



**HAL**  
open science

# HSVI Can Solve Zero-Sum Partially Observable Stochastic Games

Aurélien Delage, Olivier Buffet, Jilles Dibangoye, Abdallah Saffidine

► **To cite this version:**

Aurélien Delage, Olivier Buffet, Jilles Dibangoye, Abdallah Saffidine. HSVI Can Solve Zero-Sum Partially Observable Stochastic Games. 2022. hal-04382756v1

**HAL Id: hal-04382756**

**<https://inria.hal.science/hal-04382756v1>**

Preprint submitted on 11 Jan 2023 (v1), last revised 16 Jan 2024 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# HSVI CAN SOLVE ZERO-SUM PARTIALLY OBSERVABLE STOCHASTIC GAMES

---

A PREPRINT

**Aurélien Delage**  
CITI  
INSA Lyon  
Villeurbanne  
aurelien.delage@insa-lyon.fr

**Olivier Buffet**  
INRIA - CNRS  
Université de Lorraine  
Villers-lès-Nancy  
olivier.buffet@inria.fr

**Jilles S. Dibangoye**  
CITI  
INSA Lyon  
Villeurbanne  
jilles-steeve.dibangoye@inria.fr

**Abdallah Saffidine**  
University of New South Wales  
Sidney  
abdallahs@cse.unsw.edu.au

October 27, 2022

## ABSTRACT

State-of-the-art methods for solving 2-player zero-sum imperfect information games rely on linear programming or regret minimization, though not on dynamic programming (DP) or heuristic search (HS), while the latter are often at the core of state-of-the-art solvers for other sequential decision-making problems. In partially observable or collaborative settings (*e.g.*, POMDPs and Dec-POMDPs), DP and HS require introducing an appropriate statistic that induces a fully observable problem as well as bounding (convex) approximators of the optimal value function. This approach has succeeded in some subclasses of 2-player zero-sum partially observable stochastic games (zs-POSGs) as well, but how to apply it in the general case still remains an open question. We answer it by (i) rigorously defining an equivalent game to work with, (ii) proving mathematical properties of the optimal value function that allow deriving bounds that come with solution strategies, (iii) proposing for the first time an HSVI-like solver that provably converges to an  $\epsilon$ -optimal solution in finite time, and (iv) empirically analyzing it. This opens the door to a novel family of promising approaches complementing those relying on linear programming or iterative methods.

## 1 Introduction

Solving imperfect information sequential games is a challenging field with many applications from playing Poker [Kuhn, 1950] to security games [Basilico et al., 2016]. We focus on finite-horizon 2-player zero-sum partially observable stochastic games (zs-POSGs), an important class of games that is equivalent to that of zero-sum extensive-form games (zs-EFGs) [Oliehoek and Vlassis, 2006]<sup>1</sup>. From the viewpoint of (maximizing) player 1, we aim at finding a strategy with a worst-case expected return (*i.e.*, whatever player 2’s strategy) within  $\epsilon$  of the Nash equilibrium value (NEV).

A first approach to solving a zs-POSG is to turn it into a zs-EFG addressed as a *sequence form* linear program (SFLP) [Koller et al., 1996, von Stengel, 1996, Bošanský et al., 2014], giving rise to an exact algorithm. A second approach is to use an iterative game solver, *i.e.*, either a counterfactual-regret-based method (CFR) [Zinkevich et al., 2007, Brown and Sandholm, 2018], or a first-order method [Hoda et al., 2010, Kroer et al., 2020], both coming with asymptotic convergence properties. CFR-based approaches now incorporate deep reinforcement learning and search, some of them winning against top human players at heads-up no limit hold’em poker [Moravčík et al., 2017, Brown and Sandholm, 2018, Brown et al., 2020]. A third approach, proposed by Wiggers [2015], is to use two parallel searches in strategy space, one per player, so that the gap between both strategies’ security levels (*i.e.*, the values of their opponent’s best responses) bounds the distance to the NEV.

---

<sup>1</sup>Note: POSGs are equivalent to the large class of “well-behaved” EFGs as defined by Kovařík et al. [2019].

In contrast, dynamic programming and heuristic search have not been applied to general zs-POSGs, while often at the core of state-of-the-art solvers in other problem classes that involve Markovian dynamics, partial observability and multiple agents (POMDP [Åström, 1965, Smith, 2007], Dec-POMDP [Szer et al., 2005, Dibangoye et al., 2016], or subclasses of zs-POSGs with simplifying observability assumptions [Ghosh et al., 2004, Chatterjee and Doyen, 2014, Basu and Stettner, 2015, Horák et al., 2017, Cole and Kocherlakota, 2001, Horák and Bošanský, 2019]). They all rely on some statistic that induces a fully observable problem whose value function ( $V^*$ ) exhibits continuity properties that allow deriving bounding approximations. Wiggers et al. [2016b,a] progress in this direction for zs-POSGs by demonstrating an important continuity property of the optimal value function, and proposing a reformulation as a particular equivalent game. We work in a similar direction, 1. using a game with different observability hypotheses, 2. proving theoretical results they implicitly rely on, and 3. building on some of their results to derive an HSVI-like algorithm solving the zs-POSG.

Section 2 presents some necessary background, including the concept of *occupancy state* [Dibangoye et al., 2016, Wiggers et al., 2016a] (*i.e.*, the probability distribution over the players' past action-observation histories), and properties that rely on it. Then, Section 3 describes theoretical contributions. First, we rigorously reformulate the problem as a non-observable game, and demonstrate that the Nash equilibrium value can be expressed with a recursive formula, which is a required tool for DP and HS (Sec. 3.1). Second, we exhibit novel continuity properties of optimal value functions and derive bounding approximators, a second tool made necessary due to the continuous state space of the new game, before showing that these approximators come with valid solution strategies for the zs-POSG (Sec. 3.2). Third, we adapt Smith and Simmons' [Smith and Simmons, 2005] HSVI's algorithmic scheme to  $\epsilon$ -optimally solve the problem in finitely many iterations (Sec. 3.3). Section 4 presents an empirical analysis of the approach. Section 5 discusses similarities and differences of our work with CFR-based continual resolving methods before concluding.

## 2 Background

Here, we first give basic definitions about zs-POSGs, including the solution concept at hand. Then we introduce an equivalent game where a state corresponds to a statistic summarizing past behaviors, which leads to some important properties of the game's optimal value.

### 2.1 zs-POSGs

**Definition 2.1** (zs-POSGs). As illustrated through a dynamic influence diagram in Figure 1, a (2-player) zero-sum partially observable stochastic game (zs-POSG) is defined by a tuple  $\langle \mathcal{S}, \mathcal{A}^1, \mathcal{A}^2, \mathcal{Z}^1, \mathcal{Z}^2, P, r, H, \gamma, b_0 \rangle$ , where

- $\mathcal{S}$  is a finite set of states;
- $\mathcal{A}^i$  is (player)  $i$ 's finite set of actions;
- $\mathcal{Z}^i$  is  $i$ 's finite set of observations;
- $P_{a^1, a^2}^{z^1, z^2}(s'|s)$  is the probability to transition to state  $s'$  and receive observations  $z^1$  and  $z^2$  when actions  $a^1$  and  $a^2$  are performed while in state  $s$ ;
- $r(s, a^1, a^2)$  is a (scalar) reward function;
- $H \in \mathbb{N}$  is a (finite) temporal horizon;
- $\gamma \in [0, 1]$  is a discount factor; and
- $b_0$  is the initial belief state, *i.e.*, a probability distribution over states at  $t = 0$ .

From the Dec-POMDP, POSG and EFG literature, we use the following concepts and definitions:

$\theta_\tau^i = (a_0^i, z_1^i, \dots, a_{\tau-1}^i, z_\tau^i)$  is a length- $\tau$  *action-observation history* (AOH) for  $i$ . We note  $\Theta_\tau^i$  the set of all AOHs for player  $i$  at horizon  $\tau$  such that any AOH  $\theta_\tau^i$  is in  $\cup_{t=0}^{H-1} \Theta_t^i$ .

$\beta_\tau^i$  is a (*behavioral*) *decision rule* (DR) at  $\tau$  for  $i$ , *i.e.*, a mapping from private AOHs in  $\Theta_\tau^i$  to *distributions* over private actions.  $\beta_\tau^i(\theta_\tau^i, a^i)$  is the probability to pick  $a^i$  when facing  $\theta_\tau^i$ .

$\beta_{\tau:\tau'}^i = (\beta_\tau^i, \dots, \beta_{\tau'}^i)$  is a *behavioral strategy* for  $i$  from time step  $\tau$  to  $\tau'$  (included).

When considering both players, the last 3 concepts become:

$\theta_\tau = (\theta_\tau^1, \theta_\tau^2) \in \Theta = \cup_{t=0}^{H-1} \Theta_t$ , a *joint* AOH at  $\tau$ ,

$\beta_\tau = \langle \beta_\tau^1, \beta_\tau^2 \rangle \in \mathcal{B} = \cup_{t=0}^{H-1} \mathcal{B}_t$ , a *decision rule profile*, and

$\beta_{\tau:\tau'} = \langle \beta_{\tau:\tau'}^1, \beta_{\tau:\tau'}^2 \rangle$ , a *behavioral strategy profile*.

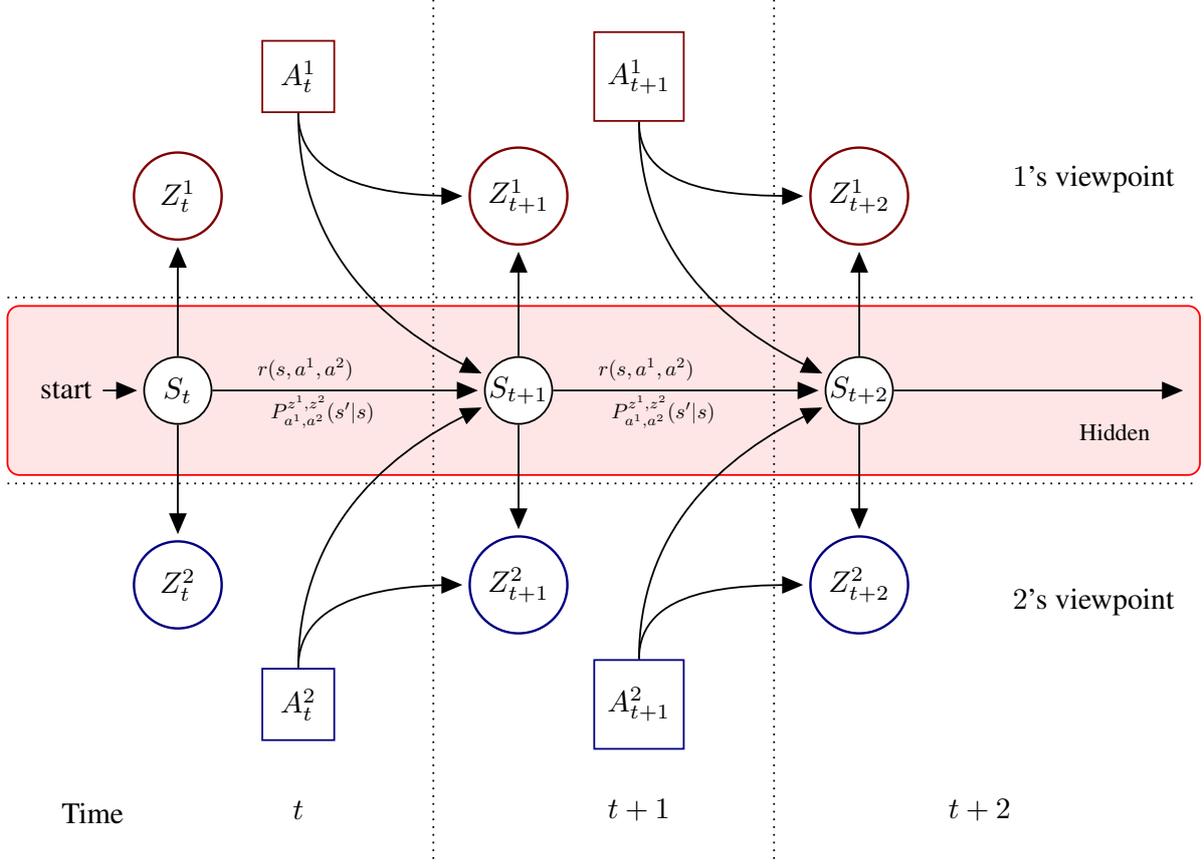


Figure 1: Dynamic influence diagram representing the evolution of a zs-POSG

### Nash Equilibria

Here, player 1 (respectively 2) wants to maximize (resp. minimize) the expected return, or *value*, of strategy profile  $\beta_{0:H-1}$ , defined as the discounted sum of future rewards, *i.e.*,

$$V_0(\beta_{0:H-1}) = E \left[ \sum_{t=0}^{H-1} \gamma^t R_t \mid \beta_{0:H-1} \right],$$

where  $R_t$  is the random variable associated to the instant reward at  $t$ . This leads to the solution concept of Nash equilibrium strategy (NES).

**Definition 2.2** (Nash Equilibrium). The strategy profile  $\beta_{0:H-1}^* = \langle \beta_{0:H-1}^{1*}, \beta_{0:H-1}^{2*} \rangle$  is a NES if no player has an incentive to deviate, which can be written:

$$\begin{aligned} \forall \beta_{0:H-1}^1, V_0(\beta_{0:H-1}^{1*}, \beta_{0:H-1}^{2*}) &\geq V_0(\beta_{0:H-1}^1, \beta_{0:H-1}^{2*}) \text{ and} \\ \forall \beta_{0:H-1}^2, V_0(\beta_{0:H-1}^{1*}, \beta_{0:H-1}^{2*}) &\leq V_0(\beta_{0:H-1}^{1*}, \beta_{0:H-1}^2). \end{aligned}$$

In such a game, all NESs have the same Nash-equilibrium value (NEV),  $V_0^* \stackrel{\text{def}}{=} V_0(\beta_{0:H-1}^{1*}, \beta_{0:H-1}^{2*})$ . Our specific objective is to find an  $\epsilon$ -NES, *i.e.*, a behavioral strategy profile such that any player would gain at most  $\epsilon$  by deviating.

### Why writing a Bellman Optimality Equation is Hard

Our approach requires writing Bellman optimality equations. The main obstacle to achieve this is to find an appropriate characterization of a *subproblem* that allows

1. *predicting* both the immediate reward and the next possible subproblems given an immediate decision;
2. *connecting* a subproblem's solution with solutions of its own (lower-level) subproblems; and

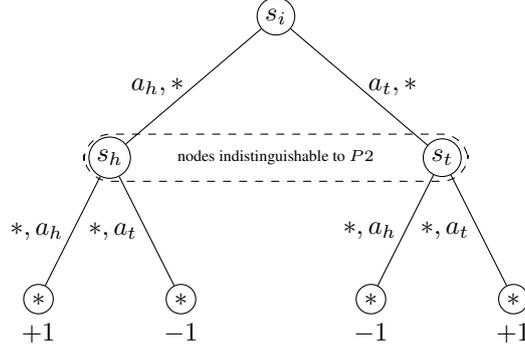


Figure 2: Simplified tree representation of the sequentialized Matching Pennies game. Irrelevant actions, noted  $*$ , allow merging edges with the same action for (i) player 2 at  $t = 0$ , and (ii) player 1 at  $t = 1$ . Notes: (a) Due to irrelevant actions, this game can be seen as an Extensive Form Game, despite players acting simultaneously. (b) Players only know about their past action history (in this observation-free game).

### 3. *prescribing* a solution strategy for the subproblem built on solutions of lower-level subproblems.

In our setting, a player's AOH does not characterize a subproblem since her opponent's strategy is also required to predict the expected reward and the next AOHs. For their part, joint AOHs allow predicting next joint AOHs given both player's immediate decision rules, but would not be appropriate either, since player  $i$  cannot decide how to act when facing some individual AOH  $\theta_i^t$  without considering all possible AOHs of his opponent  $-i$ .

Partial behavioral strategy profiles (sequences of behavioral decision rule profiles from  $t = 0$  to some  $\tau$ ) contain enough information to completely describe the situation at  $\tau$ , and are thus necessarily *predictive*. We still need to demonstrate that they are *connected*, despite decision rules not being public, and *prescriptive*, despite the need to address global-consistency issues illustrated in the following example.

*Example 1.* Matching pennies is a well-known 2-player zero-sum game in which each player has a penny and secretly chooses one side (head or tail). Then, both penny's sides are revealed, and player 1 wins (payoff of  $+1$ ) if both chosen sides match and loses (payoff of  $-1$ ) if not.

We here formalize this game as a zs-POSG (as illustrated in Figure 2) where player 1 actually picks his action at  $t = 0$ , and player 2 at  $t = 1$ . Hence the tuple  $\langle \mathcal{S}, \mathcal{A}^1, \mathcal{A}^2, \mathcal{Z}^1, \mathcal{Z}^2, P, r, H, \gamma, b_0 \rangle$  where:

- $\mathcal{S} = \{s_i, s_h, s_t\}$ , where  $s_i$  is the initial state, and  $s_h$  and  $s_t$  represent a memory of 1's last move: respectively "head" or "tail";
- $\mathcal{A}^1 = \mathcal{A}^2 = \{a_h, a_t\}$  for playing "head" ( $a_h$ ) or "tail" ( $a_t$ );
- $\mathcal{Z}^1 = \mathcal{Z}^2 = \{z_n\}$  a "none" trivial observation;
- $P_a^z(s'|s) = T(s, a, s') \cdot \mathcal{O}(a, s', z)$ , using the next two definitions;
- $T$  is deterministic and such that ( $\cdot$  is used to denote "for all")
  - $T(\cdot, \cdot, a_h) = s_h$ ,
  - $T(\cdot, \cdot, a_t) = s_t$ ;
- $\mathcal{O}$  is deterministic and always returns " $z_n$ ";
- $r$  is such that
  - $r(s_i, \cdot, \cdot) = 0$ ,
  - $r(s_t, \cdot, a_t) = r(s_h, \cdot, a_h) = +1$ ,
  - $r(s_t, \cdot, a_h) = r(s_h, \cdot, a_t) = -1$ ;
- $H = 2$ ;
- $\gamma = 1$ ;
- $b_0$  is such that the initial state is  $s_i$  with probability 1.

Let us then assume that both players' DRs at  $t = 0$  are fixed, with  $\beta_0^1$  randomly picking  $a_t$  or  $a_h$  (i.e., it induces a NES whatever his DR at  $t = 1$ ). Then, we face a "subgame" at  $t = 1$  where any strategy profile  $\langle \beta_{1:1}^1, \beta_{1:1}^2 \rangle$  is a NES

profile with Nash equilibrium value 0. In particular, 2 can pick a deterministic strategy  $\beta_{1:1}^2$ , which will be said to be *locally consistent*. Yet, for 2, such a NES in the subgame at  $\tau = 1$  is not necessarily *globally consistent*, *i.e.*, it may not be part of a NES for the original game (*i.e.*, starting from  $\tau = 0$ ). Intuitively, in such global-consistency issues ?Schmid [2021] (also called *safety issues* Burch et al. [2014]), the choices made at latter time steps do not account for possible deviations from the opponent at earlier time steps.

As detailed in the next section, we will characterize a subproblem not with the raw data of partial strategy profiles, but with a sufficient statistic, and this characterization will be used as the state of a new dynamic game equivalent to the zs-POSG.

## 2.2 Occupancy State and Occupancy Markov Game

We now introduce an equivalent game, in which trajectories correspond to behavioral strategy profiles, and which we will be able to decompose temporally (and recursively), a first key tool for DP and HS.

To cope with the necessarily continuous nature of its state space, we will set this game in occupancy space, *i.e.*, a statistic that sums up past DR profiles. This will let us derive continuity properties on which to build point-based approximators.

As Wiggers et al. [2016a], let us formally define an *occupancy state* (OS)  $\sigma_{\beta_{0:\tau-1}}$  as the probability distribution over joint AOHs  $\theta_\tau$  given partial strategy profile  $\beta_{0:\tau-1}$ . This statistic exhibits the usual Markov and sufficiency properties:

**Proposition 2.3** (Adapted from Dibangoye et al. [Dibangoye et al., 2016, Thm. 1] – Proof in App. B.1).  $\sigma_{\beta_{0:\tau-1}}$ , together with  $\beta_\tau$ , is a sufficient statistic to compute (i) the next OS,  $T(\sigma_{\beta_{0:\tau-1}}, \beta_\tau) \stackrel{\text{def}}{=} \sigma_{\beta_{0:\tau}}$ , and (ii) the expected reward at  $\tau$ :  $r(\sigma_{\beta_{0:\tau-1}}, \beta_\tau) \stackrel{\text{def}}{=} \mathbb{E} [R_\tau | \beta_{0:\tau-1} \oplus \beta_\tau]$ , where  $\oplus$  denotes a concatenation.

Writing from now on  $\sigma_\tau$ , as short for  $\sigma_{\beta_{0:\tau-1}}$ , the OS associated with some prefix strategy profile  $\beta_{0:\tau-1}$ , the proof essentially relies on deriving the following formulas:  $\forall \theta_\tau^1, a^1, z^1, \theta_\tau^2, a^2, z^2$ ,

$$\begin{aligned} T(\sigma_\tau, \beta_\tau)((\theta_\tau^1, a^1, z^1), (\theta_\tau^2, a^2, z^2)) & \quad (1) \\ & \stackrel{\text{def}}{=} Pr((\theta_\tau^1, a^1, z^1), (\theta_\tau^2, a^2, z^2) | \sigma_\tau, \beta_\tau) \\ & = \beta_\tau^1(\theta_\tau^1, a^1) \beta_\tau^2(\theta_\tau^2, a^2) \sigma_\tau(\theta_\tau) \sum_{s, s'} P_a^z(s' | s) b(s | \theta_\tau), \end{aligned}$$

where  $b(s | \theta_\tau)$  is a *belief state* obtained by Hidden Markov Model filtering; and

$$\begin{aligned} r(\sigma_\tau, \beta_\tau) & \stackrel{\text{def}}{=} E[r(S, A^1, A^2) | \sigma_\tau, \beta_\tau] & (2) \\ & = \sum_{s, \theta_\tau, \mathbf{a}} \sigma_\tau(\theta_\tau) b(s | \theta_\tau) \beta_\tau^1(a^1 | \theta_\tau^1) \beta_\tau^2(a^2 | \theta_\tau^2) r(s, \mathbf{a}). \end{aligned}$$

We can then derive, from a zs-POSG, a non-observable zero-sum game similar to Wiggers et al.’s *plan-time NOSG* [Wiggers et al., 2016a, Definition 4], but without assuming that the players’ past strategies are public.

**Definition 2.4** (zero-sum occupancy Markov Game (zs-oMG)). A *zero-sum occupancy Markov game* (zs-oMG)<sup>2</sup> is defined by the tuple  $\langle \mathcal{O}^\sigma, \mathcal{B}, T, r, H, \gamma \rangle$ , where:

- $\mathcal{O}^\sigma (= \cup_{t=0}^{H-1} \mathcal{O}_t^\sigma)$  is the set of OSS induced by the zs-POSG;
- $\mathcal{B}$  is the set of DR profiles of the zs-POSG;
- $T$  is the deterministic transition function in Eq. (1);
- $r$  is the reward function in Eq. (2); and
- $H$  and  $\gamma$  are as in the zs-POSG

( $b_0$  is not in the tuple but serves to define  $T$  and  $r$ ).

In this game, as in the zs-POSG, a player’s solution is a behavioral strategy. Besides, the value of a strategy profile  $\beta_{0:H-1}$  is the same for both zs-oMG and zs-POSG, so that they share the same  $\epsilon$ -NEV and  $\epsilon$ -NESs. We can thus work with zs-oMGs as a means to solve zs-POSGs.

The following aims at deriving a recursive expression of  $V_0^*$ , as well as continuity properties.

<sup>2</sup>We use (i) “Markov game” instead of “stochastic game” because the dynamics are not stochastic, and (ii) “partially observable stochastic game” to stick with the literature.

### Bellman Optimality Equation

Despite the OS at  $\tau > 0$  not being accessible to any player, let us define a *subgame* at  $\sigma_\tau$  as the restriction starting from time step  $\tau$  under this particular occupancy state, meaning that we are seeking strategies  $\beta_{\tau:H-1}^1$  and  $\beta_{\tau:H-1}^2$ .  $\sigma_\tau$  tells us which AOHs each player could be facing with non-zero probability, and are thus relevant for planning. We can then define the value function in any OS  $\sigma_\tau$  for any strategy profile  $\beta_{\tau:H-1}$  as follows:

$$V_\tau(\sigma_\tau, \beta_{\tau:H-1}) \stackrel{\text{def}}{=} E\left[\sum_{t=\tau}^{\infty} \gamma^{t-\tau} r(S_t, A_t) \mid \sigma_\tau, \beta_{\tau:H-1}\right]. \quad (3)$$

The optimal value of a subgame rooted at  $\sigma_\tau$ ,  $V^*(\sigma_\tau)$ , is thus the unique NEV for the previous criterion<sup>3</sup>. Wiggers et al. then proved key continuity properties of  $V^*$  discussed next.

### Concavity and Convexity Results

As a preliminary step, Wiggers et al. decompose an occupancy state  $\sigma_\tau$  into a *marginal term*  $\sigma_\tau^{m,1}$  and a *conditional term*  $\sigma_\tau^{c,1}$ , where

- $\sigma_\tau^{m,1}(\theta_\tau^1) = \sum_{\theta_\tau^2} \sigma_\tau(\theta_\tau^1, \theta_\tau^2)$  is the probability of 1 facing  $\theta_\tau^1$  under  $\sigma_\tau$ , and
- $\sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) = \frac{\sigma_\tau(\theta_\tau^1, \theta_\tau^2)}{\sigma_\tau^{m,1}(\theta_\tau^1)}$  is the probability of 2 facing  $\theta_\tau^2$  under  $\sigma_\tau$  given that 1 faces  $\theta_\tau^1$ ,

so that  $\sigma_\tau(\theta_\tau^1, \theta_\tau^2) = \sigma_\tau^{m,1}(\theta_\tau^1) \cdot \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1)$ . (Symmetric definitions apply by swapping players 1 and 2.) In addition, let us denote  $T_m^1(\sigma_\tau, \beta_\tau)$  and  $T_c^1(\sigma_\tau, \beta_\tau)$  the marginal and conditional terms associated to  $T(\sigma_\tau, \beta_\tau)$ .

Now, if 1 faces AOH  $\theta_\tau^1$ , knows 2's future strategy  $\beta_{\tau:H-1}^2$ , and has access to  $\sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1)$  for any  $\theta_\tau^2$ , then she faces a POMDP whose optimal value we denote  $\nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2(\theta_\tau^1)$ . This leads to defining the best-response value vector  $\nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2$ , which contains one component per AOH  $\theta_\tau^1$ , and writing the value of 1's best response against  $\beta_{\tau:H-1}^2$  under  $\sigma_\tau$  as  $\sigma_\tau^{m,1} \cdot \nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2$ . But then, because 2 also knows  $\sigma_\tau$ , she can in fact pick  $\beta_{\tau:H-1}^2$  to minimize this value, so that we get the following theorem.

**Theorem 2.5** ([Wiggers et al., 2016a, Thm. 2]). *For any  $\tau \in \{0 \dots H-1\}$ ,  $V_\tau^*$  is (i) concave w.r.t.  $\sigma_\tau^{m,1}$  for a fixed  $\sigma_\tau^{c,1}$ , and (ii) convex w.r.t.  $\sigma_\tau^{c,2}$  for a fixed  $\sigma_\tau^{c,1}$ . More precisely,*

$$V_\tau^*(\sigma_\tau) = \min_{\beta_{\tau:H-1}^2} \left[ \sigma_\tau^{m,1} \cdot \nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2 \right] = \max_{\beta_{\tau:H-1}^1} \left[ \sigma_\tau^{m,2} \cdot \nu_{[\sigma_\tau^{c,2}, \beta_{\tau:H-1}^1]}^1 \right], \text{ where}$$

$$\nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2(\theta_\tau^1) \stackrel{\text{def}}{=} \max_{\beta_{\tau:H-1}^1} \mathbb{E}_{\theta_\tau^2 \sim \sigma_\tau^{c,1}(\theta_\tau^1)} \left\{ \sum_{t=\tau}^{H-1} \gamma^{t-\tau} r(S_t, A_t^1, A_t^2) \mid \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2 \right\}.$$

*Proof.* (Sketch) We start from von Neumann's Minimax theorem [von Neumann, 1928] giving the following equation:

$$\begin{aligned} V_\tau^*(\sigma_\tau) &= \min_{\beta_{\tau:H-1}^2} \max_{\beta_{\tau:H-1}^1} [V_\tau(\sigma_\tau, \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2)] \\ &= \min_{\beta_{\tau:H-1}^2} \max_{\beta_{\tau:H-1}^1} \left[ \mathbb{E} \left\{ \sum_{t=\tau}^{H-1} \gamma^{t-\tau} r(S_t, A_t^1, A_t^2) \mid \theta_\tau^1, \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2, \sigma_\tau^{c,1} \right\} \right], \end{aligned}$$

then, observing that 1's best response to  $\beta_{\tau:H-1}^2$  can be computed for each AOH  $\theta_\tau^1$  independently, we can swap the max operator and part of the expectation one ( $\mathbb{E}$ ) as follows:<sup>4</sup>

$$= \min_{\beta_{\tau:H-1}^2} \mathbb{E}_{\theta_\tau^1 \sim \sigma_\tau^{m,1}} \left\{ \underbrace{\max_{\beta_{\tau:H-1}^1} \mathbb{E}_{\theta_\tau^2 \sim \sigma_\tau^{c,1}(\theta_\tau^1)} \left[ \sum_{t=\tau}^{H-1} \gamma^{t-\tau} r(S_t, A_t^1, A_t^2) \mid \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2 \right]}_{\text{best-response of 1 to } \beta_{\tau:H-1}^2 \text{ under } \theta_\tau^1} \right\}$$

and, recognizing the components of vector  $\nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2$  and writing the expectation over AOHs  $\theta_\tau^1$  as a scalar product:

$$= \min_{\beta_{\tau:H-1}^2} \left[ \sigma_\tau^{m,1} \cdot \nu_{[\sigma_\tau^{c,1}, \beta_{\tau:H-1}^2]}^2 \right].$$

□

<sup>3</sup>We will come back to the validity of this point in Section 3.1.

An important observation that ensues from this theorem is that  $V_\tau^*$  is concave in  $\sigma_\tau^{m,1}$  and convex in  $\sigma_\tau^{m,2}$ . In practice, however, such continuity properties alone only allow upper-bounding  $V_\tau^*$  for finitely many conditional terms  $\sigma_\tau^{c,i}$ , thus *not* for the whole occupancy space, as required to enable DP and HS in our game.

In the following, we complement Wiggers et al.’s results with properties of  $V^*$  in subgames, plus continuity properties that help designing bounding approximators, which will lead us to an HSVI-like solver.

**Note:** To help the reader, Appendix A provides two synthetic tables: Table 3 (p. 23) to sum up various theoretical properties that are stated in this paper (assuming a finite temporal horizon), and Table 4 (p. 24) to sum up the notations used in this paper, including some that are used only in the appendix.

Also, for convenience, we may replace in the following: (i) subscript “ $\tau : H - 1$ ” with “ $\tau$  :”, (ii) any function  $f(\mathbf{x})$  linear in vector  $\mathbf{x}$  with either  $f(\cdot) \cdot \mathbf{x}$  or  $\mathbf{x}^\top \cdot f(\cdot)$ , (iii) a full tuple with its few elements of interest, and (iv) an element (a “field”)  $x$  of a specific tuple  $t$  by  $x[t]$ .

### 3 Theoretical Contributions

In this section, we demonstrate how to implement dynamic programming and heuristic search by (1) rigorously showing that Bellman optimality equation (Sec. 3.1) holds, (2) deriving bounding approximators of two novel optimal value functions, which come with solution strategies (Sec. 3.2), and (3) proposing a variant of HSVI that computes (in finite time) a player’s strategy whose value is within  $\epsilon$  of the zs-POSG’s NEV (Sec. 3.3).

#### 3.1 The Optimal Value Function $V^*$ and its Recursive Expression

Let us first recall that, contrary to Wiggers et al. [Wiggers et al., 2016a, Section 5, Lemma 4], we do not make the strong assumption that past decision rules can be considered as public (and, thus, we do not assume that any player knows  $\sigma_\tau$ ). Indeed, while it is valid in Dec-POMDPs because the players are willing to coordinate their behaviors, it is *a priori* not valid in zs-POSGs, since players are, in the contrary, willing to deceive one another. Safety issues as presented in Example 1 illustrate the possible flaws of such an assumption.

We now discuss the existence of an optimal value function  $V_\tau^*$  and its properties. These results are implicitly used by Wiggers et al., but it seems important to state and demonstrate them. A first step is to demonstrate that von Neumann’s minimax theorem [von Neumann, 1928] applies when in  $\sigma_\tau$ , thus justifying the definition of the optimal (Nash equilibrium) value of a subgame.

**Theorem 3.1** (Minimax theorem – Proof in App. C.1.2). *The subgame defined in Eq. (3) admits a unique NEV*

$$V_\tau^*(\sigma_\tau) \stackrel{\text{def}}{=} \max_{\beta_{\tau:H-1}^1} \min_{\beta_{\tau:H-1}^2} V_\tau(\sigma_\tau, \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2). \quad (4)$$

$V(\sigma_\tau, \cdot, \cdot)$  not being bilinear in the space of behavioral strategies (Appendix C.1.1), the proof requires reasoning with mixed strategies (for which the bilinearity holds), *i.e.*, distributions over pure (deterministic) strategies defined over *all* time steps. Yet, when in a subgame, we have to reason only on mixed strategies *compatible* with the associated occupancy state  $\sigma_\tau$  (*i.e.*, which ensure that the OS at  $\tau$  is  $\sigma_\tau$ ), one step being to extend Kuhn’s equivalence results between behavioral and mixed strategies [Kuhn, 1950] to the subgames.

Then, defining the optimal action-value function:

$$Q_\tau^*(\sigma_\tau, \beta_\tau) \stackrel{\text{def}}{=} r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau)), \quad (5)$$

we can now prove that a Bellman optimality equation exists, which justifies reasoning on subgames despite the non-observability.

**Theorem 3.2** (Bellman optimality equation – Proof in App. C.1.2).  *$V_\tau^*(\sigma_\tau)$  satisfies the following functional equation:*

$$V_\tau^*(\sigma_\tau) = \max_{\beta_\tau^1} \min_{\beta_\tau^2} r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau)) = \max_{\beta_\tau^1} \min_{\beta_\tau^2} Q_\tau^*(\sigma_\tau, \beta_\tau).$$

The proof requires decomposing min and max operators over different time steps before swapping them appropriately to end up recognizing the optimal value function at the next time step ( $V_{\tau+1}^*$ ).

Theorems 3.1 and 3.2 together show that Theorem 2.5 holds even without player’s strategies being public so that we can now build on the convex-concave property to solve zs-oMGs.

<sup>4</sup>Note that this property is well known in Bayesian games, where AOHs correspond to *types*, *cf.* [Harsanyi, 1968, Th. 1, p. 321].

### 3.2 Towards Solving zs-OMGs

This section aims at providing the second tool for DP and HS with continuous state spaces, *i.e.*, bounding approximators of optimal value functions which will allow generalization across occupancy space. Their update and selection operators are written as linear programs, and they turn out to come with solution strategies.

#### 3.2.1 Bounding value functions

So far, several issues prevented to apply the HSVI scheme to zs-POSGs, starting with the continuous spaces of 1. occupancy states (zs-OMG states) and 2. decision rules (zs-OMG actions). One can address (1) by introducing the bounding functions  $\bar{V}_\tau(\sigma_\tau)$  and  $\underline{V}_\tau(\sigma_\tau)$  of  $V_\tau^*(\sigma_\tau)$  (*cf.* App. D.2), for instance:

$$\bar{V}_\tau(\sigma_\tau) = \min_{\langle \tilde{\sigma}_\tau^{c,1}, \langle \bar{v}_\tau^2, \beta_\tau^2 \rangle \rangle \in \bar{\mathcal{I}}_\tau} [\sigma_\tau^{m,1} \cdot \bar{v}_\tau^2 + \lambda_\tau \|\sigma_\tau - \sigma_\tau^{m,1} \tilde{\sigma}_\tau^{c,1}\|_1],$$

where  $\bar{v}_\tau^2$  component-wise upper-bounds  $\nu_{[\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2]}^2$  for some  $\beta_\tau^2$ . They allow generalizing knowledge from the subgame rooted at  $\sigma_\tau$  to any other one rooted at  $\tilde{\sigma}_\tau$ . To do so, we use  $V^*$ 's Lipschitz-Continuity proven below.

**Theorem 3.3** (Lipschitz-Continuity of  $V^*$  - proof in App. D.1.3). *Let  $h_\tau \stackrel{\text{def}}{=} \frac{1-\gamma^{H-\tau}}{1-\gamma}$  (or  $h_\tau \stackrel{\text{def}}{=} H - \tau$  if  $\gamma = 1$ ). Then  $V_\tau^*(\sigma_\tau)$  is  $\lambda_\tau$ -Lipschitz continuous in  $\sigma_\tau$  at any depth  $\tau \in \{0 \dots H - 1\}$ , where  $\lambda_\tau = \frac{1}{2}h_\tau(r_{\max} - r_{\min})$ .*

Yet, this yields (generally non-convex) Lipschitz-continuous functions whose max-min optimization would be intractable, so that (2) remains an issue. Also, we do not know how to retrieve valid solution strategies. In particular, and as illustrated in Example 1, simply concatenating decision rules backwards from  $\tau = H - 1$  to 0 would not guarantee globally-consistent solutions, and could result in exploitable strategies.

But then, combining Theorems 2.5 and 3.2 leads to introducing a novel value function (denoted  $W_\tau^{1,*}$ ) through writing, for any OS  $\sigma_\tau$ :

$$V_\tau^*(\sigma_\tau) = \max_{\beta_\tau^1} \min_{\beta_{\tau:H-1}^2 \in \mathcal{B}_\tau^2} \underbrace{\left[ r(\sigma_\tau, \beta_\tau) + \gamma \sigma_{\tau+1}^{m,1} \cdot \nu_{[\sigma_{\tau+1}^{c,1}, \beta_{\tau+1:H-1}^2]}^2 \right]}_{\stackrel{\text{def}}{=} W_\tau^{1,*}(\sigma_\tau, \beta_\tau^1)}.$$

Assuming that player 2 can only respond with one of finitely many stored strategies, the concavity and  $\lambda_\tau$ -Lipschitz-continuity of  $W_\tau^{1,*}$  allow upper-bounding it with finitely many tuples  $w = \langle \tilde{\sigma}_\tau, \beta_\tau^2, \langle \bar{v}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle$  stored in sets  $\bar{\mathcal{I}}_\tau$ , and where  $\bar{v}_{\tau+1}^2$  upper-bounds  $\nu_{[\tilde{\sigma}_{\tau+1}^{c,1}, \beta_{\tau+1}^2]}^2$ .

**Proposition 3.4** (proof in App. D.2.2). *Let  $\bar{\mathcal{I}}_\tau$  be a set of tuples  $w = \langle \tilde{\sigma}_\tau, \beta_\tau^2, \langle \bar{v}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle$ . Then,*

$$\begin{aligned} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1) \stackrel{\text{def}}{=} & \min_{\langle \tilde{\sigma}_\tau, \beta_\tau^2, \langle \bar{v}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle \in \bar{\mathcal{I}}_\tau} \left[ r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) \cdot \bar{v}_{\tau+1}^2 \right. \\ & \left. + \lambda_{\tau+1} \|T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) - T_m^1(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) T_c^1(\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1 \right] \end{aligned} \quad (6)$$

upper-bounds  $W_\tau^{1,*}$  over the whole space  $\mathcal{O}_\tau^\sigma \times \mathcal{B}_\tau^1$ .

Symmetrically, we define  $\underline{W}_\tau$  as the lower bound of the symmetrically defined  $W_\tau^{2,*}$ . As explained in the next two sections,  $\bar{W}_\tau$  will be easier to deal with compared to  $\bar{V}_\tau$ , allowing 1 to seek for decision rules optimistically, and providing valid solution strategies for 2 for the subgame at  $\tau$ , *i.e.*, ignoring consistency with higher-level subgames.

#### 3.2.2 Action Selection and Backup Operators

We now detail the decision rule selection for 1 using  $\bar{W}_\tau$  to optimistically guide a trajectory in occupancy space, and how to update  $\bar{W}_\tau$  by providing backup operations.

To that end, first note that linearities in  $\beta_\tau^1$  within Eq. (6) allow writing  $\bar{W}_\tau(\sigma_\tau, \beta_\tau^1) = \min_{w \in \bar{\mathcal{I}}_\tau} \beta_\tau^{1,T} \cdot M_{(\cdot, w)}^{\sigma_\tau}$ , where  $\beta_\tau^1$  and  $M_{(\cdot, w)}^{\sigma_\tau}$  (for each  $w$ ) are column vectors of dimension  $|\Theta^1 \times \mathcal{A}^1|$ .  $M^{\sigma_\tau}$  (see developed formula in App. D.3.1) is thus a  $|\Theta_\tau^1 \times \mathcal{A}^1| \times |\bar{\mathcal{I}}_\tau|$  matrix. Then,  $\bar{W}_\tau$  being a lower envelope of hyperplanes leads to a convenient way of computing  $\max_{\beta_\tau^1} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1)$ .

**Proposition 3.5** (Proof in App. D.3.1). *For any given  $\sigma_\tau$  and any set  $\bar{\mathcal{I}}_\tau$  of tuples  $w = \langle \bar{\sigma}_\tau, \beta_\tau^2, \langle \bar{v}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle$ ,  $\max_{\beta_\tau^1} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1)$  is equivalent to the LP and dual LP:*

$$\begin{aligned}
 \text{LP } \bar{W}_\tau(\sigma_\tau) : \quad & \max_{\beta_\tau^1, v} v \quad \text{s.t.} \quad (i) \quad \forall w \in \bar{\mathcal{I}}_\tau, \quad v \leq \beta_\tau^{1\top} \cdot M_{(\cdot, w)}^{\sigma_\tau} \quad \text{and} \\
 & (ii) \quad \forall \theta_\tau^1 \in \Theta_\tau^1, \quad \sum_{a^1} \beta_\tau^1(a^1 | \theta_\tau^1) = 1, \\
 \text{DLP } \bar{W}_\tau(\sigma_\tau) : \quad & \min_{\psi_\tau^2, v} v \quad \text{s.t.} \quad (i) \quad \forall (\theta_\tau^1, a^1), \quad v \geq M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2 \quad \text{and} \\
 & (ii) \quad \sum_{w \in \bar{\mathcal{I}}_\tau} \psi_\tau^2(w) = 1.
 \end{aligned} \tag{7}$$

**Remark 3.6** (Outcomes of this game). Since  $\bar{W}_\tau$  upper-bounds  $W_\tau^{1,*}$ , solving this LP provides 1 with an *optimistically* selected immediate decision rule  $\beta_\tau^1$ . For 2,  $\psi_\tau^2$  is a probability distribution over tuples containing strategies  $\beta_\tau^2 \oplus \beta_{\tau+1:H-1}^2$ , thus recursively induces a strategy, as illustrated by Fig. 3, which can be turned into a behavioral strategy  $\beta_{\tau:H-1}^2$  (more details in App. D.3.3) whose value is *at worst* (from 2's viewpoint) the LP's value, *i.e.*, against 1's best response to it.

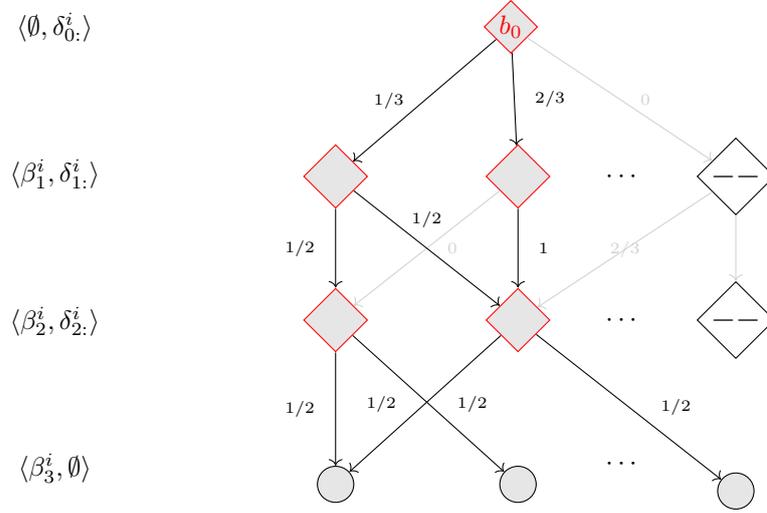


Figure 3: Representation of the strategy recursively induced by some  $\psi_0^1$ . At each time step  $\tau$ , one must (i) sample a next tuple/node  $w_\tau^1$  from current distribution  $\psi_\tau^1$ , (ii) apply DR  $\beta_\tau^1[w_\tau^1]$ , and (iii) make  $\psi_{\tau+1}^1[w_\tau^1]$  the new current distribution (unless reaching a leaf).

Then, the following properties allow performing backups, *i.e.*, filling up the set  $\bar{\mathcal{I}}_{\tau-1}$  with new tuples  $w$  containing, in particular, vectors  $\bar{v}_\tau^2$ .

**Lemma 1** (Proof in App. D.3.2). *For any  $\psi_\tau^2 = \text{DLP } \bar{W}_\tau(\sigma_\tau)$ , the vector  $v_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2$  is component-wise upper-bounded by*

$$\bar{v}_\tau^2 \stackrel{\text{def}}{=} \frac{1}{\sigma_\tau^{m,1}} M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2.$$

**Proposition 3.7** (update). *Let us assume that*

- a transition  $\sigma_{\tau-1} \rightarrow \sigma_\tau$  has been performed through playing  $\langle \beta_{\tau-1}^1, \beta_{\tau-1}^2 \rangle$ , and
- solving  $\text{DLP } \bar{W}_\tau(\sigma_\tau)$  provides both
  - a tree strategy  $\psi_\tau^2$  (as the main solution of the DLP), and
  - a vector  $\bar{v}_\tau^2 = \frac{1}{\sigma_\tau^{m,1}} M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2$  (as a by-product).

Then,

1.  $\bar{\mathcal{I}}_{\tau-1} \leftarrow \bar{\mathcal{I}}_{\tau-1} \cup \{ \langle \sigma_{\tau-1}^{c,1}, \beta_{\tau-1}^2, \langle \bar{v}_\tau^2, \psi_\tau^2 \rangle \rangle \}$  is a valid update operator in the sense that it preserves  $\bar{W}_\tau$ 's upper-bounding property, and
2. similarly,  $\bar{\mathcal{J}}_\tau \leftarrow \bar{\mathcal{J}}_\tau \cup \{ \langle \sigma_\tau^{c,1}, \langle \bar{v}_\tau^2, \psi_\tau^2 \rangle \rangle \}$  is a valid update operator for  $\bar{V}_\tau$ .

### 3.2.3 Initialization

To initialize the bounds  $\overline{W}_\tau$  and  $\overline{V}_\tau$  for any time step, we begin by generating a trajectory in a forward phase. At each time step, a uniform decision rule is picked for both players to derive a sequence of occupancy states  $\sigma_0, \dots, \sigma_{H-1}$ . Then, during a backward phase, for each time step  $\tau = H-1, \dots, 1$ , we create a tuple  $w_{\tau-1, init} = \langle \sigma_{\tau-1}^{c,1}, \beta_{\tau-1}^2, \langle \overline{\nu}_\tau^2, \psi_\tau^2 \rangle \rangle$ , where

- $\sigma_{\tau-1}^{c,1}$  is the conditional term associated to  $\sigma_{\tau-1}$ ;
- $\beta_{\tau-1}^2$  is a uniform decision rule;
- $\psi_\tau^2$  is
  - a degenerate distribution over the only next tuple  $w_{\tau+1}$  if  $\tau < H-1$  (which induces a concatenation of uniform decision rules for all future time steps);
  - undefined if  $\tau = H-1$ ;

and

- $\overline{\nu}_\tau^2(\theta_\tau^1) = r_{max} \cdot (H - \tau)$  for any history  $\theta_\tau^1$  that player 1 could face.

Tuples  $w_{\tau-1, init}$  are added to sets  $\overline{\mathcal{I}}_{\tau-1}$ . For any time step  $\tau \geq 0$ , we similarly create tuples  $\langle \sigma_\tau^{c,1}, \langle \overline{\nu}_\tau^2, \psi_\tau^2 \rangle \rangle$  and add them to sets  $\overline{\mathcal{I}}_\tau$ . The lower bounds are initialized symmetrically.

We now show that occupancy states can also be prescriptive, allowing one to retrieve an  $\epsilon$ -NES for the subgame at occupancy state  $\sigma_\tau$  once the bounds are within  $\epsilon$  from each other, in particular at  $\tau = 0$ .

### 3.2.4 Extracting a NES

Vectors  $\overline{\nu}_0^2$  upper bounding the value of their associated strategies, the following result tells when and how to extract an  $\epsilon$ -optimal solution strategy for this player.

**Theorem 3.8.** *If sets  $\overline{\mathcal{I}}_0$  and  $\underline{\mathcal{I}}_0$  are such that  $\overline{V}_0(\sigma_0) - \underline{V}_0(\sigma_0) \leq \epsilon$ , then  $\arg \max_{w \in \underline{\mathcal{I}}_0} \underline{\nu}_0^2$  and  $\arg \min_{\overline{w} \in \overline{\mathcal{I}}_0} \overline{\nu}_0^2$  respectively provide strategies  $\psi_0^1$  and  $\psi_0^2$  that form an  $\epsilon$ -NES of the zs-POSG.*

*Proof.* First, let us notice that, at  $\tau = 0$ , the occupancy-state space is reduced to a singleton,  $\{\sigma_0 = \langle 1 \rangle\}$ , because of the single (empty) joint AOH. The value vectors  $\nu$  are thus one-dimensional, and here considered as scalar numbers.

Let us assume that sets  $\overline{\mathcal{I}}_0$  and  $\underline{\mathcal{I}}_0$  are such that

$$\overline{V}_0(\sigma_0) - \underline{V}_0(\sigma_0) \leq \epsilon,$$

and let  $\underline{w}^* = \langle \sigma_0^{c,1}, \langle \underline{\nu}_0^*, \psi_0^{1,*} \rangle \rangle$  and  $\overline{w}^* = \langle \sigma_0^{c,1}, \langle \overline{\nu}_0^*, \psi_0^{2,*} \rangle \rangle$  be the tuples returned by  $\arg \max_{w \in \underline{\mathcal{I}}_0} \underline{\nu}_0^2$  and  $\arg \min_{\overline{w} \in \overline{\mathcal{I}}_0} \overline{\nu}_0^1$ . Then, noting that  $\sigma_0 = \langle 1 \rangle$ ,

$$\begin{aligned} \nu_{[\sigma_0^{c,2}, \psi_0^{1,*}]}^1 - \nu_{[\sigma_0^{c,1}, \psi_0^{2,*}]}^2 &\leq \overline{\nu}_0^* - \underline{\nu}_0^* \\ &= \max_{w \in \underline{\mathcal{I}}_0} \underline{\nu}_0^2 - \min_{\overline{w} \in \overline{\mathcal{I}}_0} \overline{\nu}_0^1 \\ &= \overline{V}_0(\sigma_0) - \underline{V}_0(\sigma_0) \\ &\leq \epsilon. \end{aligned}$$

Thus,  $\psi_0^1$  and  $\psi_0^2$  are two strategies whose security levels (values against best-responding opponents) are  $\epsilon$ -close, and thus form an  $\epsilon$ -NES of the zs-POSG.  $\square$

Note: This result can be generalized to any  $\sigma_\tau$  at later time steps, but this generalization is not used in practice.

Distributions  $\psi_0^2$  are stored and could be executed as is. Appendix D.3.3 still presents a conversion process to retrieve a behavioral strategy  $\beta_{0:H-1}^2$  from a distribution  $\psi_0^2$  over tuples  $w \in \overline{\mathcal{I}}_0$ . Next, we see how to design a practical HSVI-based algorithm that provably returns sets  $\overline{\mathcal{I}}_0$  and  $\underline{\mathcal{I}}_0$  satisfying Theorem 3.8 after finitely many iterations.

## 3.3 HSVI for zs-POSGs

This section details our adaptation of the general HSVI scheme for  $\epsilon$ -optimally solving zs-POSGs, and presents a theoretical finite-time convergence property.

### 3.3.1 Algorithm

HSVI for zs-POSGs is described in Algorithm 1. As vanilla HSVI, it relies on (i) generating trajectories while acting optimistically (lines 10+11), *i.e.*, player 1 (resp. 2) acting “greedily” w.r.t.  $\overline{W}_\tau$  (resp.  $\underline{W}_\tau$ ), and (ii) locally updating the upper and lower bounds (lines 17+18). Both phases rely on solving the same games described by LP (7). At  $\tau = H - 1$ , line 14 selects DRs by solving an exact game, and line 20 returns a distribution reduced to the single element added in line 15.

A key difference with Smith and Simmons’s HSVI algorithm [Smith and Simmons, 2005] lies in the criterion for stopping trajectories. The branching factor for zs-OMGs being infinite, we make use of  $V^*$ ’s Lipschitz-continuity to implement the same adaptations as Horák et al. [2017] used for zs-OS-POSGs. The Lipschitz-continuity allows controlling the variations of the value function within small balls of radius  $\rho$  around a previously visited occupancy-state. A finite number of such balls is sufficient to cover the whole space. Then, Theorem 3.9 (below) ensures  $\epsilon$ -optimality in finite time if stopping trajectories when  $\overline{V}_\tau(\sigma_\tau) - \underline{V}_\tau(\sigma_\tau) \leq \text{thr}(\tau)$ , with the threshold function  $\text{thr}(\tau) \stackrel{\text{def}}{=} \gamma^{-\tau} \epsilon - \sum_{i=1}^{\tau} 2\rho\lambda_{\tau-i}\gamma^{-i}$ .

---

**Algorithm 1:** zs-OMG-HSVI( $b_0, [\epsilon, \rho]$ )

[here returning a tuple  $w_0$  containing a solution strategy  $\psi_0^1$  for player 1]

---

```

1 Fct zs-OMG-HSVI( $b_0 \simeq \sigma_0$ )
2   foreach  $\tau \in 0 \dots H - 1$  do
3     Initialize  $\overline{V}_\tau, \underline{V}_\tau, \overline{W}_\tau$ , &  $\underline{W}_\tau$ 
4     while  $[\overline{V}_0(\sigma_0) - \underline{V}_0(\sigma_0) > \text{thr}(0)]$  do
5       Explore( $\sigma_0, 0, -, -$ )
6     return  $\arg \max_{w_0 \in \mathcal{J}_0} v_0^1$ 

7 Fct Explore( $\sigma_\tau, \tau, \sigma_{\tau-1}, \beta_{\tau-1}$ )
8   if  $[\overline{V}_\tau(\sigma_\tau) - \underline{V}_\tau(\sigma_\tau) > \text{thr}(\tau)]$  then
9     if  $\tau < H - 1$  then
10       $\overline{\beta}_\tau^1 \leftarrow \text{LP} \overline{W}_\tau(\sigma)$ 
11       $\underline{\beta}_\tau^2 \leftarrow \text{LP} \underline{W}_\tau(\sigma)$ 
12      Explore( $T(\sigma_\tau, \overline{\beta}_\tau^1, \underline{\beta}_\tau^2), \tau + 1, \sigma_\tau, \langle \overline{\beta}_\tau^1, \underline{\beta}_\tau^2 \rangle$ )
13    else ( $\tau = H - 1$ )
14       $(\overline{\beta}_\tau^1, \underline{\beta}_\tau^2) \leftarrow \text{NES}(r(\sigma, \beta_\tau^1, \beta_\tau^2))$ 
15       $\overline{\mathcal{I}}_\tau^1 \leftarrow \overline{\mathcal{I}}_\tau \cup \{\langle \sigma_\tau^{c,1}, \beta_\tau^2, - \rangle\}$ 
16       $\underline{\mathcal{I}}_\tau^2 \leftarrow \underline{\mathcal{I}}_\tau \cup \{\langle \sigma_\tau^{c,2}, \overline{\beta}_\tau^1, - \rangle\}$ 
17      Update( $\overline{W}_{\tau-1}, \langle \sigma_\tau, \sigma_{\tau-1}^{c,1}, \underline{\beta}_{\tau-1}^2 \rangle$ )
18      Update( $\underline{W}_{\tau-1}, \langle \sigma_\tau, \sigma_{\tau-1}^{c,2}, \overline{\beta}_{\tau-1}^1 \rangle$ )

19 Fct Update( $\overline{W}_{\tau-1}, \langle \sigma_\tau, \sigma_{\tau-1}^{c,1}, \underline{\beta}_{\tau-1}^2 \rangle$ )
20    $\langle \overline{v}_\tau^2, \psi_\tau^2 \rangle \leftarrow \text{DLP} \overline{W}_\tau(\sigma_\tau)$ 
21    $\overline{\mathcal{I}}_{\tau-1} \leftarrow \overline{\mathcal{I}}_{\tau-1} \cup \{\langle \sigma_{\tau-1}^{c,1}, \beta_{\tau-1}^2, \langle \overline{v}_\tau^2, \psi_\tau^2 \rangle \rangle\}$ 
22    $\overline{\mathcal{J}}_\tau \leftarrow \overline{\mathcal{J}}_\tau \cup \{\langle \sigma_\tau^{c,1}, \langle \overline{v}_\tau^2, \psi_\tau^2 \rangle \rangle\}$ 

```

---

**Setting  $\rho$**  As can be observed, this threshold function should always return positive values, which requires a small enough (but  $> 0$ )  $\rho$ . For a given problem (*cf.* Prop. E.1, App. E.1.1), the maximum possible value  $\rho_{\max}$  depends on the Lipschitz constants at each time step, which themselves depend on the initial upper and lower bounds of the optimal value function. Setting  $\rho \in (0, \rho_{\max})$  means making a trade-off between generating many trajectories (small  $\rho$ ) and long ones (large  $\rho$ ).

### 3.3.2 Finite-Time Convergence

**Theorem 3.9** (Proof in App. E.2.1). *zs-OMG-HSVI (Alg. 1) terminates in finite time with an  $\epsilon$ -approximation of  $V_0^*(\sigma_0)$  that satisfies Theorem 3.8.*

The finite time complexity suffers from the same combinatorial explosion as for Dec-POMDPs, and is even worse as we have to handle “infinitely branching” trees of possible futures. More precisely, the bound on the number of iterations depends on the number of balls of radius  $\rho$  required to cover occupancy simplexes at each depth.

Also, the following proposition allows solving infinite horizon problems as well (when  $\gamma < 1$ ) by bounding the length of HSVI’s trajectories using the boundedness of  $\bar{V} - \underline{V}$  and the exponential growth of  $\text{thr}(\tau)$ .

**Proposition 3.10** (Proof in App. E.2.2). *When  $\gamma < 1$ , the length of trajectories is upper bounded by  $T_{\max} \stackrel{\text{def}}{=} \left\lceil \log_{\gamma} \frac{\epsilon - \frac{2\rho\lambda^{\infty}}{1-\gamma}}{W - \frac{2\rho\lambda^{\infty}}{1-\gamma}} \right\rceil$ , where  $\lambda^{\infty}$  is a depth-independent Lipschitz constant and  $W \stackrel{\text{def}}{=} \|\bar{V}^{(0)} - \underline{V}^{(0)}\|_{\infty}$  is the maximum width between initializations.*

## 4 Experiments

Experiments presented in this section aim at validating the proposed approach and comparing its behavior to the behavior of some reference algorithms.

### 4.1 Setup

#### Benchmark Problems

Five benchmark problems were used. Adversarial Tiger and Competitive Tiger were introduced by Wiggers [2015]. Mabc and Recycling Robot are well-known Dec-POMDP benchmark problems (*cf.* <http://masplan.org>) and were adapted to our competitive setting by making player 2 minimize (rather than maximize) the objective function. The fifth benchmark is the adaptation of the well-known Matching Pennies game detailed in Example 1, with a small difference in that  $r(s_h, \cdot, a_h) = +2$  instead of  $+1$ ; this change breaks the symmetry in the optimal strategy, so that HSVI can not find the NES by “chance” by trying uniform strategies. We only consider finite horizons  $H$  and  $\gamma = 1$ . Table 1 gives the cardinal of the state, action and observation sets for each of these problems.

Table 1: Number of states/actions/observations for each benchmark problem

	$\mathcal{S}$	$\mathcal{A}^1$	$\mathcal{A}^2$	$\mathcal{O}^1$	$\mathcal{O}^2$
Competitive Tiger	2	4	4	3	3
Adversarial Tiger	2	3	2	2	2
Recycling Robot	4	3	3	2	2
Mabc	4	2	2	2	2
Matching Pennies	3	2	2	1	1

#### Algorithms

For conciseness, Algorithm 1 is here denoted HSVI, and compared against *Random* search and *Informed* search Wiggers [2015] (both using Wiggers’s implementation (unlicensed and unreleased)), SFLP [Koller et al., 1996], and CFR+ Tammelin [2014] (both using `open_spiel` [Lanctot et al., 2019] (Apache license)).

All algorithms (but SFLP, which is exact) used a target error  $\epsilon = 1\%$  of the initial gap  $H \cdot (r_{\max} - r_{\min})$ . HSVI ran with  $\lambda_{\tau} = (H - \tau) \cdot (r_{\max} - r_{\min})$ , and  $\rho$  the middle of its feasible interval. We also use FB-HSVI’s LPE lossless compression of probabilistically equivalent action-observation histories in occupancy states, so as to reduce their dimensionality [Dibangoye et al., 2016]. Experiments ran on an Ubuntu machine with i7-10810U 1.10 GHz Intel processor and 16 GB available RAM, and the code is available under MIT license at <https://gitlab.com/aureliendelage1/hsviforzspogs>.

Random and Informed, only ran once, providing fairly representative results.

### 4.2 Results

**Performance Measures** A common performance measure in 2-player zero-sum games is the *exploitability* of a strategy  $\beta_0^i$ , *i.e.*, the difference between the strategy’s *security level* (the value of  $-i$ ’s best response to  $\beta_0^i$ ) and the Nash equilibrium value  $V_0^*(\sigma_0)$ :

$$\begin{aligned} \text{exploitability}(\beta_0^i) &= |V^*(\sigma_0) - \sigma_0^{m,1} \cdot \nu_{[\sigma_0^c, \neg i, \beta_0^i]}^i| \\ &= |V^*(\sigma_0) - \nu_{[\sigma_0^c, \neg i, \beta_0^i]}^i|, \end{aligned}$$

Table 2: Comparison of different solvers on various benchmark problems. Reported values are the running times until the algorithm’s error gap (based on bounds for HSVI) is lower than 1 %, or, if the timeout limit is reached, the security-level gap percentages (100 % if  $\text{gap} = H \cdot (R_{\max} - R_{\min})$ ). Notes: (1) Horizons with a star exponent ( $H^*$ ) are those for which the security-level computations ran out of time so that, for HSVI, we give the gap between the pessimistic bounds. (2) Even though Random and Informed contain randomness, we ran them only once, getting fairly representative results.

Domain	H	Wiggers		HSVI	SFLP	CFR+
		Rand.	Inf.			
Comp Tiger	2	2.6 %	8.3 %	6 s	1 s	18 s
	3	7.0 %	6.1 %	3.8 %	48 s	30 m
	4	12.1 %	7.7 %	4.8 %	14 m	[OOT]
	5*	[OOT]	[OOT]	53.3 %	[OOT]	[OOT]
Rec. Robot	2	3.4 %	5.1 %	5 s	1 s	30 s
	3	9.2 %	15.2 %	4 m	1 s	13 m
	4	14.1 %	19.6 %	4.9 %	13 s	1.5 %
	5	[OOT]	[OOT]	<b>10.6981981981982 %</b>	[OOT]	[OOT]
	6*	[OOT]	[OOT]	<b>45.51426426426426 %</b>	[OOT]	[OOT]
Adv Tiger	2	1 s	3.7 %	1 s	1 s	1 s
	3	1.5 %	4.4 %	2 m	1 s	8 s
	4	2.9 %	5.6 %	2.6 %	8 s	13 m
MABC	2	45 s	18.8 %	8 s	3 s	5 s
	3	4.2 %	9.2 %	27 s	1 s	1 m
	4	18.1 %	36.3 %	4.4 %	3 s	47 m
MP	4	2 m	46.7 %	5 s	1 s	2 s
	5	9 m	45.8 %	1 m	1 s	10 s
	6	2.2 %	44.6 %	8 m	2 s	1 m

noting that  $\sigma_0$  is a degenerate distribution over a single element, the pair of empty action-observation histories. In our setting, it will be convenient to look at the (average) *exploitability of a strategy profile*  $\langle \beta_{0:\cdot}^1, \beta_{0:\cdot}^2 \rangle$ :

$$\begin{aligned} \text{exploitability}(\beta_{0:\cdot}^1, \beta_{0:\cdot}^2) &= \frac{(V^*(\sigma_0) - \nu_{[\sigma_0^{c,2}, \beta_{0:\cdot}^1]}^1) + (\nu_{[\sigma_0^{c,1}, \beta_{0:\cdot}^2]}^2 - V^*(\sigma_0))}{2} \\ &= \frac{\nu_{[\sigma_0^{c,1}, \beta_{0:\cdot}^2]}^2 - \nu_{[\sigma_0^{c,2}, \beta_{0:\cdot}^1]}^1}{2}. \end{aligned}$$

This quantity is a more concise statistic than both individual exploitabilities, and can be obtained by solving two POMDPs (fixing one player’s strategy or the other) without requiring to know the actual NEV.

This exploitability can also be defined as half of the *gap between security levels* (SL-gap). To analyze the convergence of algorithms with respect to the initial gap, we will look at the *SL-gap percentage*, *i.e.*,

$$\begin{aligned} \text{SL-gap percentage}(\beta_{0:\cdot}^1, \beta_{0:\cdot}^2) &= \frac{\nu_{[\sigma_0^{c,1}, \beta_{0:\cdot}^2]}^2 - \nu_{[\sigma_0^{c,2}, \beta_{0:\cdot}^1]}^1}{H \cdot (R_{\max} - R_{\min})} \\ &= \frac{2 \cdot \text{exploitability}(\beta_{0:\cdot}^1, \beta_{0:\cdot}^2)}{H \cdot (R_{\max} - R_{\min})}. \end{aligned}$$

#### 4.2.1 Comparison with the state of the art

Table 2 gives the convergence time of Wiggers’s two heuristic algorithms, CFR+, SFLP, and HSVI on the benchmark problems with various horizons, or the SL-gap percentage when reaching a 1 h time limit. Executions not returning any result (*i.e.*, for Random, Informed and CFR+, not performing a single iteration) are noted out-of-time [OOT].

This table first shows that HSVI always outperforms the heuristic baseline provided by Wiggers’s algorithms, thus proving the interest of an HSVI scheme. However, HSVI is outperformed by both SFLP and CFR+, unless they run out of time. As can be noted, HSVI is able to keep improving even when the horizon grows thanks to the LPE compression, taking advantage of underlying structure in some games (*e.g.*, Recycling Robot, a problem with transition+observation independence (TOI), when scaling to larger horizons).

We now study the dynamic behavior of the algorithms at hand by providing and analyzing the bounds and exploitability graphs for the same benchmarks.

### 4.2.2 Bounding Graphs

Left-side graphs in Figures 4 to 9 show how the computed upper- and lower-bounding values  $\bar{V}_0(\sigma_0)$  and  $\underline{V}_0(\sigma_0)$  (respectively the *dotted* dark and light green curves) evolve as a function of computation time (always given in seconds). The *solid* dark and light green curves show the security levels  $\nu_{[\sigma_0, \psi_0^1]}^1$  and  $\nu_{[\sigma_0, \psi_0^2]}^2$  of the current returned strategies  $\psi_0^2$  and  $\psi_0^1$ . Note that, when best-response computations to obtain security levels are expensive (*e.g.*, for the competitive tiger problem, with  $H = 4$ ), they are performed either periodically (*e.g.* every 10 iterations) or only once, at the end. In the captions, we indicate the (arbitrary) frequency of the POMDP evaluations. For example,  $(1, 1, \textit{once})$  means that, for the first two horizons, the POMDP evaluations were done after each iteration, and, for the last one, only once (at the end).

Overall, we observe consistent curves with (i) security levels in-between bounds and around the NEV, and (ii) bounds converging monotonically. Note that HSVI stops when the gap between bounds is small enough, while the gap between SLs (used by Informed, Random and CFR, and whose computation can be time-consuming) can be much smaller. As a matter of fact, one can notice that strategies  $\psi_0^i$  returned at each iteration by HSVI are often better (in terms of security level) than their pessimistic lower- or upper-bounding guarantees  $\underline{\nu}_{[\sigma_0, \psi_0^1]}^1$  and  $\bar{\nu}_{[\sigma_0, \psi_0^2]}^2$ .

### 4.2.3 Exploitability Graphs

Right-side graphs in Figures 4 to 9 show the exploitability of the returned strategy profile as a function of computation time for HSVI, Random, Informed, and CFR+ for the different benchmarks considered. A limit precision of  $10^{-7}$  (chosen empirically, according to the LP solver’s precision) was applied to HSVI’s exploitability.

As can be observed, Random and Informed tend to produce reasonable strategies quickly, but struggle to improve them so as to converge towards an  $\epsilon$ -NES with  $\epsilon \simeq 0$ . In contrast, our algorithm keeps improving as computation time increases. The exploitation graphs support the observed behavior in Table 2 that HSVI converges in reasonable time compared to Wiggers’s algorithms. However, the graphs also show that CFR+ essentially outperforms HSVI when the problems are difficult enough (*i.e.*, when the temporal horizon grows) but the traversal of the whole tree still remains tractable (thus allowing CFR+ to perform iterations). An interesting observation is that, on small enough problems, HSVI achieves very low exploitabilities earlier than CFR+.

Finally, HSVI’s exploitability graph shares strong similarities with those of Bořanský et al.’s double-oracle algorithms [Bořanský et al., 2014, Fig. 8 and 11]. This can be understood as HSVI iteratively building two sets of strategies, one per player, until they are sufficient to support NES profiles, so that the average exploitability is almost zero. But note that Bořanský et al. construct LPs using pure strategies (deterministic best responses), while HSVI’s strategies are stochastic.

Having empirically studied the behavior of HSVI compared to other basic offline solvers, we now provide insight about the connections between HSVI and continual (thus online) resolving methods.

## 5 Comparison with Continual Resolving

Continual Resolving techniques share some similarities with our approach, but also important differences. The purpose of this section is to clarify these points. It starts with a quick description of CFR, on which Continual Resolving is built.

**Counterfactual Regret Minimization (CFR)** Zinkevich et al. [2007] belongs to the self-play family of algorithms, which gave rise to several CFR-based approaches [Tammelin, 2014, Burch et al., 2019, Lanctot et al., 2009, Brown et al., 2017]). It iteratively traverses the whole game tree and applies, in each private history, a regret-matching update rule based on a specific type of regret called counterfactual regrets. Iteratively updating an initially uniform strategy asymptotically converges towards a NES. However, the tree traversal becomes intractable when the tree size is large.

Built on top of the CFR framework, approaches based on *limited-lookahead continual resolving* (LLCR) (inspired by Burch et al.’s decomposition Burch et al. [2014]) such as DeepStack Moravčík et al. [2017], Libratus Brown and Sandholm [2018], ReBeL Brown et al. [2020] and Player of Games Schmid et al. [2021], perform well by exploiting a temporal decomposition in subgames, which are specified through knowledge about both players’ past strategies. Our approach thus shares similarities with these works.

Yet, as we will see in the next sections, a closer look at LLCR [Schmid, 2021] demonstrates how fundamentally different they are, starting with the fact that LLCR is an *online* search algorithm, *i.e.*, is meant to make good decisions

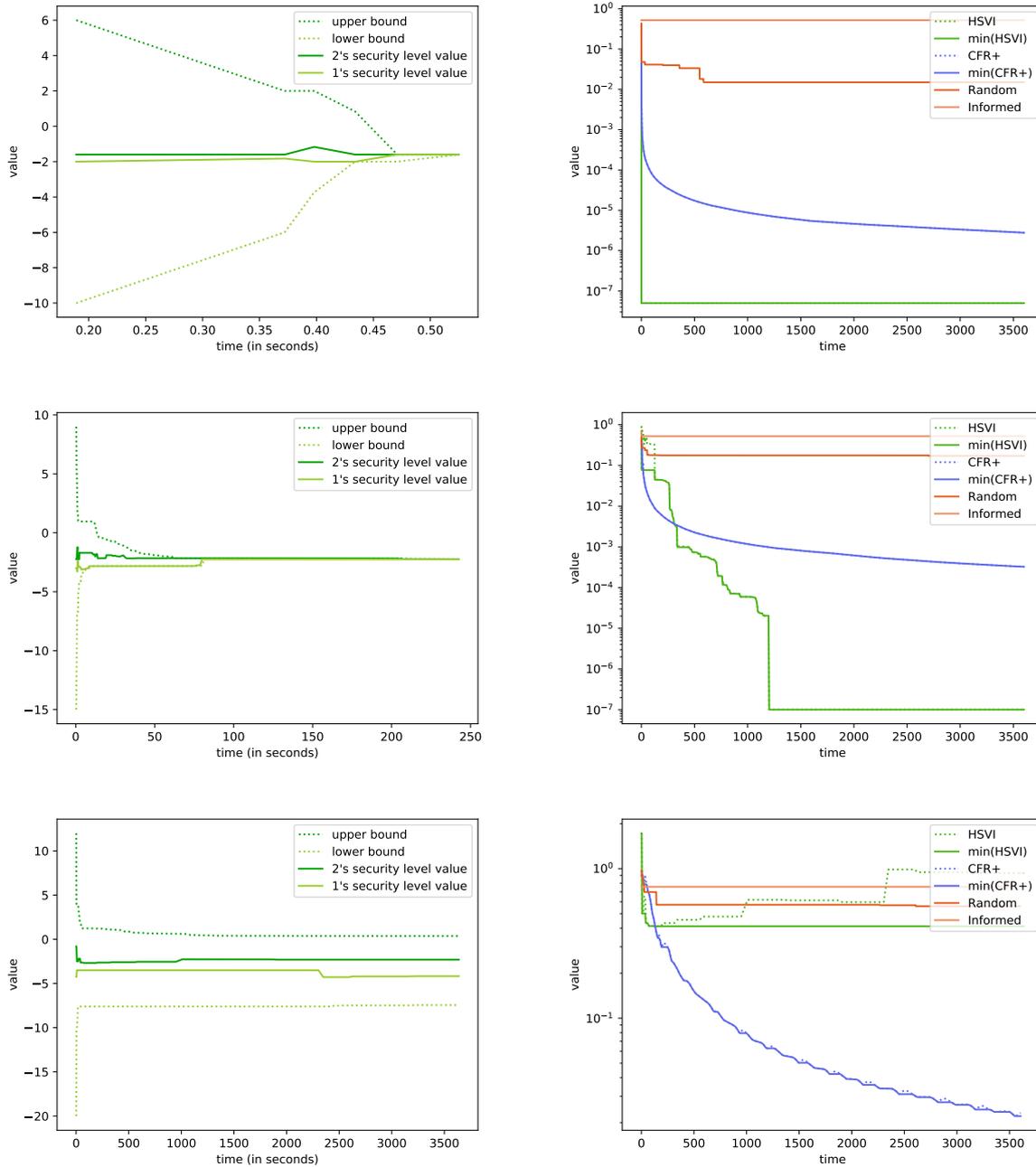


Figure 4: **Adversarial Tiger** ( $H = 2, 3, 4$ ) **(1,1,10)**: **(left)** Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms). **(right)** Exploitability ( $= \frac{SL\text{-gap}}{2}$ ) as a function of time (s) for Random, Informed, CFR+, and HSVI.

at each time step, based on the current knowledge about the state of the game, while HSVI, as SFLP or CFR (and its variants), is an *offline* algorithm returning a complete solution strategy.

### 5.1 Continual Resolving

(Continual) Resolving techniques have been the first ingredient to adapt CFR to online settings. They address the problem of solving the complete subgame (down to its end) starting in the current situation at  $\tau$ , while maintaining the global consistency (aka safety) of the whole strategy, *i.e.*, not making choices that could encourage the opponent to deviate *in the past*, before  $\tau$ . This is achieved by introducing constraints, called *gadgets*, in a preliminary stage of the subgame that represent possible deviations and their values, but increase the size of the game tree so that

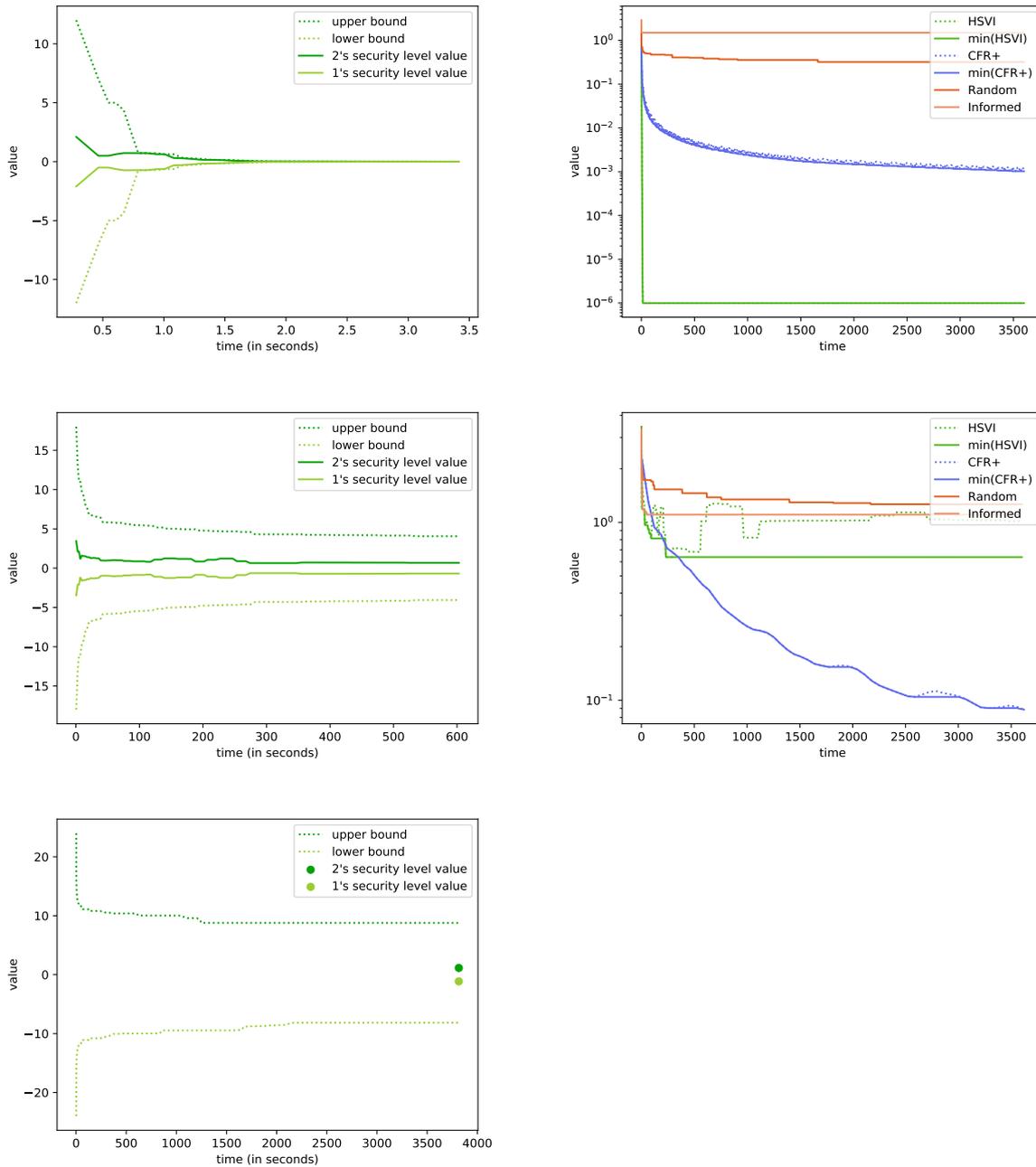


Figure 5: **Competitive Tiger** ( $H = 2, 3$ ) **(1,1)**: **(left)** Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms). **(right)** Exploitability ( $= \frac{SL\text{-gap}}{2}$ ) as a function of time (s) for Random, Informed, CFR+, and HSVI.

it is practically intractable Moravčík et al. [2017]. For its part, HSVI solves similar subgames, but ensuring only local consistency, *i.e.*, only considering the subgame. Global consistency comes from the way HSVI combines "lower-level" solutions in its backtracking process, without adding any gadget modifying the game.

For both Resolving and HSVI, solving a subgame requires *sufficient statistics* that represent a prefix strategy profile from  $t = 0$  to  $\tau$ . In Resolving's online setting, this may seem surprising, since the opponent's actual strategy is not public. Yet, Resolving does not actually require knowing or guessing the opponent's actual strategy. In Resolving, any opponent Nash equilibrium strategy is appropriate, since the purpose is to verify that the opponent has no incentive to deviate from the Nash equilibrium. A requirement for Resolving is for the sufficient statistics to represent *complete* strategies (given the current public information), so that decisions are anticipated for player  $i$  even in AOHs (infostates) not reachable given player  $\neg i$ 's strategy. This leads to using *ranges* Kovařík et al. [2019], a vector that

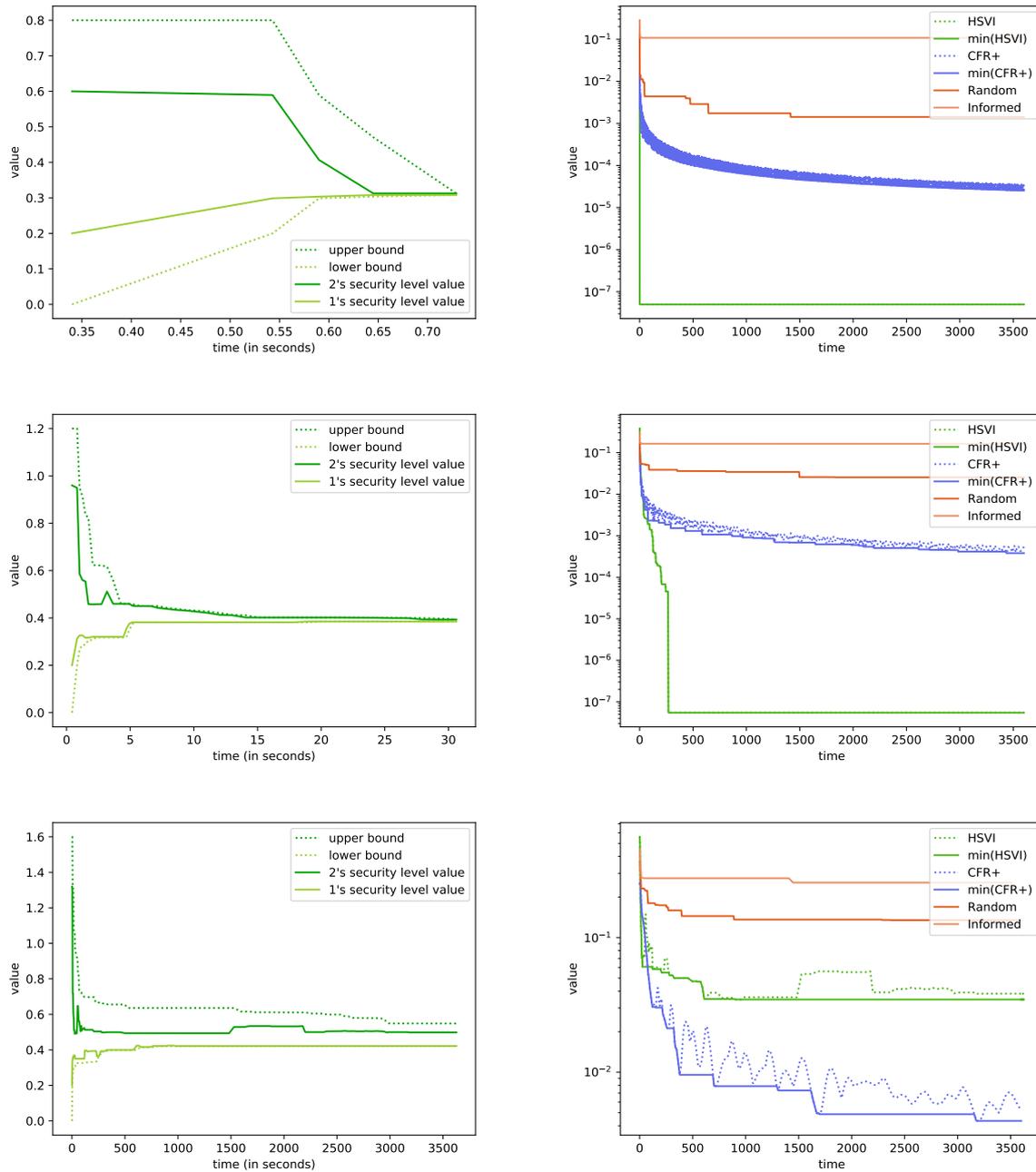


Figure 6: **Mabc** ( $H = 2, 3, 4$ ) **(1,1,10)**: (left) Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms). (right) Exploitability ( $= \frac{SL\_gap}{2}$ ) as a function of time (s) for Random, Informed, CFR+, and HSVI.

gives, for each player, her contribution to the probability of any history she could face at time step  $\tau$ . In contrast, HSVI's occupancy state is not necessarily related to a Nash equilibrium strategy in any manner, and leads to ignoring unreachable AOHs, which helps to reduce the size of the decision-making (sub)problem.

## 5.2 Limited Lookahead

Continual Resolving alone solves complete subgames, thus larger problems at early stages of the game than at the end, which is not appropriate in an online setting. To address this issue through limiting the lookahead of subgames, one needs to estimate the value of the leaves of any truncated subgame. This is achieved through learning offline, for each player  $i$ , deep networks that, given the current public belief state, map each AOH to its value under some Nash

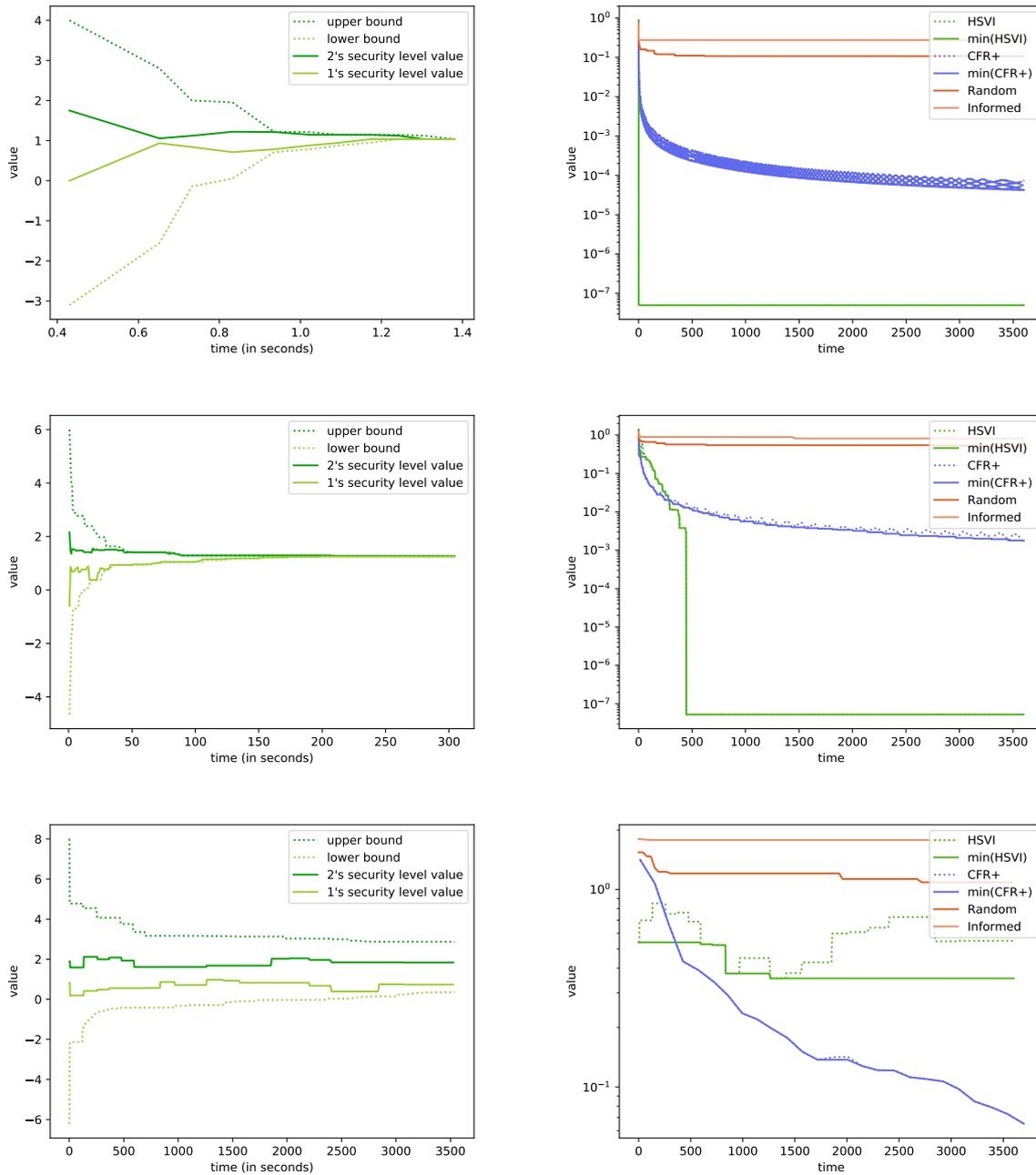


Figure 7: **Recycling Robot** ( $H = 2, 3, 4$ ) **(1,1,10)**: **(left)** Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms). **(right)** Exploitability ( $= \frac{SL\text{-}gap}{2}$ ) as a function of time (s) for Random, Informed, CFR+, and HSVI.

equilibrium strategy profile. Note that the target function is not unique [Kovařík et al., 2019, Proposition A.1], since each NES profile maps to different value vectors. Still, according to Kovařík et al. [2019], this does not seem to cause problems in practice.

In contrast, the individual value functions HSVI considers (the " $\nu$ " functions) are uniquely defined since they correspond to the best responses to given (not necessarily Nash equilibrium) strategies of the opponent.

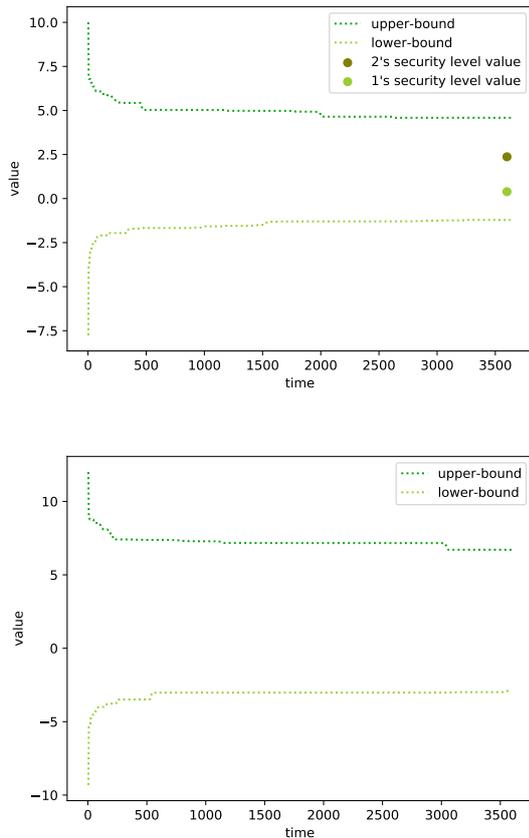


Figure 8: **Recycling Robot** ( $H = 5, 6$ ) (**once,none**): (left) Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms).

### 5.3 Limited Lookahead Continual Resolving as a General Scheme?

The previous subsections highlight to what extent HSVI and LLCR are fundamentally different, in particular because they are not on the same algorithmic level. LLCR should be seen as a general scheme in which the subgame solver used, namely CFR, could be replaced by other "basic" offline algorithms such as HSVI or SFLP. But we leave further investigation on this topic for future work.

## 6 Discussion

This paper addresses the problem of  $\epsilon$ -optimally solving zs-POSGs. In contrast to SFLP or CFR+, we provide the necessary foundational building blocks to apply dynamic programming (in tandem with heuristic search) to solve zs-POSGs. We introduce Bellman optimality equations and uniform-continuity properties of the optimal value function. Next, we exhibit rules for updating value functions while preserving uniform continuity and the ability to extract globally-consistent solutions. Finally, we describe the first effective DP algorithm for zs-POSGs, zs-oMG-HSVI, with finite-time convergence to an  $\epsilon$ -optimal solution. Experiments support our theoretical findings.

We believe our approach complements existing ones, e.g., SFLP and CFR+, in two dimensions. First, it breaks the original zs-POSG into subgames. Second, it generalizes values from visited subgames to unvisited ones. Our performances are as good as or better than those from SFLP and CFR+ for small-dimensional subgames (e.g., with TOI structure). Unfortunately, the advantage of breaking the original problem into subgames and exploiting uniform continuity properties often fails to fully manifest in the overall computational time.

Despite some similarities, our (offline) approach is fundamentally different from (online) continual resolving approaches. The latter could even possibly be adapted to use other offline methods than CFR-based ones, including HSVI.

We hope that this approach will lay the foundation for further work in the area of both exact and approximate DP solutions for zs-POSGs. In the short term, we shall investigate pruning techniques, better Lipschitz constants, and

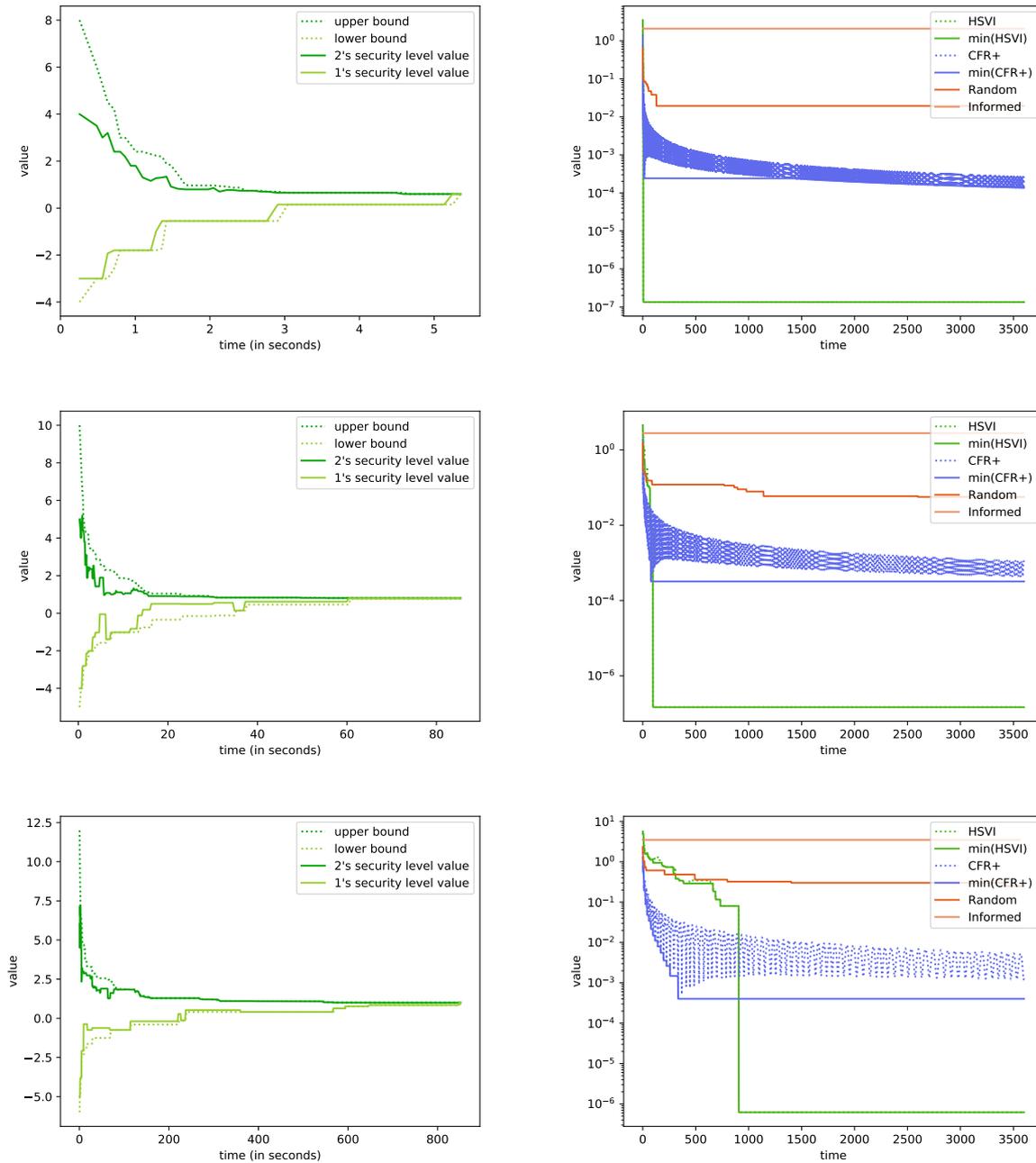


Figure 9: **Matching Pennies** ( $H = 4, 5, 6$ ) **(1,1,1)**: **(left)** Evolution of (in dotted lines) the upper- and lower-bound values, and (in solid lines) the security levels of the returned strategies for HSVI as a function of time (ms). **(right)** Exploitability ( $= \frac{SL\text{-gap}}{2}$ ) as a function of time (s) for Random, Informed, CFR+, and HSVI.

improved initial bounding approximators using solutions from relaxations of zs-POSGs, e.g., zs-OS-POSGs. In the long term, we shall investigate (deep) RL for zs-POSGs, similarly to a recent approach for Dec-POMDPs [Bono et al., 2018]. The latter shall investigate the trade-off between the update-rule accuracy and the computational efficiency when facing high-dimensional subgames, hence providing competitive solvers.

## References

Nicola Basilico, Giuseppe De Nittis, and Nicola Gatti. A security game combining patrolling and alarm-triggered responses under spatial and detection uncertainties. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

- Arnab Basu and Lukasz Stettner. Finite- and infinite-horizon Shapley games with nonsymmetric partial observation. *SIAM Journal on Control and Optimization*, 53(6):3584–3619, 2015.
- Guillaume Bono, Jilles Dibangoye, Laëtitia Maignon, Florian Pereyron, and Olivier Simonin. Cooperative multi-agent policy gradient. In *Proceedings of the Twenty-Eight European Conference on Machine Learning*, 2018.
- Branislav Bošanský, Christopher Kiekintveld, Viliam Lisý, and Michal Pěchouček. An exact double-oracle algorithm for zero-sum extensive-form games with imperfect information. *Journal of Artificial Intelligence Research*, 51: 829–866, 2014. doi: 10.1613/jair.4477.
- Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.
- Noam Brown, Christian Kroer, and Tuomas Sandholm. Dynamic thresholding and pruning for regret minimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Noam Brown, Anton Bakhtin, Adam Lerer, and Qucheng Gong. Combining deep reinforcement learning and search for imperfect-information games. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 17057–17069, 2020.
- Neil Burch, Michael Johanson, and Michael Bowling. Solving imperfect information games using decomposition. In *Twenty-eighth AAAI conference on artificial intelligence*, 2014.
- Neil Burch, Matej Moravcik, and Martin Schmid. Revisiting cfr+ and alternating updates. *Journal of Artificial Intelligence Research*, 64:429–443, 2019.
- Krishnendu Chatterjee and Laurent Doyen. Partial-observation stochastic games: How to win when belief fails. *ACM Transactions on Computational Logic*, 15(2):16, 2014.
- Harold L. Cole and Narayana Kocherlakota. Dynamic games with hidden actions and hidden states. *Journal of Economic Theory*, 98(1):114–126, 2001.
- Jilles Dibangoye, Chris Amato, Olivier Buffet, and François Charpillet. Optimally solving Dec-POMDPs as continuous-state MDPs. *Journal of Artificial Intelligence Research*, 55:443–497, 2016.
- Mrinal K. Ghosh, David R. McDonald, and Sagnik Sinha. Zero-sum stochastic games with partial information. *Journal of Optimization Theory and Applications*, 121(1):99–118, April 2004.
- John C. Harsanyi. Games with incomplete information played by "Bayesian" players, I-III. part II. Bayesian equilibrium points. *Management Science*, 14(5):320–334, January 1968. URL <http://www.jstor.org/stable/2628673>.
- Samid Hoda, Andrew Gilpin, Javier Peña, and Tuomas Sandholm. Smoothing techniques for computing nash equilibria of sequential games. *Mathematics of Operations Research*, 35(2):494–512, 2010. ISSN 0364765X, 15265471. URL <http://www.jstor.org/stable/40801239>.
- Karel Horák. *Scalable Algorithms for Solving Stochastic Games with Limited Partial Observability*. PhD thesis, Czech Technical University in Prague, Faculty of Electrical Engineering, 2019.
- Karel Horák and Branislav Bošanský. Solving partially observable stochastic games with public observations. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence*, pages 2029–2036, 2019.
- Karel Horák, Branislav Bošanský, and Michal Pěchouček. Heuristic search value iteration for one-sided partially observable stochastic games. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pages 558–564, 2017.
- Daphne Koller, Nimrod Megiddo, and Bernhard von Stengel. Fast algorithms for finding randomized strategies in game trees. In *Proceedings of the 26th ACM Symposium on the Theory of Computing (STOC'94)*, pages 750–759, 1994.
- Daphne Koller, Nimrod Megiddo, and Bernhard von Stengel. Efficient computation of equilibria for extensive two-person games. *Games and Economic Behavior*, 14(51):220–246, 1996.
- Vojtěch Kovařík, Dominik Seitz, Viliam Lisý, Jan Rudolf, Shuo Sun, and Karel Ha. Value functions for depth-limited solving in imperfect-information games. *arXiv preprint arXiv:1906.06412*, 2019.
- Vojtěch Kovařík, Martin Schmid, Neil Burch, Michael Bowling, and Viliam Lisý. Rethinking formal models of partially observable multiagent decision making. *CoRR*, abs/1906.11110, 2019.
- Christian Kroer, Kevin Waugh, Fatma Kilinç-Karzan, and Tuomas Sandholm. Faster algorithms for extensive-form game solving via improved smoothing functions. *Mathematical Programming*, 179:385–417, 2020. doi: 10.1007/s10107-018-1336-7.
- Harold W. Kuhn. Simplified two-person Poker. In H. W. Kuhn and A. W. Tucker, editors, *Contributions to the Theory of Games*, volume 1, 1950.
- Marc Lanctot, Kevin Waugh, Martin Zinkevich, and Michael Bowling. Monte carlo sampling for regret minimization in extensive games. *Advances in neural information processing systems*, 22, 2009.

- Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel Hennes, Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Faulkner, János Kramár, Bart De Vylder, Brennan Saeta, James Bradbury, David Ding, Sebastian Borgeaud, Matthew Lai, Julian Schrittwieser, Thomas Anthony, Edward Hughes, Ivo Danihelka, and Jonah Ryan-Davis. OpenSpiel: A framework for reinforcement learning in games. *CoRR*, abs/1908.09453, 2019. URL <http://arxiv.org/abs/1908.09453>.
- Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. DeepStack: Expert-level artificial intelligence in heads-up no-limit Poker. *Science*, 356(6337):508–513, 2017.
- Frans Oliehoek and Nikos Vlassis. Dec-POMDPs and extensive form games: equivalence of models and algorithms. Technical Report IAS-UVA-06-02, Intelligent Systems Laboratory Amsterdam, University of Amsterdam, 2006.
- Karl Åström. Optimal control of Markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications*, 10(1):174 – 205, 1965. ISSN 0022-247X.
- Martin Schmid. *Search in Imperfect Information Games*. PhD thesis, Charles University - Univerzita Karlova, Prague, 2021. URL <https://arxiv.org/pdf/2111.05884.pdf>.
- Martin Schmid, Matej Moravcik, Neil Burch, Rudolf Kadlec, Joshua Davidson, Kevin Waugh, Nolan Bard, Finbarr Timbers, Marc Lanctot, Zach Holland, Elnaz Davoodi, Alden Christianson, and Michael Bowling. Player of games. *CoRR*, abs/2112.03178, 2021. URL <https://arxiv.org/abs/2112.03178>.
- Trey Smith. *Probabilistic Planning for Robotic Exploration*. PhD thesis, The Robotics Institute, Carnegie Mellon University, 2007.
- Trey Smith and R.G. Simmons. Point-based POMDP algorithms: Improved analysis and implementation. In *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*, pages 542–549, 2005.
- Daniel Szer, François Charpillet, and Shlomo Zilberstein. MAA\*: A heuristic search algorithm for solving decentralized POMDPs. In *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*, pages 576–583, 2005.
- Oskari Tammelin. Solving large imperfect information games using cfr+. *arXiv preprint arXiv:1407.5042*, 2014.
- John von Neumann. Zur Theorie der Gesellschaftsspiele. *Mathematische Annalen*, 100, 1928.
- Bernhard von Stengel. Efficient computation of behavior strategies. *Games and Economic Behavior*, 14(50):220–246, 1996.
- Auke Wiggers. Structure in the value function of two-player zero-sum games of incomplete information. Master’s thesis, University of Amsterdam, 2015.
- Auke Wiggers, Frans Oliehoek, and Diederik Roijers. Structure in the value function of two-player zero-sum games of incomplete information. *Computing Research Repository*, abs/1606.06888, 2016a.
- Auke Wiggers, Frans Oliehoek, and Diederik Roijers. Structure in the value function of two-player zero-sum games of incomplete information. In *Proceedings of the Twenty-Second European Conference on Artificial Intelligence*, pages 1628–1629, 2016b. doi: 10.3233/978-1-61499-672-9-1628.
- Martin Zinkevich, Michael Johanson, Michael Bowling, and Carmelo Piccione. Regret minimization in games with incomplete information. In *Advances in Neural Information Processing Systems 20*, 2007.

## A Synthetic Tables

For convenience, we provide two synthetic tables: Table 3 to sum up various theoretical properties that are stated in this paper (assuming a finite temporal horizon), and Table 4 to sum up the notations used in this paper, adding some notations that appear only in the appendix.

More precisely, Table 3 indicates, for various functions  $f$  and variables  $x$ , properties that  $f$  is known to exhibit with respect to  $x$ . We denote by

- a function with no known (or used) property (see also comment below);
- N/A a non-applicable case;
- Lin* a linear function;
- LC* a Lipschitz-continuous function;
- Cv* (resp. *Cc*) a convex (resp. concave) function;
- PWLCv* (resp. *PWLCc*) a piecewise linear and convex (resp. concave) function;
- $\perp\!\!\!\perp$  the function being independent of the variable;
- $\neg P$  the negation of some property  $P$  (i.e.,  $P$  is known not to hold).

Note also that, as  $\sigma_\tau = \sigma_\tau^{c,1} \sigma_\tau^{m,1}$ , the linearity or Lipschitz-continuity properties of any function w.r.t.  $\sigma_\tau$  extends to both  $\sigma_\tau^{c,1}$  and  $\sigma_\tau^{m,1}$ . Reciprocally, related negative results extend from  $\sigma_\tau^{c,1}$  or  $\sigma_\tau^{m,1}$  to  $\sigma_\tau$ . In these three columns, we just indicate results that cannot be derived from one of the two other columns.

Table 3: Known properties of various functions appearing in this work

	$\sigma_\tau$	$\sigma_\tau^{m,1}$	$\sigma_\tau^{c,1}$	$\beta_\tau^i$	$\beta_\tau^{-i}$
$T(\sigma_\tau, \beta_\tau)$	<i>Lin</i> (prop. 2.3, p. 5)	-	-	<i>Lin</i> (prop. 2.3, p. 5)	<i>Lin</i> (prop. 2.3, p. 5)
$T_m^i(\sigma_\tau, \beta_\tau)$	<i>Lin</i> (lem. 3, p. 28)	-	-	<i>Lin</i> (lem. 3, p. 28)	<i>Lin</i> (lem. 3, p. 28)
$T_c^i(\sigma_\tau, \beta_\tau)$	-	$\perp\!\!\!\perp$ (lem. 4, p. 29)	-	$\perp\!\!\!\perp$ (lem. 4, p. 29)	-
$V_\tau^*(\sigma_\tau)$	<i>LC</i> (app. D.1.3, p. 30)	<i>PWLCv</i> (thm. 2.5, p. 6)	-	N/A	N/A
$W_\tau^{i,*}(\sigma_\tau, \beta_\tau^i)$	<i>LC</i> (from $Q_{\tau+1}^*$ LC)	-	-	$\neg$ <i>Lin</i> (from $Q^* \neg$ <i>Lin</i> )	N/A
$\nu_{[\sigma_\tau^{c,1}, \beta_\tau^2]}^2$	N/A	N/A	<i>LC</i> (lem. 7, p. 31)	N/A	-

## B Background

### B.1 Occupancy States

The following result shows that the occupancy state is (i) Markovian, i.e., its value at  $\tau$  only depends on its previous value  $\sigma_{\tau-1}$ , the system dynamics  $P_{a^1, a^2}^{z^1, z^2}(s'|s)$ , and the last behavioral decision rules  $\beta_{\tau-1}^1$  and  $\beta_{\tau-1}^2$ , and (ii) sufficient to estimate the expected reward. Note that it holds for general-sum POSGs with any number of agents, and as many reward functions; similar results have already been established, e.g., for Dec-POMDPs (cf. [Dibangoye et al., 2016, Theorem 1]).

**Proposition 2.3.** (originally stated on page 5)  $\sigma_{\beta_{0:\tau-1}}$ , together with  $\beta_\tau$ , is a sufficient statistic to compute (i) the next OS,  $T(\sigma_{\beta_{0:\tau-1}}, \beta_\tau) \stackrel{\text{def}}{=} \sigma_{\beta_{0:\tau}}$ , and (ii) the expected reward at  $\tau$ :  $r(\sigma_{\beta_{0:\tau-1}}, \beta_\tau) \stackrel{\text{def}}{=} \mathbb{E}[R_\tau | \beta_{0:\tau-1} \oplus \beta_\tau]$ , where  $\oplus$  denotes a concatenation.

*Proof.* Let us first derive a recursive way of computing  $\sigma_{\beta_{0:\tau}}(\theta_\tau, \mathbf{a}_\tau, \mathbf{z}_{\tau+1})$ :

$$\begin{aligned}
 \sigma_{\beta_{0:\tau}}(\theta_\tau, \mathbf{a}_\tau, \mathbf{z}_{\tau+1}) &\stackrel{\text{def}}{=} Pr(\theta_\tau, \mathbf{a}_\tau, \mathbf{z}_{\tau+1} | \beta_{0:\tau}) \\
 &= \sum_{s_\tau, s_{\tau+1}} Pr(\theta_\tau, \mathbf{a}_\tau, \mathbf{z}_{\tau+1}, s_\tau, s_{\tau+1} | \beta_{0:\tau}) \\
 &= \sum_{s_\tau, s_{\tau+1}} Pr(\mathbf{z}_{\tau+1}, s_{\tau+1} | \theta_\tau, \mathbf{a}_\tau, s_\tau, \beta_{0:\tau}) Pr(\mathbf{a}_\tau | \theta_\tau, s_\tau, \beta_{0:\tau}) \\
 &\quad Pr(s_\tau | \theta_\tau, \beta_{0:\tau}) Pr(\theta_\tau | \beta_{0:\tau})
 \end{aligned}$$

Table 4: Various notations used in this work

---

$\neg i \stackrel{\text{def}}{=} i$ 's opponent. Thus:  $\neg 1 = 2$ , and  $\neg 2 = 1$ .

Histories and occupancy states

$\theta_\tau^i \stackrel{\text{def}}{=} (a_0^i, z_1^i, \dots, a_{\tau-1}^i, z_\tau^i) (\in \Theta^i = \cup_{t=0}^{H-1} \Theta_t^i)$  is a length- $\tau$  *action-observation history* (AOH) for  $i$ .  
 $\theta_\tau \stackrel{\text{def}}{=} (\theta_\tau^1, \theta_\tau^2) (\in \Theta = \cup_{t=0}^{H-1} \Theta_t)$  is a *joint* AOH at  $\tau$ .  
 $\sigma_\tau(\theta_\tau) \stackrel{\text{def}}{=} \text{Occupancy state}$  (OS)  $\sigma_\tau (\in \mathcal{O}^\sigma = \cup_{t=0}^{H-1} \mathcal{O}_t^\sigma$ , where  $\mathcal{O}_t^\sigma \stackrel{\text{def}}{=} \Delta(\Theta_\tau)$ ), *i.e.*, probability distribution over joint AOHs  $\theta_\tau$  (typically for some applied  $\beta_{0:\tau-1}$ ).  
 $\sigma_\tau^{m,i}(\theta_\tau^i) \stackrel{\text{def}}{=} \text{Marginal term}$  of  $\sigma_\tau$  from player  $i$ 's point of view ( $\sigma_\tau^{m,i} \in \Delta(\Theta_\tau^i)$ ).  
 $\sigma_\tau^{c,i}(\theta_\tau^i | \theta_\tau^{\neg i}) \stackrel{\text{def}}{=} \text{Conditional term}$  of  $\sigma_\tau$  from  $i$ 's point of view ( $\sigma_\tau^{c,i} : \Theta_\tau^i \mapsto \Delta(\Theta_\tau^{\neg i})$ ).  
 $b(s|\theta_\tau) \stackrel{\text{def}}{=} \text{Belief state}$ , *i.e.*, probability distribution over states given a joint AOH ( $b(s|\theta_\tau) : \mathcal{S} \times \Theta_\tau \mapsto \mathbb{R}$ ). Can be computed by an HMM filtering process.  
 $o_\tau \stackrel{\text{def}}{=} \text{Full occupancy state}$   $o_\tau (\in \Delta(\mathcal{S} \times \Theta_\tau))$ , *i.e.*,  $Pr(s, \theta_\tau)$  for the current  $\beta_{0:\tau-1}$ , and thus verifies  $\sigma_\tau(\theta_\tau) = \sum_{s \in \mathcal{S}} o_\tau(s, \theta_\tau)$ . Is used in the implementation to simplify computations (*e.g.*, of  $r_t$  and  $\sigma_{\tau+1}$  through  $b$ ).

Decision rules and strategies

$\pi_{0:\tau}^i \stackrel{\text{def}}{=} \text{A pure strategy}$  for  $i$  is a mapping  $\pi_{0:\tau}^i$  from private histories in  $\Theta_t^i (\forall t \in \{0 \dots \tau\})$  to *single* private actions in  $\mathcal{A}^i$ . By default,  $\pi^i \stackrel{\text{def}}{=} \pi_{0:H-1}^i$ .  
 $\mu_{0:\tau}^i \stackrel{\text{def}}{=} \text{A mixed strategy}$   $\mu_{0:\tau}^i$  for  $i$  is a probability distribution over pure strategies. It is used by first sampling one of the pure strategies (at  $t = 0$ ), and then executing it until  $t = \tau$ .  
 $\mu_{0:\tau'}^i | \sigma_\tau \stackrel{\text{def}}{=} (\tau \leq \tau')$  is a *mixed strategy compatible* with some OS  $\sigma_\tau$ , *i.e.*, that could induce this OS at  $\tau$  (assuming an appropriate complementary  $\mu_{0:\tau'}^{\neg i} | \sigma_\tau$ ).  
 $\beta_\tau^i \stackrel{\text{def}}{=} \text{A (behavioral) decision rule}$  (DR) at time  $\tau$  for  $i$  is a mapping  $\beta_\tau^i$  from private AOHs in  $\Theta_\tau^i$  to *distributions* over private actions. We note  $\beta_\tau^i(\theta_\tau^i, a^i)$  the probability to pick  $a^i$  when facing  $\theta_\tau^i$ .  
 $\beta_{\tau:\tau'}^i \stackrel{\text{def}}{=} (\beta_\tau^i, \dots, \beta_{\tau'}^i)$  is a *behavioral strategy* for  $i$  from time step  $\tau$  to  $\tau'$  (included).  
 $rw^i(\theta_\tau^i, a_\tau^i) \stackrel{\text{def}}{=} \prod_{t=0}^{\tau} \beta_{0:t}^i(a_t^i | a_0^i, z_1^i, a_1^i, \dots, z_t^i)$  is the *realization weight* (RW) of sequence  $a_0^i, z_1^i, a_1^i, \dots, a_\tau^i (= \theta_\tau^i, a_\tau^i)$  under strategy  $\beta_{0:\tau}^i$ .  
 $rw^i(\phi_{\tau:\tau'}^i | \theta_\tau^i) \stackrel{\text{def}}{=} \prod_{t=\tau}^{\tau'} \beta_{0:t}^i(a_t^i | \theta_\tau^i, a_\tau^i, \dots, z_t^i)$  is the RW of a *suffix sequence*  $\phi_{\tau:\tau'}^i = a_\tau^i, \dots, a_{\tau'}^i$ , “conditioned” on a *prefix sequence/AOH*  $\theta_\tau^i$ .  
 $\pi_{0:\tau} \stackrel{\text{def}}{=} \text{is a pure strategy profile}$ .  
 $\mu_{0:\tau} \stackrel{\text{def}}{=} \text{is a mixed strategy profile}$ .  
 $\mu_{0:\tau'} | \sigma_\tau \stackrel{\text{def}}{=} (\tau \leq \tau')$  is a *mixed strategy profile compatible* with some OS  $\sigma_\tau$ , *i.e.*, that could induce this OS at  $\tau$ .  
 $\beta_\tau \stackrel{\text{def}}{=} \langle \beta_\tau^1, \beta_\tau^2 \rangle (\in \mathcal{B} = \cup_{t=0}^{H-1} \mathcal{B}_t)$  is a *decision rule profile*.  
 $\beta_{\tau:\tau'} \stackrel{\text{def}}{=} \langle \beta_{\tau:\tau'}^1, \beta_{\tau:\tau'}^2 \rangle$  is a *behavioral strategy profile*.

$$= \sum_{s_\tau, s_{\tau+1}} \underbrace{Pr(z_{\tau+1}, s_{\tau+1} | \mathbf{a}_\tau, s_\tau)}_{=P_{\mathbf{a}_\tau}^{z_{\tau+1}}(s_{\tau+1}|s_\tau)} \underbrace{Pr(\mathbf{a}_\tau | \theta_\tau, \beta_\tau)}_{=\beta(\theta_\tau, \mathbf{a}_\tau)} \underbrace{Pr(s_\tau | \theta_\tau, \beta_{0:\tau})}_{=b(s_\tau | \theta_\tau)} \underbrace{Pr(\theta_\tau | \beta_{0:\tau-1})}_{=\sigma_{\beta_{0:\tau-1}}(\theta_\tau)}$$

(where  $b(s | \theta_\tau)$  is the belief over states obtained by a usual HMM filtering process)

$$= \sum_{s_\tau, s_{\tau+1}} P_{\mathbf{a}_\tau}^{z_{\tau+1}}(s_{\tau+1}|s_\tau) \beta(\theta_\tau, \mathbf{a}_\tau) b(s_\tau | \theta_\tau) \sigma_{\beta_{0:\tau-1}}(\theta_\tau).$$

$\sigma_{\beta_{0:\tau}}$  can thus be computed from  $\sigma_{\beta_{0:\tau-1}}$  and  $\beta_\tau$  without explicitly using  $\beta_{0:\tau-1}$  or earlier occupancy states.

Then, let us compute the expected reward at  $\tau$  given  $\beta_{0:\tau}$ :

$$\begin{aligned} & E[r(S_\tau, A_\tau^1, A_\tau^2) | \beta_{0:\tau}] \\ &= \sum_{s_\tau, \mathbf{a}_\tau} r(s_\tau, \mathbf{a}_\tau) Pr(s_\tau, \mathbf{a}_\tau | \beta_{0:\tau}) \\ &= \sum_{s_\tau, \mathbf{a}_\tau} \sum_{\theta_\tau} r(s_\tau, \mathbf{a}_\tau) Pr(s_\tau, \mathbf{a}_\tau, \theta_\tau | \beta_{0:\tau}) \end{aligned}$$

Rewards and value functions

$$\begin{aligned}
 r_{\max} &\stackrel{\text{def}}{=} \max_{s, \mathbf{a}} r(s, \mathbf{a}) && \text{Maximum possible reward.} \\
 r_{\min} &\stackrel{\text{def}}{=} \min_{s, \mathbf{a}} r(s, \mathbf{a}) && \text{Minimum possible reward.} \\
 V_{\tau}(\sigma_{\tau}, \beta_{\tau}) &\stackrel{\text{def}}{=} E[\sum_{t=\tau}^{H-1} \gamma^t R_t \mid \sigma_{\tau}, \beta_{\tau}], && \text{Value of } \beta_{\tau:H-1} \text{ in OS } \sigma_{\tau}. \\
 &&& \text{where } R_t \text{ is the random var. for the reward at } t. \\
 V_{\tau}^*(\sigma_{\tau}) &\stackrel{\text{def}}{=} \max_{\beta_{\tau}^1} \min_{\beta_{\tau}^2} V_{\tau}(\sigma_{\tau}, \beta_{\tau}) && \text{Optimal value function} \\
 Q_{\tau}^*(\sigma_{\tau}, \beta_{\tau}) &\stackrel{\text{def}}{=} r(\sigma_{\tau}, \beta_{\tau}) + \gamma V_{\tau+1}^*(T(\sigma_{\tau}, \beta_{\tau})) && \text{Opt. (joint) action-value fct.} \\
 W_{\tau}^{i,*}(\sigma_{\tau}, \beta_{\tau}) &\stackrel{\text{def}}{=} \text{opt}_{\beta_{\tau}^i} Q_{\tau}^*(\sigma_{\tau}, \beta_{\tau}), && \text{Opt. (individual) action-value fct.} \\
 &&& \text{where } \text{opt} = \max \text{ if } i = 1, \text{ min otherwise.} \\
 \nu_{[\sigma_{\tau}^{c,1}, \beta_{\tau}^2]} &\stackrel{\text{def}}{=} \text{Vector of values (one component per AOH } \theta_{\tau}^1) \text{ for } 1\text{'s best response to } \beta_{\tau}^2. \text{ Assuming } \sigma_{\tau}^{c,1}. \text{ This} \\
 &&& \text{solution of a POMDP allows computing } V_{\tau}^* \text{ (see Theorem 2.5).}
 \end{aligned}$$

Approximations

$$\begin{aligned}
 \bar{V}_{\tau}(\sigma_{\tau}) &\stackrel{\text{def}}{=} \text{Upper bound approximation of } V_{\tau}^*(\sigma_{\tau}); \text{ relies on data set } \bar{\mathcal{I}}_{\tau-1}. \\
 \underline{V}_{\tau}(\sigma_{\tau}) &\stackrel{\text{def}}{=} \text{Lower bound approximation of } V_{\tau}^*(\sigma_{\tau}); \text{ relies on data set } \underline{\mathcal{I}}_{\tau-1}. \\
 \bar{W}_{\tau}(\sigma_{\tau}, \beta_{\tau}^1) &\stackrel{\text{def}}{=} \text{Upper bound approximation of } W_{\tau}^{*,1}(\sigma_{\tau}, \beta_{\tau}^1); \text{ relies on data set } \bar{\mathcal{I}}_{\tau}. \\
 \underline{W}_{\tau}(\sigma_{\tau}, \beta_{\tau}^2) &\stackrel{\text{def}}{=} \text{Lower bound approximation of } W_{\tau}^{*,2}(\sigma_{\tau}, \beta_{\tau}^2); \text{ relies on data set } \underline{\mathcal{I}}_{\tau}. \\
 \bar{\nu}_{\tau}^2 &\stackrel{\text{def}}{=} \text{Vector (with one component per AOH } \theta_{\tau}^1) \text{ used in } \bar{V}_{\tau} \text{ and } \bar{W}_{\tau-1} \text{ (if } \tau \geq 1).
 \end{aligned}$$

Miscellaneous

$$\begin{aligned}
 w_{\tau} &\stackrel{\text{def}}{=} \text{Denotes a triplet } \langle \sigma_{\tau-1}^{c,1}, \beta_{\tau-1}^1, \bar{\nu}_{\tau}^2 \rangle \in \bar{\mathcal{I}}_{\tau} \text{ (or a triplet in } \underline{\mathcal{I}}_{\tau}). \\
 \psi_{\tau}^2 &\stackrel{\text{def}}{=} \text{Distribution over triplets } w_{\tau+1} \in \bar{\mathcal{I}}_{\tau+1} \text{ (inducing a recursively defined strategy from } \tau \text{ to } H-1). \\
 &&& \text{Often denotes the strategy it induces.} \\
 x^{\top} &\stackrel{\text{def}}{=} \text{The transpose of a (usually column) vector } x \text{ of } \mathbb{R}^n. \\
 c[y] &\stackrel{\text{def}}{=} \text{Denotes field } c \text{ of object/tuple } y. \\
 \text{Supp}(d) &\stackrel{\text{def}}{=} \text{Support of distribution } d, \text{ i.e., set of its non-zero probability elements.}
 \end{aligned}$$

$$\begin{aligned}
 &= \sum_{s_{\tau}, \mathbf{a}_{\tau}} \sum_{\theta_{\tau}} r(s_{\tau}, \mathbf{a}_{\tau}) Pr(s_{\tau}, \mathbf{a}_{\tau} \mid \theta_{\tau}, \beta_{0:\tau}) Pr(\theta_{\tau} \mid \beta_{0:\tau}) \\
 &= \sum_{s_{\tau}, \mathbf{a}_{\tau}} \sum_{\theta_{\tau}} r(s_{\tau}, \mathbf{a}_{\tau}) \underbrace{Pr(\mathbf{a}_{\tau} \mid \theta_{\tau}, \beta_{0:\tau})}_{\beta_{\tau}(\theta_{\tau}, \mathbf{a}_{\tau})} \underbrace{Pr(s_{\tau} \mid \theta_{\tau}, \beta_{0:\tau})}_{b(s_{\tau} \mid \theta_{\tau})} \underbrace{Pr(\theta_{\tau} \mid \beta_{0:\tau})}_{\sigma_{\beta_{0:\tau-1}}(\theta_{\tau})} \\
 &= \sum_{s_{\tau}, \mathbf{a}_{\tau}} \sum_{\theta_{\tau}} r(s_{\tau}, \mathbf{a}_{\tau}) \beta_{\tau}(\theta_{\tau}, \mathbf{a}_{\tau}) b(s_{\tau} \mid \theta_{\tau}) \sigma_{\beta_{0:\tau-1}}(\theta_{\tau}).
 \end{aligned}$$

The expected reward at  $\tau$  can thus be computed from  $\sigma_{\beta_{0:\tau-1}}$  and  $\beta_{\tau}$  without explicitly using  $\beta_{0:\tau-1}$  or earlier occupancy states.  $\square$

## C Occupancy Markov Games: Definition and Preliminary Properties

### C.1 Properties of $V^*$

Before proving the postulates implicitly used by [Wiggers et al., 2016a], we need to show that we can reason with mixed strategies in subgames as is usually done on full games.

#### C.1.1 $V^*$ is not linear in $\beta$ (behavioral strategies)

Let us consider the following (finite-horizon, deterministic) Non-Observable MDP:

$$\begin{aligned}
 \mathcal{S} &\stackrel{\text{def}}{=} \{-2, -1, 0, +1, +2\}, && b_0(0) = 1, && \text{(always start in } s = 0) \\
 \mathcal{A} &\stackrel{\text{def}}{=} \{-1, +1\}, && && \text{(moves = add or subtract 1)} \\
 T(s, a) &\stackrel{\text{def}}{=} \min\{+2, \max\{-2, s + a\}\}, && && \text{(dép. de } +1 \text{ ou } -1 \text{ dans } \mathcal{S}) \\
 \mathcal{Z} &\stackrel{\text{def}}{=} \{none\}, && O(none) \stackrel{\text{def}}{=} 1, && \text{(no observation)}
 \end{aligned}$$

$$r(s) \stackrel{\text{def}}{=} \begin{cases} +1 & \text{si } s \in \{-2, +2\} \\ 0 & \text{sinon,} \end{cases} \quad (|s| = 2 : \text{victoire !})$$

$$\gamma \stackrel{\text{def}}{=} 1, \quad H \stackrel{\text{def}}{=} 2.$$

Let us then consider two particular behavioral strategies:

$$\begin{aligned} \forall \theta, \beta^+(A = +1|\theta) &= 1 && \text{(always +1), and} \\ \forall \theta, \beta^-(A = -1|\theta) &= 1 && \text{(always -1).} \end{aligned}$$

These two strategies are optimal, with an expected return of +1, because, at  $t = H = 2$ ,  $\beta^+$  reaches +2 w.p. 1, and  $\beta^-$  reaches -2 w.p. 1:

$$V(\beta^+) = V(\beta^-) = +1.$$

Let us now consider their linear combination  $\beta^\pm \stackrel{\text{def}}{=} \frac{1}{2}\beta^+ + \frac{1}{2}\beta^-$ :

$$\begin{aligned} \forall \theta, \beta^\pm(A = -1|\theta) &= 0.5, \\ \forall \theta, \beta^\pm(A = +1|\theta) &= 0.5. \end{aligned}$$

Here, the probability to reach  $s = -2$  or  $s = +2$  at the last time step is much lower, and gives the value of that strategy:

$$\begin{aligned} V(\beta^\pm) &= Pr(s_2 = +2|\beta^\pm) + Pr(s_2 = -2|\beta^\pm) \\ &= Pr(a_0 = +1|\beta^\pm) \cdot Pr(a_1 = +1|\beta^\pm) \\ &\quad + Pr(a_0 = -1|\beta^\pm) \cdot Pr(a_1 = -1|\beta^\pm) \\ &= \underbrace{0.5 \cdot 0.5}_{0.25} + \underbrace{0.5 \cdot 0.5}_{0.25} = 0.5. \end{aligned}$$

### C.1.2 Back to Mixed Strategies

We now generalize mixed strategies as a mathematical tool to handle subgames of a zs-OMG as normal-form games, and give some preliminary results.

First, for a given  $\sigma_\tau$  and  $\tau \leq \tau'$ , let  $\mu_{0:\tau'-1|\sigma_\tau}$  denote a mixed strategy profile that is defined over  $0 : \tau' - 1$ , and induces (is compatible with)  $\sigma_\tau$  at time  $\tau$ . Then, to complete a given mixed *prefix* strategy  $\mu_{0:\tau'-1|\sigma_\tau}$  (here  $\tau = \tau'$ ), the solver should provide each player with a different *suffix* strategy to execute for each  $\theta_\tau^i$  it could be facing. We now detail how to build an equivalent set of mixed *full* strategies for  $i$ . Each of the pure *prefix* strategies  $\pi_{0:\tau-1}^i$  used in  $\mu_{0:\tau-1|\sigma_\tau}^i$  (belonging to a set denoted  $\Pi_{0:\tau-1|\sigma_\tau}^i$ ) can be extended by appending a different pure *suffix* strategy  $\pi_{\tau:H-1}^i$  at each of its leaf nodes, which leads to a large set of pure strategies  $\Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)$ . Then, let  $M_{0:H-1|\sigma_\tau}^i$  be the set of mixed *full* strategies  $\mu_{0:H-1|\sigma_\tau}^i$  obtained by considering the distributions over  $\bigcup_{\pi_{0:\tau-1}^i \in \Pi_{0:\tau-1|\sigma_\tau}^i} \Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)$  that verify,  $\forall \pi_{0:\tau-1}^i$ ,

$$\sum_{\substack{\pi_{0:H-1}^i \in \\ \Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)}} \mu_{0:H-1|\sigma_\tau}^i(\pi_{0:H-1}^i) = \mu_{0:\tau-1|\sigma_\tau}^i(\pi_{0:\tau-1}^i). \quad (8)$$

This is the set of mixed strategies compatible with  $\sigma_\tau$ .

**Lemma 2.**  $M_{0:H-1|\sigma_\tau}^i$  is convex and equivalent to the set of behavioral strategies  $\beta_{0:H-1|\sigma_\tau}^i$ , thus sufficient to search for a Nash equilibrium in  $\sigma_\tau$ .

*Proof.* Let  $\mu_{0:H-1|\sigma_\tau}^i$  and  $\nu_{0:H-1|\sigma_\tau}^i$  be two mixed strategies in  $M_{0:H-1|\sigma_\tau}^i$ , i.e., which are both full and compatible with occupancy state  $\sigma_\tau$  at time step  $\tau$ , and  $\alpha \in [0, 1]$ . Then, for any  $\pi_{0:\tau-1}^i$ ,

$$\begin{aligned} &\sum_{\pi_{0:H-1}^i \in \Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)} \left[ \alpha \cdot \mu_{0:H-1|\sigma_\tau}^i(\pi_{0:H-1}^i) + (1 - \alpha) \cdot \nu_{0:H-1|\sigma_\tau}^i(\pi_{0:H-1}^i) \right] \\ &= \alpha \left[ \sum_{\pi_{0:H-1}^i \in \Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)} \mu_{0:H-1|\sigma_\tau}^i(\pi_{0:H-1}^i) \right] \end{aligned}$$

$$+ (1 - \alpha) \left[ \sum_{\pi_{0:H-1}^i \in \Pi_{0:H-1}^i(\pi_{0:\tau-1}^i)} \nu_{0:H-1|\sigma_\tau}^i(\pi_{0:H-1}^i) \right]$$

(because both mixed strategies are compatible with  $\sigma_\tau$  (eq. 8, p. 26):)

$$\begin{aligned} &= \alpha \cdot \mu_{0:\tau-1|\sigma_\tau}^i(\pi_{0:\tau-1}^i) + (1 - \alpha) \cdot \nu_{0:\tau-1|\sigma_\tau}^i(\pi_{0:\tau-1}^i) \\ &= \mu_{0:\tau-1|\sigma_\tau}^i(\pi_{0:\tau-1}^i). \end{aligned}$$

Eq. 8 thus also applies to  $\alpha \cdot \mu_{0:H-1|\sigma_\tau}^i + (1 - \alpha) \cdot \nu_{0:H-1|\sigma_\tau}^i$ , proving that it belongs to  $M_{0:H-1|\sigma_\tau}^i$  and, as a consequence, that this set is convex.

The equivalence with the set of behavioral strategies simply relies on the fact that all mixed strategies over  $\tau : H - 1$  can be independently generated at each action-observation history  $\theta_{0:\tau-1}^i$ .  $\square$

While only future rewards are relevant when making a decision at  $\tau$ , reasoning with mixed strategies defined from  $t = 0$  will be convenient because  $V_\tau(\sigma_\tau, \cdot, \cdot)$  is linear in  $\mu_{0:H-1|\sigma_\tau}^i$ , which allows coming back to a standard normal-form game and applying known results.

In the remaining, we simply note  $\mu^i$  (without index) the mixed strategies in  $M_{0:H-1|\sigma_\tau}^i$ , set which we now note  $M_{|\sigma_\tau}^i$ . Also, since we shall work with local game  $Q_\tau^*(\sigma_\tau, \beta_\tau)$ , let us define:  $M_{|\sigma_\tau, \beta_\tau}^i$  the set of  $i$ 's mixed strategies compatible with occupancy states reachable given  $\sigma_\tau$  and  $\beta_\tau^j$  (with either  $j = i$  or  $j = -i$ ). Then,  $M_{|T(\sigma_\tau, \beta_\tau)}^i \subseteq M_{|\sigma_\tau, \beta_\tau}^i \subseteq M_{|\sigma_\tau}^i$  (inclusion due to the latter sets being less constrained in their definition). As a consequence, if maximizing some function  $f$  over  $i$ 's mixed strategies compatible with a given  $\sigma_\tau$ :

$$\max_{\mu^i \in M_{|\sigma_\tau}^i} f(\sigma_\tau, \mu^i, \dots) \geq \max_{\mu^i \in M_{|\sigma_\tau, \beta_\tau}^i} f(\sigma_\tau, \mu^i, \dots) \geq \max_{\mu^i \in M_{|T(\sigma_\tau, \beta_\tau)}^i} f(\sigma_\tau, \mu^i, \dots).$$

### C.1.3 Von Neumann's Minimax Theorem for Subgames and a Bellman Optimality Equation

Using the previous results, one can show that von Neumann's minimax theorem applies in any subgame, allowing to swap operators max and min.

**Theorem 3.1.** (originally stated on page 7) *The subgame defined in Eq. (3) admits a unique NEV*

$$V_\tau^*(\sigma_\tau) \stackrel{\text{def}}{=} \max_{\beta_{\tau:H-1}^1} \min_{\beta_{\tau:H-1}^2} V_\tau(\sigma_\tau, \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2). \quad (4)$$

*Proof.* For any occupancy state  $\sigma_\tau$ ,

$$\max_{\beta_\tau^1} \min_{\beta_\tau^2} V(\sigma_\tau, \beta_\tau) = \max_{\mu_\tau^1} \min_{\mu_\tau^2} V(\sigma_\tau, \mu_\tau) \quad (\text{Kuhn's theorem (generalized)}) \quad (9)$$

$$= \min_{\mu_\tau^2} \max_{\mu_\tau^1} V(\sigma_\tau, \mu_\tau) \quad (\text{von Neumann's theorem}) \quad (10)$$

$$= \min_{\beta_\tau^2} \max_{\beta_\tau^1} V(\sigma_\tau, \beta_\tau) \quad (\text{again Kuhn's theorem (generalized)}) \quad (11)$$

$\square$

One can also show that a Bellman optimality equation allows relating optimal values in subgames at  $\tau$  and  $\tau + 1$ , leading to a recursive expression of  $V_0^*$ .

**Theorem 3.2.** (originally stated on page 7)  *$V_\tau^*(\sigma_\tau)$  satisfies the following functional equation:*

$$V_\tau^*(\sigma_\tau) = \max_{\beta_\tau^1} \min_{\beta_\tau^2} r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau)) = \max_{\beta_\tau^1} \min_{\beta_\tau^2} Q_\tau^*(\sigma_\tau, \beta_\tau).$$

*Proof.* Focusing, without loss of generality, on player 1, we have (complementary explanations follow for numbered lines in particular):

$$\max_{\beta_\tau^1} \min_{\beta_\tau^2} Q_\tau^*(\sigma_\tau, \beta_\tau) = \max_{\beta_\tau^1} \min_{\beta_\tau^2} [r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau))]$$

( $V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau))$  being the Nash equilibrium value of normal-form game  $V_{\tau+1}(T(\sigma_\tau, \beta_\tau), \mu^1, \mu^2)$ ):

$$= \max_{\beta_\tau^1} \min_{\beta_\tau^2} \left[ r(\sigma_\tau, \beta_\tau) + \gamma \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau}^1} \min_{\mu^2 \in M_{|\sigma_\tau, \beta_\tau}^2} V_{\tau+1}(T(\sigma_\tau, \beta_\tau), \mu^1, \mu^2) \right]$$

$$= \max_{\beta_\tau^1} \min_{\beta_\tau^2} \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau}^1} \min_{\mu^2 \in M_{|\sigma_\tau, \beta_\tau}^2} [r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2), \mu^1, \mu^2)]$$

(using the equivalence between maximin and minimax values for the (constrained normal-form) game at  $\tau + 1$ , the last two max and min operators can be swapped:)

$$= \max_{\beta_\tau^1} \min_{\beta_\tau^2} \min_{\mu^2 \in M_{|\sigma_\tau, \beta_\tau}^2} \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau}^1} [r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2), \mu^1, \mu^2)]$$

(merging both mins (and with explanations thereafter):)

$$= \max_{\beta_\tau^1} \min_{\mu^2 \in M_{|\sigma_\tau, \beta_\tau^1}^2} \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)}^1} [r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)), \mu^1, \mu^2)] \quad (12)$$

(since ignoring the opponent's decision rule does not influence the expected return:)

$$= \max_{\beta_\tau^1} \min_{\mu^2 \in M_{|\sigma_\tau}^2} \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau^1}^1} [r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)), \mu^1, \mu^2)]$$

(using again the minimax theorem's equivalence between maximin and minimax on an appropriate game:)

$$= \max_{\beta_\tau^1} \max_{\mu^1 \in M_{|\sigma_\tau, \beta_\tau^1}^1} \min_{\mu^2 \in M_{|\sigma_\tau}^2} [r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2(\mu^2)), \mu^1, \mu^2)] \quad (13)$$

(merging both maxs (and with explanations thereafter):)

$$= \max_{\mu^1 \in M_{|\sigma_\tau}^1} \min_{\mu^2 \in M_{|\sigma_\tau}^2} [r(\sigma_\tau, \beta_\tau^1(\mu^1), \beta_\tau^2(\mu^2)) + \gamma V_{\tau+1}(T(\sigma_\tau, \beta_\tau^1(\mu^1), \beta_\tau^2(\mu^2)), \mu^1, \mu^2)] \quad (14)$$

(again with the equivalence property discussed before the lemma:)

$$\begin{aligned} &= \max_{\mu^1 \in M_{|\sigma_\tau}^1} \min_{\mu^2 \in M_{|\sigma_\tau}^2} V_\tau(\sigma_\tau, \mu^1, \mu^2) \\ &= \max_{\beta_{\tau:H-1}^1} \min_{\beta_{\tau:H-1}^2} V_\tau(\sigma_\tau, \beta_{\tau:H-1}^1, \beta_{\tau:H-1}^2) \\ &\stackrel{\text{def}}{=} V_\tau^*(\sigma_\tau). \end{aligned}$$

More precisely, line 12 (and, similarly, line 14) is obtained by observing that

- minimizing over both (i)  $\beta_\tau^2$  and (ii)  $\mu^2$  constrained by  $\sigma_\tau$  and  $\beta_\tau$  is equivalent to minimizing over  $\mu^2$  constrained by  $\sigma_\tau$  and  $\beta_\tau^1$ ; and
- in the remainder of the formula, decision rule  $\beta_\tau^2$  at time  $\tau$  can be retrieved as a function of  $\mu^2$  (noted  $\beta_\tau^2(\mu^2)$ ).

Also, line 13 results from the observation that, while  $M_{|\sigma_\tau, \beta_\tau^1}^1$  and  $M_{|\sigma_\tau}^2$  allow to actually make decisions over different time intervals, we are here minimizing over  $\mu^2$  while maximizing over  $\mu^1$  a function that is linear in both input spaces. This amounts to solving some 2-player zero-sum normal-form game, hence the applicability of von Neumann's minimax theorem.

The above derivation tells us that the maximin value (the best outcome player 1 can guarantee whatever player 2's strategy) in the one-time-step game is thus the Nash equilibrium value (NEV) for the complete subgame from  $\tau$  onwards.  $\square$

## D Solving zs-OMGs

### D.1 Preliminary Properties

#### D.1.1 Properties of $T_m^1$ and $T_c^1$

The first two lemmas below present properties of  $T_m^1$  and  $T_c^1$  that will be useful afterwards.

**Lemma 3.**  $T_m^1(\sigma_\tau, \beta_\tau)$  is linear in  $\sigma_\tau$ ,  $\beta_\tau^1$ , and  $\beta_\tau^2$ .

*Proof.*

$$T_m^1(\sigma_\tau, \beta_\tau)(\theta_\tau^1, a^1, z^1) \quad (15)$$

$$\begin{aligned}
 &= \sum_{\theta_\tau^2, a^2, z^2} T(\sigma_\tau, \beta_\tau)((\theta_\tau^1, a^1, z^1), (\theta_\tau^2, a^2, z^2)) \quad (\text{from Eq. (1)}) \\
 &= \sum_{s', \theta_\tau^2, a^2, z^2} \beta_\tau^1(\theta_\tau^1, a^1) \beta_\tau^2(\theta_\tau^2, a^2) \sum_s P_{a^1, a^2}^{z^1, z^2}(s'|s) b(s|\theta_\tau^1, \theta_\tau^2) \sigma_\tau(\theta_\tau^1, \theta_\tau^2) \\
 &= \beta_\tau^1(\theta_\tau^1, a^1) \sum_{\theta_\tau^2, a^2} \beta_\tau^2(\theta_\tau^2, a^2) \sum_{s, s', z^2} P_{a^1, a^2}^{z^1, z^2}(s'|s) b(s|\theta_\tau^1, \theta_\tau^2) \sigma_\tau(\theta_\tau^1, \theta_\tau^2). \tag{16}
 \end{aligned}$$

□

**Lemma 4.**  $T_c^1(\sigma_\tau, \beta_\tau)$  is independent of  $\beta_\tau^1$  and  $\sigma_\tau^{m,1}$ .

*Proof.* See Wiggers [2015], Lemma 4.2.3. □

### D.1.2 Linearity and Lipschitz-continuity of $T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2)$

**Lemma 5.** At depth  $\tau$ ,  $T(\sigma_\tau, \beta_\tau)$  is linear in  $\beta_\tau^1$ ,  $\beta_\tau^2$ , and  $\sigma_\tau$ , where  $\beta_\tau = \langle \beta_\tau^1, \beta_\tau^2 \rangle$ . It is more precisely 1-Lipschitz-continuous (1-LC) in  $\sigma_\tau$  (in 1-norm), i.e., for any  $\sigma_\tau, \sigma'_\tau$ :

$$\|T(\sigma'_\tau, \beta_\tau) - T(\sigma_\tau, \beta_\tau)\|_1 \leq 1 \cdot \|\sigma'_\tau - \sigma_\tau\|_1.$$

*Proof.* Let  $\sigma$  be an occupancy state at time  $\tau$  and  $\beta_\tau$  be a decision rule. Then, as seen in the proof of Proposition 2.3, the next occupancy state  $\sigma' = T(\sigma, \beta_\tau)$  satisfies, for any  $s'$  and  $(\theta, \mathbf{a}, \mathbf{z})$ :

$$\begin{aligned}
 \sigma'(\theta, \mathbf{a}, \mathbf{z}) &\stackrel{\text{def}}{=} Pr(\theta, \mathbf{a}, \mathbf{z} | \sigma, \beta_\tau^1, \beta_\tau^2) \\
 &= \beta_\tau^1(\theta^1, a^1) \beta_\tau^2(\theta^2, a^2) \left[ \sum_{s', s \in \mathcal{S}} P_{\mathbf{a}}^{z'}(s'|s) b(s|\theta) \right] \sigma(\theta).
 \end{aligned}$$

$b(s|\theta)$  depending only on the model (transition function and initial belief), the next occupancy state  $\sigma'$  thus evolves linearly w.r.t. (i) *private* decision rules  $\beta_\tau^1$  and  $\beta_\tau^2$ , and (ii) the occupancy state  $\sigma$ .

The 1-Lipschitz-continuity holds because each component of vector  $\sigma_\tau$  is distributed over multiple components of  $\sigma'$ . Indeed, let us view two occupancy states as vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , and their corresponding next states under  $\beta_\tau$  as  $M\mathbf{x}$  and  $M\mathbf{y}$ , where  $M \in \mathbb{R}^{m \times n}$  is the corresponding transition matrix (i.e., which turns  $\sigma$  into  $\sigma' \stackrel{\text{def}}{=} T(\sigma_\tau, \beta_\tau)$ ). Then,

$$\begin{aligned}
 \|M\mathbf{x} - M\mathbf{y}\|_1 &\stackrel{\text{def}}{=} \sum_{j=1}^m \left| \sum_{i=1}^n M_{i,j} (x_i - y_i) \right| \\
 &\leq \sum_{j=1}^m \sum_{i=1}^n |M_{i,j} (x_i - y_i)| \quad (\text{convexity of } |\cdot|) \\
 &= \sum_{j=1}^m \sum_{i=1}^n M_{i,j} |x_i - y_i| \quad (\forall i, j, M_{i,j} \geq 0) \\
 &= \sum_{i=1}^n \underbrace{\sum_{j=1}^m M_{i,j}}_{=1} |x_i - y_i| \quad (M \text{ is a transition matrix}) \\
 &\stackrel{\text{def}}{=} \|\mathbf{x} - \mathbf{y}\|_1.
 \end{aligned}$$

□

### D.1.3 Lipschitz-Continuity of $V^*$

The next two results demonstrate that, in the finite horizon setting,  $V^*$  is Lipschitz-continuous (LC) in occupancy space, which allows defining LC upper- and lower-bound approximations.

**Lemma 6.** At depth  $\tau$ ,  $V_\tau(\sigma_\tau, \beta_\tau)$  is linear w.r.t.  $\sigma_\tau$  and  $\beta_\tau$ .

Note: This result in fact applies to any reward function of a general-sum POSG with any number of agents (here  $N$ ), e.g., to a Dec-POMDP. The following proof handles the general case (with  $\beta_\tau \stackrel{\text{def}}{=} \langle \beta_\tau^1, \dots, \beta_\tau^N \rangle$ , and  $\beta_\tau(\mathbf{a}|\theta) = \prod_{i=1}^N \beta_\tau^i(a^i, \theta^1)$ ).

*Proof.* This property trivially holds for  $\tau = H - 1$  because

$$\begin{aligned} V_{H-1}(\sigma_{H-1}, \beta_{H-1}) &= r(\sigma_{H-1}, \beta_{H-1}) \\ &= \sum_{s, \mathbf{a}} \left( \sum_{\theta} Pr(s, \mathbf{a} | \theta) \sigma_{H-1}(\theta) \right) r(s, \mathbf{a}) \\ &= \sum_{s, \mathbf{a}} \left( \sum_{\theta} b(s | \theta) \beta_{\tau}(\mathbf{a} | \theta) \sigma_{H-1}(\theta) \right) r(s, \mathbf{a}) \\ &= \sum_{s, \theta} b(s | \theta) \sigma_{H-1}(\theta) \left( \sum_{\mathbf{a}} \beta_{\tau}(\mathbf{a} | \theta) r(s, \mathbf{a}) \right). \end{aligned}$$

Now, let us assume that the property holds for  $\tau + 1 \in \{1 \dots H - 1\}$ . Then,

$$\begin{aligned} V_{\tau}(\sigma_{\tau}, \beta_{\tau}) &= \sum_{s, \mathbf{a}} \left( \sum_{\theta} b(s | \theta) \beta_{\tau}(\mathbf{a} | \theta) \sigma_{\tau}(\theta) \right) r(s, \mathbf{a}) + \gamma V_{\tau+1}(T(\sigma_{\tau}, \beta_{\tau}), \beta_{\tau+1}) \\ &= \sum_{s, \theta} b(s | \theta) \sigma_{\tau}(\theta) \left( \sum_{\mathbf{a}} \beta_{\tau}(\mathbf{a} | \theta) r(s, \mathbf{a}) \right) + \gamma V_{\tau+1}(T(\sigma_{\tau}, \beta_{\tau}), \beta_{\tau+1}). \end{aligned}$$

As

- $T(\sigma_{\tau}, \beta_{\tau})$  is linear in  $\sigma_{\tau}$  (Lemma 5) and
- $V_{\tau+1}(\sigma_{\tau+1}, \beta_{\tau+1})$  is linear in  $\sigma_{\tau+1}$  (induction hypothesis),

their composition,  $V_{\tau+1}(T(\sigma_{\tau}, \beta_{\tau}), \beta_{\tau+1})$ , is also linear in  $\sigma_{\tau}$ , and so is  $V_{\tau}(\sigma_{\tau}, \beta_{\tau})$ . Similarly,  $V_{\tau}(\sigma_{\tau}, \beta_{\tau})$  is linear in  $\beta_{\tau}$  for any  $\sigma_{\tau}$ .  $\square$

**Theorem 3.3.** (originally stated on page 7) *Let  $h_{\tau} \stackrel{\text{def}}{=} \frac{1-\gamma^{H-\tau}}{1-\gamma}$  (or  $h_{\tau} \stackrel{\text{def}}{=} H - \tau$  if  $\gamma = 1$ ). Then  $V_{\tau}^*(\sigma_{\tau})$  is  $\lambda_{\tau}$ -Lipschitz continuous in  $\sigma_{\tau}$  at any depth  $\tau \in \{0 \dots H - 1\}$ , where  $\lambda_{\tau} = \frac{1}{2} h_{\tau} (r_{\max} - r_{\min})$ .*

*Proof.* At depth  $\tau$ , the value of any behavioral strategy  $\beta_{\tau}$  is bounded, independently of  $\sigma_{\tau}$ , by

$$\begin{aligned} V_{\tau}^{\max} &\stackrel{\text{def}}{=} h_{\tau} r_{\max}, \quad \text{where } r_{\max} \stackrel{\text{def}}{=} \max_{s, \mathbf{a}} r(s, \mathbf{a}), \text{ and} \\ V_{\tau}^{\min} &\stackrel{\text{def}}{=} h_{\tau} r_{\min}, \quad \text{where } r_{\min} \stackrel{\text{def}}{=} \min_{s, \mathbf{a}} r(s, \mathbf{a}). \end{aligned}$$

Thus,  $V_{\beta_{\tau}}$  being a linear function defined over a probability simplex ( $\mathcal{O}_{\tau}^{\sigma}$ ) (cf. Appendix D.1.3) and bounded by  $[V_{\tau}^{\min}, V_{\tau}^{\max}]$ , we can apply Horák's PhD thesis' Lemma 3.5 (p. 33) 2019 to establish that it is also  $\lambda_{\tau}$ -LC, i.e.,

$$\begin{aligned} |V_{\beta_{\tau}}(\sigma) - V_{\beta_{\tau}}(\sigma')| &\leq \lambda_{\tau} \|\sigma - \sigma'\|_1 \quad (\forall \sigma, \sigma'), \\ \text{with } \lambda_{\tau} &= \frac{V_{\tau}^{\max} - V_{\tau}^{\min}}{2}. \end{aligned}$$

Considering now optimal solutions, this means that, at depth  $\tau$  and for any  $(\sigma, \sigma') \in \mathcal{O}_{\tau}^{\sigma}$ :

$$\begin{aligned} &V_{\tau}^*(\sigma) - V_{\tau}^*(\sigma') \\ &= \max_{\beta_{\tau}^1} \min_{\beta_{\tau}^2} V_{\tau}(\sigma, \beta_{\tau}^1, \beta_{\tau}^2) - \max_{\beta_{\tau}^1} \min_{\beta_{\tau}^2} V_{\tau}(\sigma', \beta_{\tau}^1, \beta_{\tau}^2) \\ &\leq \max_{\beta_{\tau}^1} \min_{\beta_{\tau}^2} [V_{\tau}(\sigma', \beta_{\tau}^1, \beta_{\tau}^2) + \lambda_{\tau} \|\sigma - \sigma'\|_1] - \max_{\beta_{\tau}^1} \min_{\beta_{\tau}^2} V_{\tau}(\sigma', \beta_{\tau}^1, \beta_{\tau}^2) \\ &= \lambda_{\tau} \|\sigma - \sigma'\|_1. \end{aligned}$$

Symmetrically,  $V_{\tau}^*(\sigma) - V_{\tau}^*(\sigma') \geq -\lambda_{\tau} \|\sigma - \sigma'\|_1$ , hence the expected result:

$$|V_{\tau}^*(\sigma) - V_{\tau}^*(\sigma')| \leq \lambda_{\tau} \|\sigma - \sigma'\|_1.$$

$\square$

As it will be used later, let us also present the following lemma.

**Lemma 7.** Let us consider  $\tau \in \{0 \dots H - 1\}$ ,  $\theta_\tau^1$ , and  $\psi_\tau^2$ . Then  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2(\theta_\tau^1)$  is  $\lambda_\tau$ -LC in  $\sigma_\tau^{c,1}(\cdot|\theta_\tau^1)$ .

Equivalently, we will also write that  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2$  is  $\lambda_\tau$ -LC in  $\sigma_\tau^{c,1}$  in vector-wise 1-norm, i.e.:

$$\overrightarrow{|\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2 - \nu_{[\tilde{\sigma}_\tau^{c,1}, \psi_\tau^2]}^2|}_1 \preceq \lambda_\tau \overrightarrow{\|\sigma_\tau^{c,1} - \tilde{\sigma}_\tau^{c,1}\|}_1,$$

where (i) the absolute value of a vector is obtained by taking the absolute value of each component; and (ii) the vector-wise 1-norm of a matrix is a vector made of the 1-norm of each of its component vectors.

*Proof.* For any  $\theta_\tau^1$ ,  $\sigma_\tau^{c,1}$  and  $\psi_\tau^2$  induce a POMDP for Player 1 from  $\tau$  on, where (i) the state at any  $t \in \{\tau \dots H - 1\}$  corresponds to a pair  $\langle s, \theta_t^2 \rangle$ , and (ii) the initial belief is derived from  $\sigma_\tau^{c,1}(\cdot|\theta_\tau^1)$ . The belief state at  $t$  thus gives:

$$b_{\theta_t^1}(s, \theta_t^2) \stackrel{\text{def}}{=} \Pr(s, \theta_t^2 | \theta_t^1) = \underbrace{\Pr(s | \theta_t^2, \theta_t^1)}_{b_{\theta_t^2, \theta_t^1}^{\text{HMM}}(s)} \cdot \underbrace{\Pr(\theta_t^2 | \theta_t^1)}_{\sigma_t^{c,1}(\theta_t^2 | \theta_t^1)}.$$

So,

- the value function of any behavioral strategy  $\beta_\tau^1$  is linear at  $t$  in  $b_{\theta_t^1}$ , thus (in particular) in  $\sigma_t^{c,1}(\cdot|\theta_t^1)$ ; and
- the optimal value function is LC at  $t$  also in  $b_{\theta_t^1}$  (with the same depth-dependent upper-bounding Lipschitz constant  $\lambda_t$  as in the proof of Theorem 3.3),<sup>5</sup> thus (in particular) in  $\sigma_t^{c,1}(\cdot|\theta_t^1)$ .

Using  $t = \tau$ , the optimal value function is  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2(\theta_\tau^1)$ , which is thus  $\lambda_\tau$ -LC in  $\sigma_\tau^{c,1}(\cdot|\theta_\tau^1)$ .  $\square$

## D.2 Bounding Approximations of $V^*$ , $W^{1,*}$ and $W^{2,*}$

### D.2.1 $\bar{V}_\tau$ and $\underline{V}_\tau$

To find a form that could be appropriate for an upper bound approximation of  $V_\tau^*$ , let us consider an OS  $\sigma_\tau$  and a single tuple  $\langle \tilde{\sigma}_\tau, \nu_{[\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2]}^2 \rangle$ , and define  $\zeta_\tau \stackrel{\text{def}}{=} \sigma_\tau^{m,1} \tilde{\sigma}_\tau^{c,1}$ . Then,

$$\begin{aligned} V^*(\sigma_\tau) &\leq V^*(\zeta_\tau) + \lambda_\tau \|\sigma_\tau - \zeta_\tau\|_1 && \text{(LC, cf. Theorem 3.3)} \\ &= V^*(\sigma_\tau^{m,1} \tilde{\sigma}_\tau^{c,1}) + \lambda_\tau \|\sigma_\tau - \zeta_\tau\|_1 \\ &\leq \sigma_\tau^{m,1} \cdot \nu_{[\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2]}^2 + \lambda_\tau \|\sigma_\tau - \sigma_\tau^{m,1} \tilde{\sigma}_\tau^{c,1}\|_1. && \text{(Cvx, cf. Theorem 2.5)} \end{aligned}$$

Notes:

- $\tilde{\sigma}_\tau^{m,1}$  does not appear in the resulting upper bound, thus will not need to be specified.
- For  $\tau = H - 1$ ,  $\nu_{[\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2]}^2$  is a simple function of  $r$ ,  $\tilde{\sigma}_\tau^{c,1}$ ,  $\beta_\tau^2$ , and the dynamics of the system, as described in Eq. (9) of Wiggers et al. [2016a].

From this, we can deduce the following appropriate forms of upper and (symmetrically) lower bound function approximations for  $V_\tau^*$ :

$$\begin{aligned} \bar{V}_\tau(\sigma_\tau) &= \min_{\langle \tilde{\sigma}_\tau^{c,1}, \langle \bar{\nu}_\tau^2, \psi_\tau^2 \rangle \rangle \in \bar{\mathcal{J}}_\tau} \left[ \sigma_\tau^{m,1} \cdot \bar{\nu}_\tau^2 + \lambda_\tau \|\sigma_\tau - \sigma_\tau^{m,1} \tilde{\sigma}_\tau^{c,1}\|_1 \right], \text{ and} \\ \underline{V}_\tau(\sigma_\tau) &= \max_{\langle \tilde{\sigma}_\tau^{c,2}, \langle \underline{\nu}_\tau^1, \psi_\tau^1 \rangle \rangle \in \underline{\mathcal{J}}_\tau} \left[ \sigma_\tau^{m,2} \cdot \underline{\nu}_\tau^1 - \lambda_\tau \|\sigma_\tau - \sigma_\tau^{m,2} \tilde{\sigma}_\tau^{c,2}\|_1 \right], \end{aligned}$$

which are respectively concave in  $\sigma_\tau^{m,1}$  and convex in  $\sigma_\tau^{m,2}$ , and which both exploit the Lipschitz continuity.

### D.2.2 $\bar{W}_\tau$ and $\underline{W}_\tau$

Note: We discuss all depths from 0 to  $H - 1$ , even though we do not need these approximations at  $\tau = H - 1$ .

Let us first see how concavity-convexity properties affect  $W_\tau^{*,1}$ .

**Lemma 8.** Considering that vectors  $\nu_{[\sigma_H^{c,1}, \beta_H^2]}^2$  are null vectors, we have, for all  $\tau \in \{0 \dots H - 1\}$ :

$$\begin{aligned} W_\tau^{1,*}(\sigma_\tau, \beta_\tau^1) &= \min_{\beta_\tau^2, \langle \beta_{\tau+1}^2, \nu_{[T_\tau^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]}^2 \rangle} \beta_\tau^1 \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) \right. \\ &\quad \left. + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \nu_{[T_\tau^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]}^2 \right]. \end{aligned}$$

<sup>5</sup>The proof process is similar. The only difference lies in the space at hand, but without any impact on the resulting formulas.

*Proof.* Considering that vectors  $\nu_{[\sigma_H^{c,1}, \beta_H^2]}$  are null vectors, we have, for all  $\tau \in \{0 \dots H-1\}$ :

$$W_\tau^{1,*}(\sigma_\tau, \beta_\tau^1) = \min_{\beta_\tau^2} Q_\tau^*(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) = \min_{\beta_\tau^2} [r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau))]$$

(Line below exploits Theorem 2.5 (p. 6) and  $T_c^1$ 's independence from  $\beta_\tau^1$  (Lemma 4).)

$$\begin{aligned} &= \min_{\beta_\tau^2} \left[ r(\sigma_\tau, \beta_\tau) + \gamma \min_{\langle \beta_{\tau+1}^2, \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]} \rangle} \left[ T_m^1(\sigma_\tau, \beta_\tau) \cdot \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]}^2 \right] \right] \\ &= \min_{\beta_\tau^2, \langle \beta_{\tau+1}^2, \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]} \rangle} \left[ r(\sigma_\tau, \beta_\tau) + \gamma T_m^1(\sigma_\tau, \beta_\tau) \cdot \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]}^2 \right] \end{aligned}$$

(Line below exploits  $r$  and  $T_m^1$ 's linearity in  $\beta_\tau^1$  (Lemma 3).)

$$\begin{aligned} &= \min_{\beta_\tau^2, \langle \beta_{\tau+1}^2, \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]} \rangle} \beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) \right. \\ &\quad \left. + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \nu_{[T_c^1(\sigma_\tau, \beta_\tau^2), \beta_{\tau+1}^2]}^2 \right]. \end{aligned}$$

□

Note that, since  $V_H^* = 0$ ,  $\tau = H-1$  is a particular case which can be simply re-written:

$$W_\tau^{1,*}(\sigma_\tau, \beta_\tau^1) = \min_{\beta_\tau^2} \beta_\tau^{1\top} \cdot r(\sigma_\tau, \cdot, \beta_\tau^2).$$

To find a form that could be appropriate for an upper bound approximation of  $W_\tau^{*,1}$ , let us now consider an OS  $\sigma_\tau$  and a single tuple  $\langle \tilde{\sigma}_\tau, \tilde{\beta}_\tau^2, \nu_{[T_c^1(\tilde{\sigma}_\tau, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]} \rangle$ . Then,

$$\begin{aligned} W_\tau^{1,*}(\sigma_\tau, \beta_\tau^1) &= \min_{\beta_\tau^2} [r(\sigma_\tau, \beta_\tau) + \gamma V_{\tau+1}^*(T(\sigma_\tau, \beta_\tau))] \\ &\leq r(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) + \gamma V_{\tau+1}^{BR,1}(T(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) | \tilde{\beta}_{\tau+1}^2) \end{aligned}$$

(Use  $\tilde{\beta}_\tau^2$  &  $\tilde{\beta}_{\tau+1}^2$ : instead of mins)

(where  $V_{\tau+1}^{BR,1}(T(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) | \tilde{\beta}_{\tau+1}^2)$  is the value of 1's best response to  $\tilde{\beta}_{\tau+1}^2$ : if in  $T(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2)$ )

$$= r(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) + \gamma T_m^1(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) \cdot \underbrace{\nu_{[T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]}^2}_{\text{value of 1's best response to } \tilde{\beta}_{\tau+1}^2}$$

(Lem. 3 of Wiggers et al. [2016a])

$$\begin{aligned} &\leq r(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) + \gamma T_m^1(\sigma_\tau, \beta_\tau^1, \tilde{\beta}_\tau^2) \cdot \left( \nu_{[T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]}^2 \right. \\ &\quad \left. + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2) - T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2)\|_1} \right) \end{aligned} \tag{17}$$

(Lem. 7:  $\lambda_{\tau+1}$ -LC of  $\nu_{[T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]}^2$ )

$$\begin{aligned} &= \beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) \cdot \left( \nu_{[T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]}^2 \right. \right. \\ &\quad \left. \left. + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2) - T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2)\|_1} \right) \right] \end{aligned} \tag{18}$$

(Linearity in  $\beta_\tau^1$ )

$$\begin{aligned} &= \beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) \cdot \nu_{[T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2), \tilde{\beta}_{\tau+1}^2]}^2 \right. \\ &\quad \left. + \gamma \lambda_{\tau+1} \cdot \|T(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) - T_m^1(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2)\|_1 \right] \end{aligned}$$

(Alternative writing)

From this, we can deduce the following appropriate forms of (i) upper bounding approximation for  $W_\tau^{1,*}$  and (ii) (symmetrically) of lower bound approximation for  $W_\tau^{2,*}$ :

$$\begin{aligned}\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) &= \min_{\langle \tilde{\sigma}_\tau^{c,1}, \beta_\tau^2, \bar{\nu}_{\tau+1}^2 \rangle \in \overline{\mathcal{I}}_\tau} \beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \bar{\nu}_{\tau+1}^2 \right. \\ &\quad \left. + \gamma \lambda_{\tau+1} \cdot \|T(\sigma_\tau, \cdot, \beta_\tau^2) - T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) T_c^1(\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1 \right], \text{ and} \\ \underline{W}_\tau(\sigma_\tau, \beta_\tau^2) &= \max_{\langle \tilde{\sigma}_\tau^{c,2}, \beta_\tau^1, \underline{\nu}_{\tau+1}^1 \rangle \in \underline{\mathcal{I}}_\tau} \beta_\tau^{2\top} \cdot \left[ r(\sigma_\tau, \beta_\tau^1, \cdot) + \gamma T_m^2(\sigma_\tau, \beta_\tau^1, \cdot) \cdot \underline{\nu}_{\tau+1}^1 \right. \\ &\quad \left. - \gamma \lambda_{\tau+1} \cdot \|T(\sigma_\tau, \beta_\tau^1, \cdot) - T_m^2(\sigma_\tau, \beta_\tau^1, \cdot) T_c^2(\tilde{\sigma}_\tau^{c,2}, \beta_\tau^1)\|_1 \right],\end{aligned}$$

where  $\bar{\nu}_{\tau+1}^2$  and  $\underline{\nu}_{\tau+1}^1$  respectively upper and lower bound the actual vectors associated to the players' future strategies (resp. of 2 and 1).

Again,  $\tau = H - 1$  is a particular case where only the reward term is preserved.

This constitutes the proof to the following proposition

**Proposition 3.4.** (originally stated on page 8) *Let  $\overline{\mathcal{I}}_\tau$  be a set of tuples  $w = \langle \tilde{\sigma}_\tau, \beta_\tau^2, \langle \bar{\nu}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle$ . Then,*

$$\begin{aligned}\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) &\stackrel{\text{def}}{=} \min_{\langle \tilde{\sigma}_\tau, \beta_\tau^2, \langle \bar{\nu}_{\tau+1}^2, \beta_{\tau+1}^2 \rangle \rangle \in \overline{\mathcal{I}}_\tau} \left[ r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) \cdot \bar{\nu}_{\tau+1}^2 \right. \\ &\quad \left. + \lambda_{\tau+1} \|T(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) - T_m^1(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) T_c^1(\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1 \right]\end{aligned}\quad (6)$$

upper-bounds  $W_\tau^{1,*}$  over the whole space  $\mathcal{O}_\tau^\sigma \times \mathcal{B}_\tau^1$ .

### D.3 Related Operators

#### D.3.1 Selection Operator: Solving for $\beta_\tau^1$ as an LP

**Proposition D.1.** *Using now a distribution  $\psi_\tau^2$  over tuples  $w = \langle \tilde{\sigma}_\tau^{c,1}, \beta_\tau^2, \bar{\nu}_{\tau+1}^2 \rangle \in \overline{\mathcal{I}}_\tau$ , the corresponding upper-bounding value for “profile”  $\langle \beta_\tau^1, \psi_\tau^2 \rangle$  when in  $\sigma_\tau$  can be written as an expectancy:*

$$\beta_\tau^{1\top} \cdot M^{\sigma_\tau} \cdot \psi_\tau^2,$$

where  $M^{\sigma_\tau}$  is an  $|\Theta_\tau^1 \times \mathcal{A}^1| \times |\overline{\mathcal{I}}_\tau^1|$  matrix.

*Proof.* From the right-hand side term in (Equation (18)), the upper-bounding value associated to  $\sigma_\tau$ ,  $\beta_\tau^1$  and a tuple  $\langle \tilde{\sigma}_\tau^{c,1}, \beta_\tau^2, \bar{\nu}_{\tau+1}^2 \rangle \in \overline{\mathcal{I}}_\tau$  can be written:

$$\beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \left( \bar{\nu}_{\tau+1}^2 + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \tilde{\beta}_\tau^2) - T_c^1(\tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2)\|_1} \right) \right].$$

Using now a distribution  $\psi_\tau^2$  over tuples  $w = \langle \tilde{\sigma}_\tau^{c,1}, \beta_\tau^2, \bar{\nu}_{\tau+1}^2 \rangle \in \overline{\mathcal{I}}_\tau$ , the corresponding upper-bounding value for “profile”  $\langle \beta_\tau^1, \psi_\tau^2 \rangle$  when in  $\sigma_\tau$  can be written as an expectancy:

$$\begin{aligned}\sum_{w \in \overline{\mathcal{W}}_\tau} \beta_\tau^{1\top} \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2[w]) + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2[w]) \cdot \left( \bar{\nu}_{\tau+1}^2[w] \right. \right. \\ \left. \left. + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2[w]) - T_c^1(\tilde{\sigma}_\tau^{c,1}[w], \beta_\tau^2[w])\|_1} \right) \right] \cdot \psi_\tau^2(w)\end{aligned}$$

(where  $x[w]$  denotes the field  $x$  of tuple  $w$ )

$$= \beta_\tau^{1\top} \cdot M^{\sigma_\tau} \cdot \psi_\tau^2,$$

where  $M^{\sigma_\tau}$  is an  $|\Theta_\tau^1 \times \mathcal{A}^1| \times |\overline{\mathcal{I}}_\tau^1|$  matrix. □

For implementation purposes, using Eqs. (2) and (16) (to develop respectively  $r(\cdot, \cdot, \cdot)$  and  $T_m^1(\cdot, \cdot, \cdot)$ ), we can derive the expression of a component, *i.e.*, the upper-bounding value if  $a^1$  is applied in  $\theta_\tau^1$  while  $w$  is chosen:

$$\begin{aligned}M_{(\theta_\tau^1, a^1), w}^{\sigma_\tau} &\stackrel{\text{def}}{=} r(\sigma_\tau, \cdot, \beta_\tau^2[w]) + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2[w]) \cdot \\ &\quad \left( \bar{\nu}_{\tau+1}^2[w] + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2[w]) - T_c^1(\tilde{\sigma}_\tau^{c,1}[w], \beta_\tau^2[w])\|_1} \right)\end{aligned}$$

$$\begin{aligned}
 &= \sum_{s, \theta_\tau^2, a^2} \sigma_\tau(\theta_\tau) b(s|\theta_\tau) \beta_\tau^2[w](a^2|\theta_\tau^2) r(s, \mathbf{a}) \\
 &\quad + \gamma \sum_{z^1} \left[ \sum_{\theta_\tau^2, a^2} \beta_\tau^2[w](a^2|\theta_\tau^2) \sum_{s, s', z^2} P_{\mathbf{a}}^z(s'|s) b(s|\theta_\tau) \sigma_\tau(\theta_\tau) \right] \\
 &\quad \left( \bar{v}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1) \right. \\
 &\quad \left. + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2[w]) - T_c^1(\tilde{\sigma}_\tau^{c,1}[w], \beta_\tau^2[w])\|_1}(\theta_\tau^1, a^1, z^1)} \right) \\
 &= \sum_{\theta_\tau^2} \sigma_\tau(\theta_\tau) \sum_{a^2} \beta_\tau^2[w](a^2|\theta_\tau^2) \\
 &\quad \cdot \left( \sum_s b(s|\theta_\tau) r(s, \mathbf{a}) \right. \\
 &\quad \left. + \gamma \sum_{z^1} \left[ \sum_{s, s', z^2} P_{\mathbf{a}}^z(s'|s) b(s|\theta_\tau) \right] \cdot \left( \bar{v}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1) \right. \right. \\
 &\quad \left. \left. + \lambda_{\tau+1} \cdot \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2[w]) - T_c^1(\tilde{\sigma}_\tau^{c,1}[w], \beta_\tau^2[w])\|_1}(\theta_\tau^1, a^1, z^1)} \right) \right).
 \end{aligned}$$

Then, solving  $\max_{\beta_\tau^1} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1)$  can be rewritten as solving a zero-sum game where pure strategies are:

- for Player 1, the choice of not 1, but  $|\Theta_\tau^1|$  actions (among  $|\mathcal{A}^1|$ ) and,
- for Player 2, the choice of 1 element of  $\bar{\mathcal{I}}_\tau^1$ .

One can view it as a Bayesian game with one type per history  $\theta_\tau^1$  for 1, and a single type for 2.

With our upper bound,  $\max_{\beta_\tau^1} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1)$  can thus be solved as the LP:

$$\begin{aligned}
 \max_{\beta_\tau^1, v} v \quad \text{s.t.} \quad & \text{(i)} \quad \forall w \in \bar{\mathcal{I}}_\tau^1, \quad v \leq \beta_\tau^1{}^\top \cdot M_{(\cdot, w)}^{\sigma_\tau} \\
 & \text{(ii)} \quad \forall \theta_\tau^1 \in \Theta_\tau^1, \quad \sum_{a^1} \beta_\tau^1(a^1|\theta_\tau^1) = 1,
 \end{aligned}$$

whose dual LP is given by

$$\begin{aligned}
 \min_{\psi_\tau^2, v} v \quad \text{s.t.} \quad & \text{(i)} \quad \forall (\theta_\tau^1, a^1) \in \Theta_\tau^1 \times \mathcal{A}^1, \quad v \geq M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2 \\
 & \text{(ii)} \quad \sum_{w \in \bar{\mathcal{I}}_\tau^1} \psi_\tau^2(w) = 1.
 \end{aligned}$$

As can be noted,  $M^{\sigma_\tau}$ 's columns corresponding to 0-probability histories  $\theta_\tau^1$  in  $\sigma_\tau^{m,1}$  are empty (full of zeros), so that the corresponding decision rules (for these histories) are not relevant and can be set arbitrarily. The actual implementation thus ignores these histories, whose corresponding decision rules also do not need to be stored.

*Remark D.2* (Interpretation of  $M^{\sigma_\tau}$ ). The content of this matrix can be interpreted by noting that, a given  $w$  containing a behavioral strategy  $\beta_{\tau:H-1}^2$  and an OS  $\tilde{\sigma}_\tau$ , a pair  $\langle w, \theta_\tau^1 \rangle$  induces a POMDP for player 1 whose state space is made of pairs  $\langle s, \theta_t^2 \rangle$ , and whose initial belief  $b_\tau$  depends on  $\tilde{\sigma}_\tau$  and  $\theta_\tau^1$ . Solving this POMDP amounts to finding a best response of player 1 to  $\beta_{\tau:H-1}^2$ . In this setting, an element  $M_{((\theta_\tau^1, a^1), w)}^{\sigma_\tau}$  is an upper-bound of the optimal (POMDP)  $Q$ -value when player 1 performs  $a^1$  while facing  $b_\tau$  ( $Q_{\text{POMDP}}^*(b_\tau, a^1)$ ).

### D.3.2 Upper Bounding $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}$

Adding a new complete tuple to  $\bar{\mathcal{I}}_\tau^1$  requires a new vector  $\bar{v}_\tau^2$  that upper bounds the vector  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2$  associated to the strategy induced by  $\psi_\tau^2$ . We can obtain one in a recursive manner (not solving the induced POMDP).

**Proposition D.3.** *For each  $\psi_\tau^2$  obtained as the solution of the aforementioned (dual) LP in  $\sigma_\tau$ , and each  $\theta_\tau^1$ ,  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2(\theta_\tau^1)$  is upper bounded by a value  $\bar{v}_\tau^2(\theta_\tau^1)$  that depends on vectors  $\bar{v}_{\tau+1}^2$  in the support of  $\psi_\tau^2$ . In particular, if  $\theta_\tau^1 \in \text{Supp}(\sigma_\tau^{m,1})$ , we have:*

$$\bar{v}_\tau^2(\theta_\tau^1) \stackrel{\text{def}}{=} \frac{1}{\sigma_{\tau,m}^1(\theta_\tau^1)} \max_{a^1 \in \mathcal{A}^1} M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2.$$

*Proof.* For a newly derived  $\psi_\tau^2$ , as  $\nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2(\theta_\tau^1)$  is the value of 1's best action ( $\in \mathcal{A}^1$ ) if 1 (i) observes  $\theta_\tau^1$  while in  $\sigma_\tau^{c,1}$  and (ii) 2 plays  $\psi_\tau^2$ , we have:

$$\begin{aligned}
 & \nu_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^2(\theta_\tau^1) \stackrel{\text{def}}{=} V_{[\sigma_\tau^{c,1}, \psi_\tau^2]}^*(\theta_\tau^1) \quad (\text{optimal POMDP value function}) \\
 &= \max_{\beta_\tau^1} \mathbb{E} \left[ \sum_{t=\tau}^H \gamma^{t-\tau} R_t \mid \beta_\tau^1, \theta_\tau^1, \sigma_\tau^{c,1}, \psi_\tau^2 \right] \\
 &= \max_{a^1} \mathbb{E} \left[ R_\tau + \gamma \max_{\beta_{\tau+1}^1} \mathbb{E} \left[ \sum_{t=\tau+1}^H \gamma^{t-(\tau+1)} R_t \mid \beta_{\tau+1}^1, \langle \theta_\tau^1, a^1, Z^1 \rangle, \sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2 \right] \right. \\
 & \quad \left. \mid a^1, \theta_\tau^1, \sigma_\tau^{c,1}, \psi_\tau^2 \right] \\
 &= \max_{a^1} \mathbb{E} \left[ R_\tau + \gamma V_{[\sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2]}^*(\theta_\tau^1, a^1, Z^1) \mid a^1, \theta_\tau^1, \sigma_\tau^{c,1}, \psi_\tau^2 \right] \\
 &= \max_{a^1} \sum_{w, \theta_\tau^2, a^2, z^1} \underbrace{Pr(w, \theta_\tau^2, z^1, a^2 \mid a^1, \theta_\tau^1, \sigma_\tau^{c,1}, \psi_\tau^2)} \\
 & \quad \cdot \left( r(\theta_\tau, \mathbf{a}) + \gamma \nu_{[\sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2[w]]}^2(\theta_\tau^1, a^1, z^1) \right)
 \end{aligned}$$

(where  $\sigma_{\tau+1}^{c,1} = T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2[w])$  (Lemma 4, p. 29))

$$\begin{aligned}
 &= \max_{a^1} \sum_{w, \theta_\tau^2, a^2, z^1} \underbrace{Pr(w \mid \psi_\tau^2)} \cdot \underbrace{Pr(\theta_\tau^2 \mid \theta_\tau^1, \sigma_\tau^{c,1})} \cdot \underbrace{Pr(a^2 \mid \beta_\tau^2[w], \theta_\tau^2)} \cdot \underbrace{Pr(z^1 \mid \theta_\tau, \mathbf{a})} \\
 & \quad \cdot \left( r(\theta_\tau, \mathbf{a}) + \gamma \nu_{[\sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2[w]]}^2(\theta_\tau^1, a^1, z^1) \right) \\
 &= \max_{a^1} \sum_w \psi_\tau^2(w) \sum_{\theta_\tau^2} \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) \sum_{a^2} \beta_\tau^2[w](a^2 \mid \theta_\tau^2) \\
 & \quad \cdot \left( r(\theta_\tau, \mathbf{a}) + \gamma \sum_{z^1} Pr(z^1 \mid \theta_\tau, \mathbf{a}) \underbrace{\nu_{[\sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2[w]]}^2(\theta_\tau^1, a^1, z^1)} \right)
 \end{aligned}$$

then, as  $\nu_{[\sigma_{\tau+1}^{c,1}, \psi_{\tau+1}^2[w]]}^2$  is  $\lambda_{\tau+1}$ -LC in (any)  $\sigma_{\tau+1}^{c,1}$  (Lemma 7),

$$\begin{aligned}
 &\leq \max_{a^1} \sum_w \psi_\tau^2(w) \sum_{\theta_\tau^2} \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) \sum_{a^2} \beta_\tau^2[w](a^2 \mid \theta_\tau^2) \cdot \left( r(\theta_\tau, \mathbf{a}) + \gamma \sum_{z^1} Pr(z^1 \mid \theta_\tau, \mathbf{a}) \right. \\
 & \quad \left. \cdot \left[ \underbrace{\nu_{[\sigma_{\tau+1}^{c,1}[w], \psi_{\tau+1}^2[w]]}^2(\theta_\tau^1, a^1, z^1)} + \lambda_{\tau+1} \overrightarrow{\|\sigma_{\tau+1}^{c,1} - \tilde{\sigma}_{\tau+1}^{c,1}[w]\|_1}(\theta_\tau^1, a^1, z^1)} \right] \right) \\
 &\leq \max_{a^1} \sum_w \psi_\tau^2(w) \sum_{\theta_\tau^2} \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) \sum_{a^2} \beta_\tau^2[w](a^2 \mid \theta_\tau^2) \cdot \left( \underbrace{r(\theta_\tau, \mathbf{a}) + \gamma \sum_{z^1} Pr(z^1 \mid \theta_\tau, \mathbf{a})}_{\bar{\nu}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1)} \right. \\
 & \quad \left. \cdot \left[ \overrightarrow{\bar{\nu}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1)} + \lambda_{\tau+1} \overrightarrow{\|\sigma_{\tau+1}^{c,1} - \tilde{\sigma}_{\tau+1}^{c,1}[w]\|_1}(\theta_\tau^1, a^1, z^1)} \right] \right) \\
 &= \max_{a^1} \sum_w \psi_\tau^2(w) \sum_{\theta_\tau^2} \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) \sum_{a^2} \beta_\tau^2[w](a^2 \mid \theta_\tau^2) \cdot \left( \sum_s \underbrace{b(s \mid \theta_\tau) r(s, \mathbf{a})}_{\bar{\nu}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1)} \right. \\
 & \quad \left. + \gamma \sum_{z^1} \left( \sum_s \underbrace{b(s \mid \theta_\tau) Pr(z^1 \mid s, \mathbf{a})}_{\bar{\nu}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1)} \right) \cdot \left[ \overrightarrow{\bar{\nu}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1)} \right. \right. \\
 & \quad \left. \left. + \lambda_{\tau+1} \overrightarrow{\|\sigma_{\tau+1}^{c,1} - \tilde{\sigma}_{\tau+1}^{c,1}[w]\|_1}(\theta_\tau^1, a^1, z^1)} \right] \right) \\
 &= \max_{a^1} \sum_w \psi_\tau^2(w) \sum_{\theta_\tau^2} \sigma_\tau^{c,1}(\theta_\tau^2 \mid \theta_\tau^1) \sum_{a^2} \beta_\tau^2[w](a^2 \mid \theta_\tau^2)
 \end{aligned}$$

$$\begin{aligned}
 & \cdot \left( \sum_s b(s|\theta_\tau) r(s, \mathbf{a}) + \gamma \sum_{z^1} \left( \sum_{s, s', z^2} b(s|\theta_\tau) P_\alpha^z(s'|s) \right) \cdot \left[ \bar{v}_{\tau+1}^2[w](\theta_\tau^1, a^1, z^1) \right. \right. \\
 & \left. \left. + \lambda_{\tau+1} \left\| T_c^1(\bar{\sigma}_\tau^{c,1}, \beta_\tau^2[w]) - T_c^1(\bar{\sigma}_\tau^{c,1}[w], \beta_\tau^2[w]) \right\|_1(\theta_\tau^1, a^1, z^1) \right] \right) \\
 & = \frac{1}{\sigma_{\tau, m}^1(\theta_\tau^1)} \max_{a^1 \in \mathcal{A}^1} M_{((\theta_\tau^1, a^1), \cdot)}^{\sigma_\tau} \cdot \psi_\tau^2.
 \end{aligned}$$

□

### D.3.3 Strategy Conversion

Firstly, we give details regarding solutions of Dual LPs (Equation (7)) inducing behavioral strategies. As suggested in Section 3.2.2, one can show by induction that for any timestep  $\tau$ , each  $\psi_\tau^2$  is actually equivalent to an element of  $\Delta(\mathcal{B}_\tau^2)$ . The following lemma shows that for any timestep  $\tau$ , each element of  $\Delta(\mathcal{B}_\tau^2)$  induces an element of  $\mathcal{B}_\tau^2$ .

**Lemma 9.** *Each  $\psi_\tau^i \in \Delta(\mathcal{B}_\tau^i)$  induces a behavioral strategy. More precisely, we prove that (i) there is a natural injection from the set  $\mathcal{B}_\tau^i$  to the set of distributions  $\Delta(\mathcal{B}_\tau^i)$  and (ii) there is a surjection from the set  $\Delta(\mathcal{B}_\tau^i)$  to  $\mathcal{B}_\tau^i$ .*

*Proof.* By induction on  $\tau \in \{0, \dots, H-1\}$ , we prove that  $\cup_{t=\tau}^H \mathcal{B}_t^i \subset \cup_{t=\tau}^H \Delta(\mathcal{B}_t^i)$ . Firstly, for  $\tau = H-1$ , for all  $\beta_{H-1} \in \mathcal{B}_{H-1}$ , one can pick the degenerate distribution  $\psi_{H-1} = \beta_{H-1}$  which is in  $\Delta(\mathcal{B}_{H-1}^i)$ . Next, assume that  $\cup_{t=\tau+1}^{H-1} \mathcal{B}_t^i \subset \cup_{t=\tau+1}^{H-1} \Delta(\mathcal{B}_t^i)$  for some  $\tau \in \{0, \dots, H-2\}$ , then for all  $\beta_{\tau:H-1}, \beta_{\tau:H-1} = \beta_\tau \oplus \beta_{\tau+1:H-1}$ . By the induction hypothesis, there is  $\psi_{\tau+1:H-1} \in \Delta(\mathcal{B}_{\tau+1:H-1}^i)$  equal to  $\beta_{\tau+1:H-1}$ . Thus, we define  $\psi_{\tau:H-1} = \beta_\tau \oplus \psi_{\tau+1:H-1} = \beta_\tau \oplus \beta_{\tau+1:H-1} = \beta_{\tau:H-1}$  which is in  $\mathcal{B}_{\tau:H-1}$ . From this follows a natural injection from the behavioral strategies' set to the set of distributions over behavioral strategies.

The surjection from  $\Delta(\mathcal{B}_\tau^i)$  to  $\mathcal{B}_\tau^i$  is given by the realization weight computation algorithm detailed in Algorithm 2. □

As discussed in Appendix D.3.3, no effort is required to extract a solution strategy for a player from the lower bound (for 1) or the upper bound (for 2), but that strategy is in an unusual recursive form. We will here see (in the finite horizon setting) how to derive a (unique) equivalent behavioral strategy  $\beta_0^i$  using realization weights [Koller et al., 1994] in intermediate steps. To that end, we first define these realization weights in the case of a behavioral strategy (rather than for a mixed strategy as done by Koller et al.) and present some useful properties.

**About Realization Weights** Let us denote  $rw^i(a_0^i, z_1^i, a_1^i, \dots, a_\tau^i)$  the *realization weight* (RW) of sequence  $a_0^i, z_1^i, a_1^i, \dots, a_\tau^i$  under strategy  $\beta_0^i$ , defined as

$$rw^i(a_0^i, z_1^i, a_1^i, \dots, a_\tau^i) \stackrel{\text{def}}{=} \prod_{t=0}^{\tau} \beta_{0:}^i(a_t^i | a_0^i, z_1^i, a_1^i, \dots, z_t^i) \quad (19)$$

$$= rw^i(a_0^i, z_1^i, a_1^i, \dots, a_{\tau-1}^i) \cdot \beta_{0:}^i(a_\tau^i | \underbrace{a_0^i, z_1^i, a_1^i, \dots, z_\tau^i}_{\theta_\tau^i}). \quad (20)$$

This definition already leads to useful results such as:

$$\beta_{0:}^i(a_\tau^i | \theta_\tau^i) = \frac{rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i, z_\tau^i, a_\tau^i)}{rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i)}, \quad (21)$$

and

$$\forall z_\tau^i, \quad rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i) = rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i) \cdot \underbrace{\sum_{a_\tau^i} \beta(a_\tau^i | \theta_{\tau-1}^i, a_{\tau-1}^i, z_\tau^i)}_{=1} \quad (22)$$

$$= \sum_{a_\tau^i} rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i) \cdot \beta(a_\tau^i | \theta_{\tau-1}^i, a_{\tau-1}^i, z_\tau^i) \quad (23)$$

$$= \sum_{a_\tau^i} rw^i(\theta_{\tau-1}^i, a_{\tau-1}^i, z_\tau^i, a_\tau^i). \quad (24)$$

We now extend Koller et al.'s definition by introducing *conditional realization weights*, where the realization weight of a *suffix sequence* is “conditioned” on a *prefix sequence*:

$$rw^i(\underbrace{a_\tau^i, \dots, a_{\tau'}^i}_{\text{suffix seq.}} | \underbrace{a_0^i, \dots, z_\tau^i}_{\text{prefix seq.}}) \stackrel{\text{def}}{=} \prod_{t=\tau}^{\tau'} \beta_{0:}^i(a_t^i | a_0^i, \dots, z_\tau^i, a_\tau^i, \dots, z_t^i) \quad (25)$$

$$= \beta_0^i(a_\tau^i | a_0^i, \dots, z_\tau^i) \cdot rw^i(a_{\tau+1}^i, \dots, a_{\tau'}^i | a_0^i, \dots, z_{\tau+1}^i). \quad (26)$$

As can be noted, this definition only requires the knowledge of a partial strategy  $\beta_{\tau'}^i$  rather than a complete strategy  $\beta_0^i$ .

**Mixing Realization Weights** Let  $\tau' \geq \tau + 1$ , and  $rw^i[w]$  denote the realization weights of some element  $w$  at  $\tau + 1$ . Then, for some  $\psi_\tau^i$ , we have

$$rw[\psi_\tau^i](a_{\tau+1}^i, \dots, a_{\tau'}^i | a_0^i, \dots, z_{\tau+1}^i) \quad (27)$$

$$= \sum_w \psi_\tau^i(w) \cdot rw[w](a_{\tau+1}^i, \dots, a_{\tau'}^i | a_0^i, \dots, z_{\tau+1}^i). \quad (28)$$

---

**Algorithm 2:** Extracting  $\beta_0^i$  from  $w_0^i$

---

```

1 Fct Extract( $w_0^i$ )
   /* Step 1., keeping only  $rw(\theta_{0:H-1}^i)$  for all  $\theta_{0:H-1}^i$  */
2    $(rw(\theta_{0:H-1}^i))_{\theta_{0:H-1}^i} \leftarrow \mathbf{RecGetRWMix}(0, w_0^i)$ 
   /* Step 2. */
3   for  $t = H - 2, \dots, 0$  do
4     forall  $\theta_{0:t}^i, a_t^i$  do
5        $z_{t+1}^i \leftarrow z^i$  s.t.  $\beta_t(\cdot | \theta_{0:t}^i, a_t^i, z^i)$  is defined
6        $rw(\theta_{0:t}^i, a_t^i) \leftarrow \sum_{a_{t+1}^i} rw(\theta_{0:t}^i, a_t^i, z_{t+1}^i, a_{t+1}^i | -)$ 
   /* Step 3. */
7   for  $t = H - 1, \dots, 0$  do
8     forall  $\theta_{0:t}^i, a_t^i$  do
9        $\beta_t^i(a_t^i | \theta_{0:t}^i) \leftarrow \frac{rw^i(\theta_{0:t-1}^i, a_{t-1}^i, z_t^i, a_t^i)}{rw^i(\theta_{0:t-1}^i, a_{t-1}^i)}$ 
10  return  $\beta_0^i$ 
11 Fct RecGetRWMix( $t, w = \langle \beta_t^i, \psi_t^i \rangle$ )
12  for  $w'$  s.t.  $\psi_t^i(w') > 0$  do
13     $rwCat[w'] \leftarrow \mathbf{RecGetRWCat}(t, w')$ 
14  forall  $(a_0^i, \dots, a_{H-1}^i)$  do
15     $rwMix[w](a_t^i, \dots, a_{H-1}^i | a_0^i, \dots, z_t^i)$ 
16     $\leftarrow \sum_{w'} \psi_t^i(w') \cdot rwCat[w'](a_{t+1}^i, \dots, a_{H-1}^i | a_0^i, \dots, z_{t+1}^i)$ 
17  return  $rwMix[w]$ 
18 Fct RecGetRWCat( $t, w = \langle \beta_t^i, \psi_t^i \rangle$ )
19  if  $t = H - 1$  then
20    forall  $(a_0^i, \dots, a_{H-1}^i)$  do
21       $rwCat[w](a_{H-1}^i | a_0^i, \dots, z_{H-1}^i) \leftarrow \beta_t^i(a_{H-1}^i | a_0^i, \dots, z_{H-1}^i)$ 
22  else
23     $rwMix[w] \leftarrow \mathbf{RecGetRWMix}(t, w)$ 
24    forall  $(a_0^i, \dots, a_{H-1}^i)$  do
25       $rwCat[w](a_t^i, \dots, a_{H-1}^i | a_0^i, \dots, z_t^i) \leftarrow$ 
26       $\beta_t^i(a_t^i | a_0^i, \dots, z_t^i) \cdot rwMix[w](a_{t+1}^i, \dots, a_{H-1}^i | a_0^i, \dots, z_{t+1}^i)$ 
26  return  $rwCat[w]$ 

```

---

**From  $w_0^i$  to  $\beta_0^i$ .** Using the above results, function **Extract** in Algorithm 2 derives a behavioral strategy  $\beta_0^i$  equivalent to the recursive strategy induced by some tuple  $w_0^i$  in 3 steps as follows:

1. **From  $w_0^i$  to  $rw(\theta_{0:H-1}^i, a_{H-1}^i)$  ( $\forall (\theta_{0:H-1}^i, a_{H-1}^i)$ )** — These (classical) realization weights are obtained by recursively going through the directed acyclic graph describing the recursive strategy, computing *full length* (conditional) realization weights  $rw(\theta_{0:t}^i, a_{H-1}^i | \theta_{0:t}^i)$  (for  $t = H - 1$  down to 0).

When in a leaf node, at depth  $H - 1$ , the initialization is given by Equation (25) when  $\tau = \tau' = H - 1$ :

$$rw^i(a_{H-1}^i | a_0^i, \dots, z_{H-1}^i) \stackrel{\text{def}}{=} \prod_{t=H-1}^{H-1} \beta^i(a_t^i | a_0^i, \dots, z_t^i)$$

$$= \beta^i(a_{H-1}^i | a_0^i, \dots, z_{H-1}^i).$$

Then, in the backward phase, we can compute full length realization weights  $rw(\theta_{t+1:H-1}^i, a_{H-1}^i | \theta_{0:t}^i)$  with increasingly longer suffixes (thus shorter prefixes) using (i) Equation (28) (in function **RecGetRWMix**, line 25) to “mix” several strategies using the distribution  $\psi_t^i$  attached to the current  $w$ , and (ii) Equation (26), with  $\tau' = H - 1$ , (in function **RecGetRWCat**, line 16) to concatenate the behavioral decision rule  $\beta_t^i$  attached to the current  $w$  in front of the strategy induced by the distribution  $\psi_t^i$  also attached to  $w$ . Note: Memoization can here be used to avoid repeating the same computations.

2. **Retrieving (classical) realization weights**  $rw(\theta_{0:t}^i, a_t^i | -)$  ( $\forall t$ ) — We can now compute realization weights  $rw(\theta_{0:t}^i, a_t^i | -)$  for all  $t$ 's using Equation (24) (line 6).
3. **Retrieving behavioral decision rules**  $\beta_t^i$  — Applying Equation (21) (line 9) then provides the expected behavioral decision rules.

In practice, lossless compressions are used to reduce the dimensionality of the occupancy state (*cf.* Section 4.1), which are currently lost in the current implementation of the conversion. Ideally, one would like to preserve compressions whenever possible or at least retrieve them afterwards, and possibly identify further compressions in the solution strategy.

## E HSVI for zs-POSGs

This section presents results that help (i) tune zs-oMG-HSVI's radius parameter  $\rho$ , ensuring that trajectories will always stop, and (ii) then demonstrate the finite time convergence of this algorithm.

### E.1 Algorithm

#### E.1.1 Setting $\rho$

**Proposition E.1.** *Bounding  $\lambda_\tau$  by  $\lambda^\infty = \frac{1}{2} \frac{1}{1-\gamma} [r_{\max} - r_{\min}]$  when  $\gamma < 1$ , and noting that*

$$\begin{aligned} thr(\tau) &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \frac{\gamma^{-\tau} - 1}{1 - \gamma} \quad \text{if } \gamma < 1 \\ &= \epsilon - \rho(r_{\max} - r_{\min})(2H + 1 - \tau)\tau \quad \text{if } \gamma = 1, \end{aligned} \quad (29)$$

*one can ensure positivity of the threshold at any  $\tau \in 1..H - 1$  by enforcing  $0 < \rho < \frac{1-\gamma}{2\lambda^\infty} \epsilon$  (or  $0 < \rho < \frac{\epsilon}{(r_{\max} - r_{\min})(H+1)H}$  if  $\gamma = 1$ ).*

*Proof.* Let us first consider the case  $\gamma < 1$ .

We have (for  $\tau \in \{1..H - 1\}$ ):

$$\begin{aligned} thr(\tau) &= \gamma^{-\tau} \epsilon - \sum_{i=1}^{\tau} 2\rho \lambda^\infty \gamma^{-i} \\ &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \sum_{i=1}^{\tau} \gamma^{-i} \\ &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty (\gamma^{-1} + \gamma^{-2} + \dots + \gamma^{-\tau}) \\ &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \gamma^{-1} (\gamma^0 + \gamma^{-1} + \dots + \gamma^{-(\tau-1)}) \\ &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \gamma^{-1} \frac{\gamma^{-\tau} - 1}{\gamma^{-1} - 1} \\ &= \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \frac{\gamma^{-\tau} - 1}{1 - \gamma}. \end{aligned}$$

Then, let us derive the following equivalent inequalities:

$$\begin{aligned} 0 &< thr(\tau) \\ 2\rho \lambda^\infty \frac{\gamma^{-\tau} - 1}{1 - \gamma} &< \gamma^{-\tau} \epsilon \\ \rho &< \frac{1}{2\lambda^\infty} \frac{1 - \gamma}{\gamma^{-\tau} - 1} \gamma^{-\tau} \epsilon \end{aligned}$$

$$\rho < \frac{1}{2\lambda^\infty} \frac{1-\gamma}{1-\gamma^\tau} \epsilon.$$

To ensure positivity of the threshold for any  $\tau \geq 1$ , one thus just needs to set  $\rho$  as a positive value smaller than  $\frac{1-\gamma}{2\lambda^\infty} \epsilon$ .

Let us now consider the case  $\gamma = 1$ .

We have (for  $\tau \in \{1, \dots, H-1\}$ ):

$$\begin{aligned} thr(\tau) &\stackrel{\text{def}}{=} \epsilon - \sum_{i=1}^{\tau} 2\rho\lambda_{\tau-i} \\ &= \epsilon - \sum_{i=1}^{\tau} 2\rho(H - (\tau - i)) \cdot (r_{\max} - r_{\min}) \\ &= \epsilon - 2\rho(r_{\max} - r_{\min}) \left[ \tau(H - \tau) + \sum_{i=1}^{\tau} i \right] \\ &= \epsilon - 2\rho(r_{\max} - r_{\min}) \left[ \tau H - \tau^2 + \frac{1}{2}\tau(\tau + 1) \right] \\ &= \epsilon - 2\rho(r_{\max} - r_{\min}) \left[ \left(H + \frac{1}{2}\right)\tau - \frac{1}{2}\tau^2 \right] \\ &= \epsilon - \rho(r_{\max} - r_{\min}) [(2H + 1)\tau - \tau^2] \\ &= \epsilon - \rho(r_{\max} - r_{\min}) [(2H + 1 - \tau)\tau]. \end{aligned}$$

Then, let us derive the following equivalent inequalities:

$$\begin{aligned} 0 &< thr(\tau) \\ \rho(r_{\max} - r_{\min})(2H + 1 - \tau)\tau &< \epsilon \quad (\text{holds when } \tau = 0 \text{ and } \tau = H + 1) \\ \rho &< \frac{\epsilon}{(r_{\max} - r_{\min})(2H + 1 - \tau)\tau} \quad (\text{when } \tau \in \{0 \dots H + 1\}). \end{aligned}$$

The function  $f : \tau \mapsto \frac{\epsilon}{(r_{\max} - r_{\min})(2H + 1 - \tau)\tau}$  reaches its minimum (for  $\tau \in (0, H + 1)$ ) when  $\tau = H + \frac{1}{2}$ . To ensure positivity of the threshold for any  $\tau \in \{1 \dots H - 1\}$ , one thus just needs to set  $\rho$  as a positive value smaller than  $\frac{\epsilon}{(r_{\max} - r_{\min})(H + 1)H}$ .  $\square$

## E.2 Finite-Time Convergence

### E.2.1 Convergence Proof

Proving the finite-time convergence of zs-oMG-HSVI to an error-bounded solution requires some preliminary lemmas.

**Lemma 10.** *Let  $(\sigma_0, \dots, \sigma_{\tau+1})$  be a full trajectory generated by zs-oMG-HSVI and  $\beta_\tau$  the behavioral DR profile that induced the last transition, i.e.,  $\sigma_{\tau+1} = T(\sigma_\tau, \beta_\tau)$ . Then, after updating  $\overline{W}_\tau$  and  $\underline{W}_\tau$ , we have that  $\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) - \underline{W}_\tau(\sigma_\tau, \beta_\tau^2) \leq \gamma thr(\tau + 1)$ .*

*Proof.* By definition,

$$\begin{aligned} \overline{W}_\tau(\sigma_\tau, \beta_\tau^1) &= \min_{\substack{\langle \tilde{\sigma}_\tau^{c,1}, \tilde{\beta}_\tau^2, \langle \tilde{\nu}_{\tau+1}^2, \psi_{\tau+1}^2 \rangle \rangle \\ \in \overline{\mathcal{I}}_\tau^1}} \beta_\tau^1 \cdot \left( r(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \tilde{\beta}_\tau^2) \right. \\ &\quad \left. \cdot \left[ \tilde{\nu}_{\tau+1}^2 + \lambda_{\tau+1} \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2) - T_c^1(\tilde{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1} \right] \right). \end{aligned}$$

Therefore, after the update ( $\beta_\tau^2$  and  $\beta_\tau^1$  being added to their respective bags  $\overline{\mathcal{I}}_\tau^1$  and  $\underline{\mathcal{I}}_\tau^2$ ) along with vectors  $\tilde{\nu}_{\tau+1}^2$  and  $\underline{\nu}_{\tau+1}^1$ ,

$$\begin{aligned} \overline{W}_\tau(\sigma_\tau, \beta_\tau^1) &\leq \beta_\tau^1 \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \tilde{\nu}_{\tau+1}^2 \right], \text{ and} \\ \underline{W}_\tau(\sigma_\tau, \beta_\tau^2) &\geq \beta_\tau^2 \cdot \left[ r(\sigma_\tau, \beta_\tau^1, \cdot) + \gamma T_m^2(\sigma_\tau, \beta_\tau^1, \cdot) \cdot \underline{\nu}_{\tau+1}^1 \right]. \end{aligned}$$

Then,

$$\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) - \underline{W}_\tau(\sigma_\tau, \beta_\tau^2) \leq \left[ r(\sigma_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma T_m^1(\sigma_\tau, \beta_\tau) \cdot \tilde{\nu}_{\tau+1}^2 \right]$$

$$\begin{aligned}
 & - \left[ r(\underline{\sigma}_\tau, \beta_\tau^1, \beta_\tau^2) + \gamma T_m^2(\sigma_\tau, \beta_\tau) \cdot \underline{\nu}_{\tau+1}^1 \right] \\
 & = \gamma \left[ \overline{V}(T(\sigma_\tau, \beta_\tau)) - \underline{V}(T(\sigma_\tau, \beta_\tau)) \right] \\
 & \leq \gamma thr(\tau + 1) \quad (\text{Holds at the end of any trajectory.})
 \end{aligned}$$

□

**Lemma 11** (Monotonic evolution of  $\overline{W}_\tau$  and  $W_\tau$ ). *Let  $K\overline{W}_\tau$  and  $KW_\tau$  be the approximations after an update at  $\sigma_\tau$  with behavioral DR  $\langle \beta_\tau^1, \beta_\tau^2 \rangle$  (respectively associated to vectors  $\overline{\nu}_{\tau+1}^2$  and  $\underline{\nu}_{\tau+1}^1$ ). Let also  $K^{(n+1)}\overline{W}_\tau$  and  $K^{(n+1)}W_\tau$  be the same approximations after  $n$  other updates (in various OSS). Then,*

$$\begin{aligned}
 \max_{\beta_\tau^1} K^{(n+1)}\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) & \leq \max_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) \leq \overline{W}_\tau(\sigma_\tau, \overline{\beta}_\tau^1) \quad \text{and} \\
 \min_{\beta_\tau^2} K^{(n+1)}W_\tau(\sigma_\tau, \beta_\tau^2) & \geq \min_{\beta_\tau^2} KW_\tau(\sigma_\tau, \beta_\tau^2) \geq W_\tau(\sigma_\tau, \underline{\beta}_\tau^2).
 \end{aligned}$$

*Proof.* Starting from the definition,

$$\begin{aligned}
 & \max_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) \\
 & = \max_{\substack{\beta_\tau^1 \in \langle \overline{\sigma}_\tau^{c,1}, \beta_\tau^2, \langle \overline{\nu}_{\tau+1}^2, \psi_{\tau+1}^2 \rangle \rangle \in \\ \overline{\mathcal{I}}_\tau^1 \cup \{ \langle \overline{\sigma}_\tau^{c,1}, \beta_\tau^2, \langle \overline{\nu}_{\tau+1}^2, \psi_{\tau+1}^2 \rangle \rangle \}}} \min_{\langle \overline{\sigma}_\tau^{c,1}, \beta_\tau^2, \langle \overline{\nu}_{\tau+1}^2, \psi_{\tau+1}^2 \rangle \rangle \in \overline{\mathcal{I}}_\tau^1} \beta_\tau^1 \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) \right. \\
 & \quad \left. + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \left( \overline{\nu}_{\tau+1}^2 + \lambda_{\tau+1} \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2) - T_c^1(\overline{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1} \right) \right] \\
 & \leq \max_{\beta_\tau^1} \min_{\langle \overline{\sigma}_\tau^{c,1}, \beta_\tau^2, \langle \overline{\nu}_{\tau+1}^2, \psi_{\tau+1}^2 \rangle \rangle \in \overline{\mathcal{I}}_\tau^1} \beta_\tau^1 \cdot \left[ r(\sigma_\tau, \cdot, \beta_\tau^2) \right. \\
 & \quad \left. + \gamma T_m^1(\sigma_\tau, \cdot, \beta_\tau^2) \cdot \left( \overline{\nu}_{\tau+1}^2 + \lambda_{\tau+1} \overrightarrow{\|T_c^1(\sigma_\tau^{c,1}, \beta_\tau^2) - T_c^1(\overline{\sigma}_\tau^{c,1}, \beta_\tau^2)\|_1} \right) \right] \\
 & = \max_{\beta_\tau^1} \overline{W}_\tau(\sigma_\tau, \beta_\tau^1) \\
 & = \overline{W}_\tau(\sigma_\tau, \overline{\beta}_\tau^1).
 \end{aligned}$$

Then, this upper bound approximation can only be refined, so that, for any  $n \in \mathbf{N}$ ,

$$\begin{aligned}
 \forall \beta_\tau^1, \quad K^{(n+1)}\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) & \leq K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1), \\
 \text{thus, } \min_{\beta_\tau^1} K^{(n+1)}\overline{W}_\tau(\sigma_\tau, \beta_\tau^1) & \leq \min_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1).
 \end{aligned}$$

The expected result thus holds for  $\overline{W}_\tau$ , and symmetrically for  $W_\tau$ . □

**Lemma 12.** *After updating, in order,  $\overline{W}_\tau$  and  $\overline{V}_\tau$ , we have*

$$K\overline{V}_\tau(\sigma_\tau) \leq \max_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1).$$

*After updating, in order,  $W_\tau$  and  $\underline{V}_\tau$ , we have*

$$K\underline{V}_\tau(\sigma_\tau) \geq \min_{\beta_\tau^2} KW_\tau(\sigma_\tau, \beta_\tau^2).$$

*Proof.* After updating  $\overline{\mathcal{I}}_\tau^1$ , the algorithm computes (Algorithm 1, line 20) a new solution  $\overline{\psi}_\tau^2$  of the dual LP (at  $\sigma_\tau^1$ ) and the associated vector  $\overline{\nu}_\tau^2$ , so that

$$\max_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau^1, \beta_\tau^1) = \sigma_\tau^{m,1} \cdot \overline{\nu}_\tau^2.$$

This vector will feed  $\overline{bagV}_\tau$  along with  $\sigma_\tau^1$ , so that

$$K\overline{V}_\tau(\sigma_\tau) \leq \sigma_\tau^{m,1} \cdot \overline{\nu}_\tau^2.$$

As a consequence,

$$K\overline{V}_\tau(\sigma_\tau) \leq \max_{\beta_\tau^1} K\overline{W}_\tau(\sigma_\tau, \beta_\tau^1).$$

The symmetric property holds for  $K\underline{V}_\tau$  and  $K\underline{W}_\tau$ , which concludes the proof. □

**Theorem 3.9.** (originally stated on page 11) *zs-oMG-HSVI (Alg. 1) terminates in finite time with an  $\epsilon$ -approximation of  $V_0^*(\sigma_0)$  that satisfies Theorem 3.8.*

*Proof.* We will prove by induction from  $\tau = H$  to 0, that the algorithm stops expanding OSS at depth  $\tau$  after finitely many iterations (/trajectories).

First, by definition of horizon  $H$ , no OS  $\sigma_H$  is ever expanded. The property thus holds at  $\tau = H$ .

Let us now assume that the property holds at depth  $\tau + 1$  after  $N_{\tau+1}$  iterations. By contradiction, let us assume that the algorithm generates infinitely many trajectories of length  $\tau + 1$ . Then, because  $\mathcal{O}_\tau^\sigma \times \mathcal{B}_\tau$  is compact, after some time the algorithm will have visited  $\langle \sigma_\tau, \beta_\tau \rangle$ , then, some iterations later,  $\langle \sigma'_\tau, \beta'_\tau \rangle$ , such that  $\|\sigma_\tau - \sigma'_\tau\|_1 \leq \rho$ . Let us also note the corresponding terminal OSS (because trajectories beyond iteration  $N_{\tau+1}$  do not go further)  $\sigma_{\tau+1} = T(\sigma_\tau, \beta_\tau)$  and  $\sigma'_{\tau+1} = T(\sigma'_\tau, \beta'_\tau)$ .

Now, we show that the second trajectory should not have happened, i.e.,  $\bar{V}(\sigma'_\tau) - \underline{V}(\sigma'_\tau) \leq thr(\tau)$ .

Combining the previous lemmas,

$$\begin{aligned} \bar{V}(\sigma'_\tau) &\leq \bar{V}(\sigma_\tau) + \lambda_\tau \|\sigma_\tau - \sigma'_\tau\|_1 && \text{(By Lipschitz-Continuity)} \\ &\leq \max_{\beta_\tau^1} \bar{W}_\tau(\sigma_\tau, \beta_\tau^1) + \lambda_\tau \|\sigma_\tau - \sigma'_\tau\|_1 && \text{(Lemma 12)} \\ &\leq \bar{W}_\tau(\sigma_\tau, \beta_\tau^1) + \lambda_\tau \|\sigma_\tau - \sigma'_\tau\|_1 && \text{(Lemma 11)} \\ &= \bar{W}_\tau(\sigma_\tau, \beta_\tau^1) + \lambda_\tau \rho. \end{aligned}$$

Symmetrically, we also have

$$\underline{V}(\sigma'_\tau) \geq \underline{W}_\tau(\sigma_\tau, \beta_\tau^2) - \lambda_\tau \rho.$$

Hence,

$$\begin{aligned} \bar{V}(\sigma'_\tau) - \underline{V}(\sigma'_\tau) &\leq (\bar{W}_\tau(\sigma_\tau, \beta_\tau^1) + \lambda_\tau \rho) - (\underline{W}_\tau(\sigma_\tau, \beta_\tau^2) - \lambda_\tau \rho) \\ &= (\bar{W}_\tau(\sigma_\tau, \beta_\tau^1) - \underline{W}_\tau(\sigma_\tau, \beta_\tau^2)) + 2\lambda_\tau \rho \\ &\leq \gamma thr(\tau + 1) + 2\lambda_\tau \rho && \text{(Lemma 10)} \\ &= \gamma \left( \gamma^{-(\tau+1)} \epsilon - \sum_{i=1}^{\tau+1} 2\rho \lambda_{\tau+1-i} \gamma^{-i} \right) + 2\lambda_\tau \rho \\ &= \gamma^{-\tau} \epsilon - \sum_{i=1}^{\tau+1} 2\rho \lambda_{\tau+1-i} \gamma^{-i+1} + 2\lambda_\tau \rho \\ &= \gamma^{-\tau} \epsilon - \sum_{j=0}^{\tau} 2\rho \lambda_{\tau-j} \gamma^{-j} + 2\lambda_\tau \rho \\ &= \gamma^{-\tau} \epsilon - \cancel{2\rho \lambda_{\tau-0} \gamma^{-0}} - \sum_{j=1}^{\tau} 2\rho \lambda_{\tau-j} \gamma^{-j} + \cancel{2\lambda_\tau \rho} = thr(\tau). \end{aligned}$$

Therefore,  $\sigma'_\tau$  should not have been expanded. This shows that the algorithm will generate only a finite number of trajectories of length  $\tau$ .  $\square$

## E.2.2 Handling Infinite Horizons

**Proposition 3.10.** (originally stated on page 12) *When  $\gamma < 1$ , the length of trajectories is upper bounded by  $T_{\max} \stackrel{\text{def}}{=} \left\lceil \log_\gamma \frac{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}}{W - \frac{2\rho\lambda^\infty}{1-\gamma}} \right\rceil$ , where  $\lambda^\infty$  is a depth-independent Lipschitz constant and  $W \stackrel{\text{def}}{=} \|\bar{V}^{(0)} - \underline{V}^{(0)}\|_\infty$  is the maximum width between initializations.*

*Proof.* (detailed version) Since  $W$  is the largest possible width, any trajectory stops in the worst case at depth  $\tau$  such that

$$\begin{aligned} thr(\tau) &< W \\ \gamma^{-\tau} \epsilon - 2\rho \lambda^\infty \frac{\gamma^{-\tau} - 1}{1 - \gamma} &< W \end{aligned} \quad \text{(from Eq. (29))}$$

$$\begin{aligned}
 \gamma^{-\tau} \epsilon - 2\rho\lambda^\infty \frac{\gamma^{-\tau}}{1-\gamma} - 2\rho\lambda^\infty \frac{-1}{1-\gamma} &< W \\
 \underbrace{\gamma^{-\tau} \left( \epsilon - \frac{2\rho\lambda^\infty}{1-\gamma} \right)}_{>0 \quad (\text{Prop. E.1})} &< W - \frac{2\rho\lambda^\infty}{1-\gamma} \\
 \gamma^{-\tau} &< \frac{W - \frac{2\rho\lambda^\infty}{1-\gamma}}{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}} \\
 \exp(-\tau \ln(\gamma)) &< \exp\left(\ln\left(\frac{W - \frac{2\rho\lambda^\infty}{1-\gamma}}{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}}\right)\right) \\
 -\tau \ln(\gamma) &< \ln\left(\frac{W - \frac{2\rho\lambda^\infty}{1-\gamma}}{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}}\right) \\
 \tau \ln(\gamma) &> \ln\left(\frac{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}}{W - \frac{2\rho\lambda^\infty}{1-\gamma}}\right) \\
 \tau &< \log_\gamma\left(\frac{\epsilon - \frac{2\rho\lambda^\infty}{1-\gamma}}{W - \frac{2\rho\lambda^\infty}{1-\gamma}}\right).
 \end{aligned}$$

□

Even if the problem horizon is infinite, trajectories will thus have bounded length. Then, everything beyond this *effective* horizon will rely on the upper- and lower-bound initializations and the corresponding strategies.