hence any sequence of application of the point-based updates of V_{LB}^Γ preserves the aforementioned properties.

- (1) By Lemma 6.6, $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ is δ -Lipschitz continuous (and thus so is the value function $V_{\text{LB}}^{\Gamma'}$ represented by the set $\Gamma' = \Gamma \cup \{\text{valcomp}(\pi_1, \bar{\alpha})\}$).
- (2) Each α -vector in Γ forms lower bound on the value of some strategy of player 1. Since $\bar{\alpha}^{\text{LB}} \in \Gamma^{A_1 \times O}$, we have that every $\alpha_{a_1, o}$ lower bounds the value of some strategy of player 1. The fact that $\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})$ is also a lower bound follows from Lemma 6.5—and hence every α -vector from the set $\Gamma' = \Gamma \cup \{\text{valcomp}(\pi_1^{\text{LB}}, \bar{\alpha}^{\text{LB}})\}$ is a lower bound on V^* . Hence also $V_{\text{LB}}^{\Gamma'}(b) = \sup_{\alpha \in \Gamma'} \alpha(b) \leq V^*(b)$.

□

Lemma 9.3. *The upper bound V_{UB}^Υ initially satisfies the following conditions, which are subsequently preserved during point-based updates:*

- (1) V_{UB}^Υ is δ -Lipschitz continuous.
- (2) V_{UB}^Υ is an upper bound on V^* .

Proof. V_{UB}^Υ has been defined as a lower δ -Lipschitz envelope of $V_{\text{HSVII}}^\Upsilon$, hence it is δ -Lipschitz continuous (Lemma 9.1). We will therefore focus only on the property (2). Since the upper bound is initialized by a solution of a perfect information variant of the game, we have that $y_i \geq V^*(b_i)$ for every (b_i, y_i) from the initial set Υ (Equation (34)). Hence, applying Lemma 9.1, V_{UB}^Υ is an upper bound on V^* .

We will now show that if $y_i \geq V^*(b_i)$ holds for $(b_i, y_i) \in \Upsilon$ (and V_{UB}^Υ is thus an upper bound on V^*), the application of a point-based update in any belief yields set Υ' such that $y_i \geq V^*(b_i)$ also holds for every $(b_i, y_i) \in \Upsilon'$ —and the resulting value function $V_{\text{UB}}^{\Upsilon'}$ is therefore upper bound on V^* as well. Since $V_{\text{UB}}^\Upsilon \geq V^*$, the utility function of any stage game satisfies $u^{V_{\text{UB}}^\Upsilon, b}(\pi_1, \pi_2) \geq u^{V^*, b}(\pi_1, \pi_2)$ for every $b \in \Delta(S)$, $\pi_1 \in \Pi_1$ and $\pi_2 \in \Pi_2$. This implies that $[HV_{\text{UB}}^\Upsilon](b) \geq [HV^*](b) = V^*(b)$. We already know that $y_i \geq V^*(b_i)$ holds for $(b_i, y_i) \in \Upsilon$, and now we have $[HV_{\text{UB}}^\Upsilon](b) \geq V^*(b)$. Therefore, for every $(b_i, y_i) \in \Upsilon \cup \{(b, [HV_{\text{UB}}^\Upsilon](b))\}$, we have $y_i \geq V^*(b_i)$, and applying the Lemma 9.1, we have that the value function $V_{\text{UB}}^{\Upsilon'}$ is an upper bound on V^* . □

Lemma 9.4. *Let b_t be a belief encountered at t -th recursion level of **Explore** procedure and assume that the corresponding action-observation pair (a_1^*, o^*) (from line 9 of Algorithm 4) satisfies*

$$\mathbb{P}_{b_t, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}}[a_1^*, o^*] \cdot \text{excess}_{t+1}(\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*)) \leq 0. \quad (40)$$

Then $\text{excess}_t(b_t) \leq -2\delta D$ after performing a point-based update at b_t . Furthermore, all beliefs $b'_t \in \Delta(S)$ such that $\|b_t - b'_t\|_1 \leq D$ satisfy $\text{excess}_t(b'_t) \leq 0$.

Proof. Since $V_{\text{LB}}^\Gamma \leq V^* \leq V_{\text{UB}}^\Upsilon$, it holds that $[HV_{\text{LB}}^\Gamma](b_t) \leq [HV_{\text{UB}}^\Upsilon](b_t)$. Applying Lemma 7.5 with $C = \rho(t+1)$ implies that when the beliefs $\tau(b_t, a_1, \pi_2^{\text{LB}}, o)$ satisfy

$$V_{\text{UB}}^\Upsilon(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) - V_{\text{LB}}^\Gamma(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) \leq \rho(t+1),$$

we have $[HV_{\text{UB}}^{\Upsilon}](b_t) - [HV_{\text{LB}}^{\Gamma}](b_t) \leq \gamma\rho(t+1)$. Luckily, this assumption is satisfied in the considered situation — indeed, otherwise there would be some $(a_1, o) \in A_1 \times O$ with

$$V_{\text{UB}}^{\Upsilon}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) - V_{\text{LB}}^{\Gamma}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) > \rho(t+1),$$

i.e., one satisfying $\text{excess}_{t+1}(\tau(b_t, a_1, \pi_2^{\text{LB}}, o)) > 0$, for which $\mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}}[a_1, o] > 0$. This would contradict the assumption

$$\mathbb{P}_{b, \pi_1^{\text{UB}}, \pi_2^{\text{LB}}}[a_1^*, o^*] \cdot \text{excess}_{t+1}(\tau(b_t, a_1^*, \pi_2^{\text{LB}}, o^*)) \leq 0.$$

Now, according to Equation (37), we have $[HV_{\text{UB}}^{\Upsilon}](b_t) - [HV_{\text{LB}}^{\Gamma}](b_t) \leq \gamma\rho(t+1) = \rho(t) - 2\delta D$. It follows that the excess gap after performing the point-based update in b_t satisfies

$$\begin{aligned} \text{excess}_t(b_t) &= V_{\text{UB}}^{\Upsilon}(b_t) - V_{\text{LB}}^{\Gamma}(b_t) - \rho(t) \leq \gamma\rho(t+1) - \rho(t) \\ &= [\rho(t) - \rho(t)] - 2\delta D = -2\delta D, \end{aligned} \quad (\text{A.18})$$

which completes the proof of the first part of the lemma.

Now since the value functions V_{LB}^{Γ} and V_{UB}^{Υ} are δ -Lipschitz continuous (Lemma 9.2 and Lemma 9.3), the difference $V_{\text{UB}}^{\Upsilon} - V_{\text{LB}}^{\Gamma}$ is 2δ -Lipschitz continuous. Thus for every belief $b'_t \in \Delta(S)$ satisfying $\|b_t - b'_t\|_1 \leq D$, we have

$$V_{\text{UB}}^{\Upsilon}(b'_t) - V_{\text{LB}}^{\Gamma}(b'_t) \leq V_{\text{UB}}^{\Upsilon}(b_t) - V_{\text{LB}}^{\Gamma}(b_t) + 2\delta\|b_t - b'_t\|_1 \leq V_{\text{UB}}^{\Upsilon}(b_t) - V_{\text{LB}}^{\Gamma}(b_t) + 2\delta D. \quad (\text{A.19})$$

Now since $\text{excess}_t(b_t) \leq -2\delta D$, we have $\text{excess}_t(b'_t) \leq 0$ which proves the second part of the lemma. \square

Lemma 10.4. *Let V be a value function that is min-justified. Then $V(b) \geq L$.*

Proof. Assume for the contradiction that $V(b) < L$ for some belief $b \in \Delta(S)$. We pick $b = \arg \min_{b' \in \Delta(S)} V(b')$ and denote $\varepsilon = L - V(b)$. Now, using the utility $u^{V,b}$ from Definition 7.4 and using our choice of b , we have

$$\begin{aligned} u^{V,b}(\pi_1, \pi_2) &= \mathbb{E}_{b, \pi_1, \pi_2}[R(s, a_1, a_2)] + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2}[a_1, o] V(\tau(b, a_1, \pi_2, o)) \\ &\geq \underline{r} + \gamma \sum_{a_1, o} \mathbb{P}_{b, \pi_1, \pi_2}[a_1, o] V(b) = \underline{r} + \gamma V(b) = \underline{r} + \gamma(L - \varepsilon) \end{aligned}$$

where \underline{r} is the minimum reward in the game. Since $L = \sum_{t=1}^{\infty} \gamma^{t-1} \underline{r} = \underline{r} + \sum_{t=2}^{\infty} \gamma^{t-1} \underline{r} = \underline{r} + \gamma L$, we also have that $u^{V,b}(\pi_1, \pi_2) \geq L - \gamma\varepsilon$. Therefore it would have to also hold that $[HV](b) = \max_{\pi_1 \in \Pi_1} \min_{\pi_2 \in \Pi_2} u^{V,b}(\pi_1, \pi_2) \geq L - \gamma\varepsilon > L - \varepsilon = V(b)$ which contradicts that V is min-justified. \square

Lemma 10.5. *Let V be a value function that is max-justified by a set of α -vectors Γ . Then for every $\alpha \in \Gamma$ we have $\alpha \leq U$.*

Proof. Let V be max-justified by Γ and let us assume for contradiction that there exists $\alpha \in \Gamma$ and $s \in S$ such that $\alpha(s) > U$. We pick α and s such that $(\alpha, s) = \arg \max_{\alpha \in \Gamma, s \in S} \alpha(s)$ and denote $\varepsilon = \alpha(s) - U$. Using Definition 6.4 and our choice of (α, s) , we get the following for every $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \Gamma^{A_1 \times O}$:

$$\begin{aligned}
\text{valcomp}(\pi_1, \bar{\alpha})(s) &= \\
&= \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1(a_1) \left[R(s, a_1, a_2) + \gamma \sum_{o, s' \in O \times S} T(o, s' \mid s, a_1, a_2) \alpha_{a_1, o}(s') \right] \\
&\leq \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1(a_1) \left[\bar{r} + \gamma \sum_{o, s' \in O \times S} T(o, s' \mid s, a_1, a_2) \alpha(s) \right] \\
&= \min_{a_2 \in A_2} [\bar{r} + \gamma \alpha(s)]
\end{aligned}$$

where $\bar{r} = \max_{(s, a_1, a_2)} R(s, a_1, a_2)$ is the maximum reward in the game. Since $U = \sum_{t=1}^{\infty} \gamma^{t-1} \bar{r} = \bar{r} + \sum_{t=2}^{\infty} \gamma^{t-1} \bar{r} = \bar{r} + \gamma U$, we have the following inequality for every $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \Gamma^{A_1 \times O}$

$$\text{valcomp}(\pi_1, \bar{\alpha})(s) \leq \min_{a_2 \in A_2} [\bar{r} + \gamma \alpha(s)] = \bar{r} + \gamma(U + \varepsilon) = U + \gamma\varepsilon < U + \varepsilon = \alpha(s). \quad (\text{A.20})$$

By Equation (A.20), no value composition can satisfy $\text{valcomp}(\pi_1, \bar{\alpha})(b_s) \geq \alpha(b_s)$ where $b_s(s) = 1$ and $b_s(s') = 0$ otherwise. Consequently, no value composition can satisfy $\text{valcomp}(\pi_1, \bar{\alpha})(b) \geq \alpha(b)$ for every belief $b \in \Delta(S)$ as required by Definition 10.3. This contradicts our assumption and concludes the proof. \square

Lemma 10.6. *Let Γ be a set of linear functions, and V a value function that is max-justified by Γ . Then V is also max-justified by $\text{Conv}(\Gamma)$.*

Proof. Recall that V is max-justified by Ω if 1) $V(b) = \sup_{\alpha \in \Omega} \alpha(b)$ and 2) for every $\alpha \in \Omega$ there exists $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \Omega^{A_1 \times O}$ such that $\text{valcomp}(\pi_1, \bar{\alpha}) \geq \alpha$. Let V be a value function and suppose that Γ satisfies 1) and 2). We will now verify that these properties hold for $\text{Conv}(\Gamma)$ as well. By Proposition 5.11, we have that $\sup_{\alpha \in \text{Conv}(\Gamma)} \alpha(b) = \sup_{\alpha \in \Gamma} \alpha(b)$. Since the property 1) holds for Γ and the value of V remains unchanged, 1) holds for $\text{Conv}(\Gamma)$ as well. We will now prove 2) by showing that for every $\alpha \in \text{Conv}(\Gamma)$, there exists $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \text{Conv}(\Gamma)^{A_1 \times O}$ such that $\text{valcomp}(\pi_1, \bar{\alpha}) \geq \alpha$.

First of all, let us write $\alpha \in \text{Conv}(\Gamma)$ as a finite convex combination $\sum_{i=1}^k \lambda_i \alpha^i$ of α -vectors $\alpha^i \in \Gamma$. Using the assumption that V is max-justified by Γ , we have that for every α^i there exists $\pi_1^{(i)} \in \Pi_1$ and $\bar{\alpha}^{(i)} \in \Gamma$, such that $\text{valcomp}(\pi_1^{(i)}, \bar{\alpha}^{(i)}) \geq \alpha^i$. Denote $\pi_1(a_1) := \sum_{i=1}^k \lambda_i \pi_1^{(i)}(a_1)$ and $\alpha_{a_1, o} := \sum_{i=1}^k \lambda_i \pi_1^{(i)}(a_1) \alpha_{a_1, o}^i / \pi_1(a_1)$.¹¹ We claim that π_1 and $\bar{\alpha}$ witness that V is max-justified by $\text{Conv}(\Gamma)$. Since $\pi_1^{(i)}$ and $\bar{\alpha}^{(i)}$ were chosen s.t. $\text{valcomp}(\pi_1^{(i)}, \bar{\alpha}^{(i)}) \geq \alpha^i$, we have $\sum_{i=1}^k \lambda_i \text{valcomp}(\pi_1^{(i)}, \bar{\alpha}^{(i)}) \geq$

¹¹Observe that $\alpha_{a_1, o} \in \text{Conv}(\Gamma)$ since $\pi_1(a_1) = \sum_{i=1}^k \lambda_i \pi_1^{(i)}(a_1)$ and the coefficients $\lambda_i \pi_1^{(i)}(a_1) / \pi_1(a_1)$ thus sum to 1.

$\sum_{i=1}^k \lambda_i \alpha^i = \alpha \in \text{Conv}(\Gamma)$. To finish the proof, we show that $\text{valcomp}(\pi_1, \bar{\alpha}) \geq \sum_{i=1}^k \lambda_i \text{valcomp}(\pi_1^i, \bar{\alpha}^i)$. By Definition 6.4, we have

$$\begin{aligned}
\text{valcomp}(\pi_1, \bar{\alpha}) &= \\
&= \min_{a_2 \in A_2} \left[\sum_{a_1 \in A_1} \pi_1(a_1) R(s, a_1, a_2) \right. \\
&\quad \left. + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1(a_1) T(o, s' | s, a_1, a_2) \alpha_{a_1, o}(s') \right] \\
&= \min_{a_2 \in A_2} \left[\sum_{i=1}^k \sum_{a_1 \in A_1} \lambda_i \pi_1^i(a_1) R(s, a_1, a_2) + \right. \\
&\quad \left. + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} T(o, s' | s, a_1, a_2) \sum_{i=1}^k \lambda_i \pi_1^i(a_1) \alpha_{a_1, o}^i(s') \right] \\
&= \min_{a_2 \in A_2} \sum_{i=1}^k \lambda_i \left[\sum_{a_1 \in A_1} \pi_1^i(a_1) R(s, a_1, a_2) \right. \\
&\quad \left. + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1^i(a_1) T(o, s' | s, a_1, a_2) \alpha_{a_1, o}^i(s') \right] \\
&\geq \sum_{i=1}^k \lambda_i \min_{a_2 \in A_2} \left[\sum_{a_1 \in A_1} \pi_1^i(a_1) R(s, a_1, a_2) \right. \\
&\quad \left. + \gamma \sum_{(a_1, o, s') \in A_1 \times O \times S} \pi_1^i(a_1) T(o, s' | s, a_1, a_2) \alpha_{a_1, o}^i(s') \right] \\
&= \sum_{i=1}^k \lambda_i \text{valcomp}(\pi_1^i, \bar{\alpha}^i).
\end{aligned}$$

□

Proposition 10.7. *Let V be a value function that is max-justified by a set of α -vectors Γ . Let $b^{\text{init}} \in \Delta(S)$ and $\rho^{\text{init}} \in \Gamma$. By playing according to $\text{Act}(b^{\text{init}}, \rho^{\text{init}})$, player 1 implicitly forms a strategy σ_1 for which $\text{val}^{\sigma_1} \geq \rho^{\text{init}}$.*

Proof. Let b^{init} and ρ^{init} be as in the proposition and assume that player 1 follows $\text{Act}(b, \rho)$ for the first K stages and then follows the uniformly-random strategy σ_1^{unif} . We denote this strategy as $\sigma_1^{b, \rho, K}$. To get to our result, we will first consider an arbitrary belief $b \in \Delta(S)$ and gadget $\rho \in \Gamma$. We will use induction to prove that the value of $\sigma_1^{b, \rho, K}$ satisfies $\text{val}^{\sigma_1^{b, \rho, K}} \geq \rho - \gamma^K \cdot (U - L)$.

First, assume that $K = 0$, i.e., player 1 plays the uniform strategy σ_1^{unif} immediately. Value of the uniform strategy σ_1^{unif} is at least $\text{val}^{\sigma_1^{\text{unif}}} \geq L$ (Proposition 5.3) while $\rho \leq U$ (Lemma 10.5). Hence $\text{val}^{\sigma_1^{b, \rho, 0}} \geq L \geq L - (U - \rho) = \rho - \gamma^0(U - L)$.

Let $K \geq 1$ and assume that $\text{val}^{\sigma_1^{b', \rho', K-1}} \geq \rho' - \gamma^{K-1}(U - L)$ for every belief $b' \in \Delta(S)$ and gadget $\rho' \in \Gamma$. Observe that due to the recursive nature of the **Act** method, we can represent the strategy $\sigma_1^{b, \rho, K}$ as a composite strategy $\sigma_1^{b, \rho, K} = \text{comp}(\pi_1^*, \bar{\zeta})$, where $\zeta_{a_1, o} = \sigma_1^{\tau(b, a_1, \pi_2, o), \alpha_{a_1, o}^*, K-1}$ and π_1^* comes from line 5 of Algorithm 5. (To ensure that $\bar{\alpha}^*$ and π_1^* are correctly defined, the algorithm requires the existence of a value composition satisfying $\text{valcomp}(\pi_1, \bar{\alpha}) \geq \rho$. This requirement holds since V is max-justified by the set Γ and $\rho \in \Gamma$.) Applying Lemma 6.3, the induction hypothesis, and Definition 6.4 (in this order), we have

$$\begin{aligned}
& \text{val}^{\text{comp}(\pi_1^*, \bar{\zeta})}(s) = \\
&= \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1^*(a_1) \left[R(s, a_1, a_2) + \gamma \sum_{(o, s') \in O \times S} T(o, s' \mid s, a_1, a_2) \text{val}^{\zeta_{a_1, o}}(s') \right] \\
&\geq \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1^*(a_1) \left[R(s, a_1, a_2) + \right. \\
&\quad \left. + \gamma \sum_{(o, s') \in O \times S} T(o, s' \mid s, a_1, a_2) [\alpha_{a_1, o}^*(s') - \gamma^{K-1}(U - L)] \right] \\
&= \min_{a_2 \in A_2} \sum_{a_1 \in A_1} \pi_1^*(a_1) \left[R(s, a_1, a_2) + \right. \\
&\quad \left. + \gamma \sum_{(o, s') \in O \times S} T(o, s' \mid s, a_1, a_2) \alpha_{a_1, o}^*(s') \right] - \gamma^K(U - L) \\
&= \text{valcomp}(\pi_1^*, \bar{\alpha}^*) - \gamma^K(U - L).
\end{aligned}$$

We thus have $\text{val}^{\sigma_1^{b, \rho, K}} = \text{val}^{\text{comp}(\pi_1^*, \bar{\zeta})} \geq \text{valcomp}(\pi_1^*, \bar{\alpha}^*) - \gamma^K(U - L)$. Moreover, according to constraint on line 5 of Algorithm 5, we also have $\text{valcomp}(\pi_1^*, \bar{\alpha}^*) \geq \rho$. As a result, we also have $\text{val}^{\sigma_1^{b, \rho, K}} \geq \rho - \gamma^K(U - L)$. This completes the induction step.

Denote by σ_1 the strategy where player 1 follows **Act**($b^{\text{init}}, \rho^{\text{init}}$) for *infinite* period of time (i.e., as $K \rightarrow \infty$). We then have

$$\text{val}^{\sigma_1} = \lim_{K \rightarrow \infty} \text{val}^{\sigma_1^{b^{\text{init}}, \rho^{\text{init}}, K}} \geq \lim_{K \rightarrow \infty} [\rho^{\text{init}} - \gamma^K(U - L)] = \rho^{\text{init}}$$

which completes the proof. \square

Proposition 10.9. *Let V be a min-justified value function and let b^{init} be the initial belief of the game. The Algorithm 6 implicitly constructs a strategy σ_2 which guarantees that the utility to player 1 will be at most $V(b^{\text{init}})$.*

Proof. For the purposes of this proof, we will use

$$\text{val}_2(\sigma'_2, b) = \sup_{\sigma_1 \in \Sigma_1} \mathbb{E}_{b, \sigma_1, \sigma'_2}[\text{Disc}^\gamma]$$

to denote the value a strategy σ'_2 of player 2 guarantees when the belief of player 1 is b . Similarly to the proof of Proposition 10.7, we will first consider

strategies $\sigma_2^{b,K}$ where player 2 plays according to $\mathbf{Act}(b)$ for K steps, and then follows an arbitrary (e.g., uniform) strategy in the rest of the game, and we show that $\text{val}_2(\sigma_2^{b,K}, b) \leq V(b) + \gamma^K(U - L)$.

First, let $K = 0$ and $b \in \Delta(S)$ be the belief of player 1. By Proposition 5.3, player 1 cannot achieve higher utility than U . Moreover, V is min-justified, so we have $V(b) \geq L$ by Lemma 10.4. Therefore, player 1 cannot achieve higher utility than $\text{val}_2(\sigma_2^{b,0}, b) \leq U \leq U + V(b) - L = V(b) + \gamma^0(U - L)$ when his belief is b .

Now let $K \geq 1$ be arbitrary. By the induction hypothesis, we have that strategy $\sigma_2^{b',K-1}$ guarantees that the utility is at most $\text{val}_2(\sigma_2^{b',K-1}, b') \leq V(b') + \gamma^{K-1}(U - L)$ when the belief of player 1 is b' . Let us evaluate the utility that $\sigma_2^{b,K}$ guarantees against arbitrary strategy σ_1 of player 1 in belief b . In the first stage of the game, player 2 plays according to π_2^* obtained on line 3 of Algorithm 6, and the expected reward from the first stage is $\mathbb{E}_{b,\sigma_1,\pi_2^*}[R(s, a_1, a_2)]$. If player 1 plays a_1 and observes o , he reaches an (a_1, o) -subgame where the belief of player 1 is $\tau(b, a_1, \pi_2^*, o)$ and player 2 plays $\sigma_2^{\tau(b, a_1, \pi_2^*, o), K-1}$. Using the induction hypothesis, we know that player 1 is able to achieve utility of at most $\text{val}_2(\sigma_2^{\tau(b, a_1, \pi_2^*, o), K-1}, \tau(b, a_1, \pi_2^*, o)) \leq V(\tau(b, a_1, \pi_2^*, o)) + \gamma^{K-1}(U - L)$. This implies that an upper bound on the utility that σ_1 achieves against $\sigma_2^{b,K}$ (i.e., the strategy corresponding to player 2 following $\mathbf{Act}(b)$ for K stages) is

$$\begin{aligned} & \mathbb{E}_{b,\sigma_1,\pi_2^*}[R(s, a_1, a_2)] + \gamma \mathbb{E}_{b,\sigma_1,\pi_2^*}[V(\tau(b, a_1, \pi_2^*, o)) + \gamma^{K-1}(U - L)] \\ &= \mathbb{E}_{b,\sigma_1,\pi_2^*}[R(s, a_1, a_2)] + \\ & \quad + \gamma \sum_{(a_1, o) \in A_1 \times O} \mathbb{P}_{b,\sigma_1,\pi_2^*}[a_1, o] \cdot [V(\tau(b, a_1, \pi_2^*, o)) + \gamma^{K-1}(U - L)] . \end{aligned}$$

By allowing player 1 to maximize over σ_1 , we get an upper bound on the value $\text{val}_2(\sigma_2^{b,K}, b)$ strategy $\sigma_2^{b,K}$ guarantees when the belief of player 1 is b .

$$\begin{aligned} & \text{val}_2(\sigma_2^{b,K}, b) \leq \\ & \leq \sup_{\sigma_1 \in \Sigma_1} \left[\mathbb{E}_{b,\sigma_1,\pi_2^*}[R(s, a_1, a_2)] + \right. \\ & \quad \left. + \gamma \sum_{(a_1, o) \in A_1 \times O} \mathbb{P}_{b,\sigma_1,\pi_2^*}[a_1, o] \cdot [V(\tau(b, a_1, \pi_2^*, o)) + \gamma^{K-1}(U - L)] \right] \\ &= \max_{\pi_1 \in \Pi_1} \left[\mathbb{E}_{b,\pi_1,\pi_2^*}[R(s, a_1, a_2)] + \right. \\ & \quad \left. + \gamma \sum_{(a_1, o) \in A_1 \times O} \mathbb{P}_{b,\pi_1,\pi_2^*}[a_1, o] \cdot V(\tau(b, a_1, \pi_2^*, o)) \right] + \gamma^K(U - L) \\ &= \max_{\pi_1 \in \Pi_1} u^{V,b}(\pi_1, \pi_2^*) + \gamma^K(U - L) \end{aligned}$$

Using the fact that π_2^* is the optimal strategy in the stage game $[HV](b)$, the

definition of the stage game's value, and the fact that V is min-justified, we get

$$\begin{aligned} \max_{\pi_1 \in \Pi_1} u^{V,b}(\pi_1, \pi_2^*) + \gamma^K(U - L) &= \min_{\pi_2 \in \Pi_2} \max_{\pi_1 \in \Pi_1} u^{V,b}(\pi_1, \pi_2) + \gamma^K(U - L) \\ &= [HV](b) + \gamma^K(U - L) \leq V(b) + \gamma^K(U - L) . \end{aligned}$$

Hence, the utility player 1 with belief b can achieve against player 2 who follows strategy $\sigma_2^{b,K}$ is at most $V(b) + \gamma^K(U - L)$, and we have $\text{val}_2(\sigma_2^{b,K}, b) \leq V(b) + \gamma^K(U - L)$ which completes the induction step.

Now, similarly to the proof of Proposition 10.7, when player 2 follows $\text{Act}(b^{\text{init}})$ for *infinitely* many stages (i.e., plays strategy σ_2 from the theorem), player 1 is able to achieve utility at most

$$\text{val}_2(\sigma_2, b^{\text{init}}) = \lim_{K \rightarrow \infty} \text{val}_2(\sigma_2^{b^{\text{init}},K}, b^{\text{init}}) \leq \lim_{K \rightarrow \infty} [V(b^{\text{init}}) + \gamma^K(U - L)] = V(b^{\text{init}})$$

which completes the proof. \square

Lemma 10.10. *Let Γ be the set of α -vectors that have been generated at any time during the execution of the HSVI algorithm for one-sided POSGs (Algorithm 4). Then the lower bound V_{LB}^Γ is max-justified by the set $\text{Conv}(\Gamma)$.*

Proof. Observe that during the execution of Algorithm 4 the set Γ is modified only by the point-based updates on lines 8 and 12 of Algorithm 4. To prove the result, it thus suffices to show that (1) the initial lower bound V_{LB}^Γ is max-justified by the set $\text{Conv}(\Gamma) = \Gamma = \{\text{val}^{\sigma_1^{\text{unif}}}\}$ and that (2) if V_{LB}^Γ is max-justified by $\text{Conv}(\Gamma)$ then any point-based update results in a value function $V_{\text{LB}}^{\Gamma'}$ that is max-justified by the set $\text{Conv}(\Gamma')$.

First, let us show that the initial lower bound V_{LB}^Γ is max-justified by the initial set of α -vectors $\Gamma = \{\text{val}^{\sigma_1^{\text{unif}}}\}$ (and therefore also by $\text{Conv}(\Gamma) = \Gamma$). Clearly, $\sigma_1^{\text{unif}} = \text{comp}(\pi_1^{\text{unif}}, \zeta^{\text{unif}})$, i.e., the uniform strategy σ_1^{unif} can be composed from a uniform stage strategy π_1^{unif} for the first stage of the game, and playing uniform strategy $\zeta_{a_1,o}^{\text{unif}} = \sigma_1^{\text{unif}}$ in every (a_1, o) -subgame after playing and observing (a_1, o) . Hence, $\text{val}^{\sigma_1^{\text{unif}}} = \text{valcomp}(\pi_1^{\text{unif}}, \bar{\alpha}^{\text{unif}})$ for $\alpha_{a_1,o}^{\text{unif}} = \text{val}^{\sigma_1^{\text{unif}}}$ and the initial V_{LB}^Γ is therefore max-justified by the set $\text{Conv}(\Gamma) = \Gamma = \{\text{val}^{\sigma_1^{\text{unif}}}\}$.

Next, consider a lower bound V_{LB}^Γ from Algorithm 4 and assume that it is max-justified by a set $\text{Conv}(\Gamma)$. The point-based update constructs a set $\Gamma' = \Gamma \cup \{\text{valcomp}(\pi_1, \bar{\alpha})\}$ for some $\pi_1 \in \Pi_1$ and $\bar{\alpha} \in \text{Conv}(\Gamma)^{A_1 \times O}$, see Equation (35). Since V_{LB}^Γ was max-justified by $\text{Conv}(\Gamma)$, we know that for every $\alpha \in \text{Conv}(\Gamma)$ there exists $\pi_1' \in \Pi_1$, $\bar{\alpha}' \in \text{Conv}(\Gamma)^{A_1 \times O}$ such that $\text{valcomp}(\pi_1', \bar{\alpha}') \geq \alpha$. The same holds for the newly constructed α vector $\text{valcomp}(\pi_1, \bar{\alpha})$, and $V_{\text{LB}}^{\Gamma'}$ is therefore max-justified by $\text{Conv}(\Gamma) \cup \{\text{valcomp}(\pi_1, \bar{\alpha})\}$. By Lemma 10.6, we also have that $V_{\text{LB}}^{\Gamma'}$ is max-justified by $\text{Conv}(\text{Conv}(\Gamma) \cup \{\text{valcomp}(\pi_1, \bar{\alpha})\}) = \text{Conv}(\Gamma')$. Every point-based update thus results in a value function $V_{\text{LB}}^{\Gamma'}$ which is max-justified by $\text{Conv}(\Gamma')$ which completes the proof. \square