# STRATEGY SYNTHESIS FOR ZERO-SUM NEURO-SYMBOLIC CONCURRENT STOCHASTIC GAMES[*]

RUI YAN[†], GABRIEL SANTOS[†], GETHIN NORMAN[‡], DAVID PARKER[§], AND MARTA KWIATKOWSKA[†]

**Abstract.** Neuro-symbolic approaches to artificial intelligence, which combine neural networks with classical symbolic techniques, are growing in prominence, necessitating formal approaches to reason about their correctness. We propose a novel modelling formalism called neuro-symbolic concurrent stochastic games (NS-CSGs), which comprise probabilistic finite-state agents interacting in a shared continuous-state environment, observed through perception mechanisms implemented as neural networks. Since the environment state space is continuous, we focus on the class of NS-CSGs with Borel state spaces. We consider the problem of zero-sum discounted cumulative rewards and prove the existence of the value of NS-CSGs under Borel measurability and piecewise-constant restrictions on the components of the model. From an algorithmic perspective, existing methods to compute values and optimal strategies for CSGs focus on finite state spaces. We present, for the first time, implementable value iteration and policy iteration algorithms to solve a class of uncountable state space CSGs, namely NS-CSGs, and prove their convergence. Our approach works by exploiting the underlying game structures and then formulating piecewise linear or constant representations of the value functions and strategies of NS-CSGs. We illustrate our approach by applying a prototype implementation of value iteration to a dynamic vehicle parking case study.

**Key words.** Stochastic games, neuro-symbolic systems, value iteration, policy iteration, Borel state spaces

**MSC codes.** 91A15, 92B20, 60J05, 93E20

**1. Introduction.** Game theory offers an attractive framework for analysing strategic interactions among agents in machine learning, with application to, for instance, the game of Go [30], autonomous driving [28] and robotics [10]. An important class of dynamic games is *stochastic games* [29], which move between states according to transition probabilities controlled jointly by multiple agents (players). Extending both strategic-form games to dynamic environments and Markov decision processes to multiple players, stochastic games have long been used to model sequential decision-making problems with more than one agent, ranging from multi-agent reinforcement learning [34], to quantitative verification and synthesis for equilibria [16].

Recent years have witnessed encouraging advances in the use of neural networks (NNs) to approximate either value functions or strategies [17] for stochastic games that model large, complex environments. Such *end-to-end* NNs directly map environment states to Q-values or actions. This means that they have a relatively complex structure and a large number of weights and biases, since they interweave multiple tasks (e.g., object detection and recognition, decision making) within a single NN. An emerging trend in autonomous and robotic systems is *neuro-symbolic* approaches, where some components that are synthesized from data (e.g., perception modules) are implemented as NNs, while others (e.g., nonlinear controllers) are formulated us-

[†]Department of Computer Science, University of Oxford, Oxford, OX1 2JD, UK (rui.yan@cs.ox.ac.uk, gabriel.santos@cs.ox.ac.uk, marta.kwiatkowska@cs.ox.ac.uk).

[‡]School of Computing Science, University of Glasgow, University Avenue, Glasgow, G12 8QQ, UK (gethin.norman@glasgow.ac.uk)

[§]School of Computer Science, University of Birmingham, Edgbaston, Birmingham, B15 2TT, UK (d.a.parker@cs.bham.ac.uk)

ing traditional symbolic methods. This can greatly simplify the design and training process, and yield smaller NNs.

However, there remains a lack of modelling and verification frameworks which can reason formally about the correctness of neuro-symbolic systems. Progress has been made on techniques for both multi-agent verification [1, 24] and safe reinforcement learning [2] in this context, but without the ability to reason formally about stochasticity, which is crucial for modelling uncertainty. Elsewhere, concurrent stochastic games (CSGs) have been widely studied [32, 31, 9, 22], but primarily in the context of finite state spaces, which are insufficient for many real-life systems; for instance, uncountable real vector spaces are usually supplied as inputs to NNs.

Stochastic games were introduced by Shapley [29], who assumed a finite state space. Since then, many researchers, have considered CSGs with uncountable state spaces, e.g., [15, 18, 20]. Maitra and Parthasarathy [18] were the first to study discounted zero-sum CSGs in this setting, assuming that the state space is a compact metric space. Following this, more general results for discounted zero-sum CSGs with Borel state spaces have been derived, e.g., [15, 19, 20, 11]. These aim at providing sufficient conditions for the *existence* of either values or optimal strategies for players.

Another important and practical problem for zero-sum CSGs with uncountable state spaces is the *computation* of values and optimal strategies. Since the seminal policy iteration (PI) methods were introduced by Hoffman and Karp [12], and by Pollatschek and Avi-Itzhak [23], a wide range of fixed-point algorithms have been developed for zero-sum CSGs with finite state spaces [32, 31, 9, 22]. Recent work by Bertsekas [5] proposed a distributed optimistic abstract PI algorithm, which inherits the attractive structure of the Pollatschek and Avi-Itzhak algorithm while resolving the convergence difficulties suffered by the latter. All these assume finite state spaces. However, to the best of our knowledge, there are no existing value iteration (VI) or PI algorithms for CSGs with uncountable, or more specifically Borel, state spaces. VI and PI algorithms for stochastic control (i.e., the one player case) with Borel state spaces can be found in [37, 36]. Other problems for zero-sum CSGs with uncountable state spaces have been studied and include information structure [13], specialized strategy spaces [4], continuous time setup [8] and payoff criteria [11].

As discussed later, this paper assumes a *fully observable* game setting; a natural extension would be partially observable stochastic games (POSGs). A variant of POSGs, called factored-observation stochastic games (FOSGs), was recently proposed [14] that distinguishes between private and public observations in a similar fashion to our model but for finite-state models without NNs. Partial observability in FOSGs is dealt with via a mechanism that converts imperfect-information games into continuous-state (public belief state) perfect-information games [7, 14], such that many techniques for perfect-information games can also be applied. The fully-observable model can arguably serve as a vehicle to later solve the more complex case with imperfect information.

**Contributions.** First, we propose a new modelling formalism called *neuro-symbolic concurrent stochastic games (NS-CSGs)*, which comprise multiple finite-state agents endowed with perception mechanisms implemented via NNs and conventional, symbolic decision-making mechanisms. The agents proceed concurrently in a continuous, shared environment, since in many applications inputs supplied to NNs are real-valued vectors. Second, we tackle the problem of optimising zero-sum discounted cumulative rewards for NS-CSGs, under the assumption that agents have full state observability. Working with Borel state spaces, we establish restrictions on the mod-

elling formalism which ensure that the NS-CSGs belong to a class of uncountable state-space CSGs [15] that are determined, and therefore prove the existence of the value of NS-CSGs.

Third, we present two algorithms to compute values and synthesise optimal strategies for these games, one using value iteration and the other using policy iteration. Our approach is based on formulating piecewise constant representations of the value functions and strategies for NS-CSGs by exploiting the underlying game structure. To the best of our knowledge, these are the first implementable algorithms for solving zero-sum CSGs over Borel state spaces with a convergence guarantee. The policy iteration approach is inspired by recent work for finite state spaces [5], which we generalise by employing piecewise constant or piecewise linear functions to ensure finite representability. This allows us to overcome the main issue that arises when solving Borel state space CSGs with policy iteration, namely that the value function may change from a Borel measurable function to a non-Borel measurable function across iterations. Finally, we illustrate our approach by modelling a dynamic vehicle parking case study with NS-CSGs, and then synthesising strategies for it using a prototype implementation of value iteration.

**2. Background.** In this section we summarise the notation used in this paper.

**2.1. Borel measurable functions.** Given a non-empty set $X$, we denote its Borel $\sigma$-algebra by $\mathcal{B}(X)$, and the sets in $\mathcal{B}(X)$ are called *Borel sets* of $X$. The pair $(X, \mathcal{B}(X))$ is a (standard) *Borel space* if there exists a metric on $X$ that makes it a complete separable metric space (unless required for clarity, $\mathcal{B}(X)$ will be omitted). For convenience we will work with real vector spaces; however, this is not essential and any complete separable metric spaces could be used. For Borel spaces $X$ and $Y$, a function $f : X \to Y$ is *Borel measurable* if $f^{-1}(B) \in \mathcal{B}(X)$ for all $B \in \mathcal{B}(Y)$ and *bimeasurable* if it is Borel measurable and $f(B) \in \mathcal{B}(Y)$ for all $B \in \mathcal{B}(X)$.

We denote by $\mathbb{F}(X)$ the space of all bounded, Borel measurable real-valued functions on a Borel space $X$, with respect to the unweighted sup-norm $\|J\| = \sup_{x \in X} |J(x)|$ for $J \in \mathbb{F}(X)$. For functions $J, K \in \mathbb{F}(X)$, we use $\max[J, K]$ and $\min[J, K]$ to denote their respective pointwise maximum and minimum functions, i.e., for each $x \in X$, we have $\mathrm{opt}[J, K](x) := \mathrm{opt}\{J(x), K(x)\}$ for $\mathrm{opt} \in \{\min, \max\}$.

We now introduce notation and definitions for concepts that are fundamental to the abstraction on which our algorithms are performed. The abstraction is based on a decomposition of the uncountable state space into finitely many abstract regions such that all concrete states in a region have the same behaviour. In the definitions below, let $X \subset \mathbb{R}^{n_1}$ and $Y \subset \mathbb{R}^{n_2}$ for $n_1, n_2 > 0$.

DEFINITION 2.1 (FCP and Borel FCP). *A finite connected partition (FCP) of $X$, denoted $\Phi$, is a finite collection of disjoint connected subsets (regions) that cover $X$. Furthermore, $\Phi$ is a Borel FCP (BFCP) if each region $\phi \in \Phi$ is a Borel set of $X$.*

DEFINITION 2.2 (PWC Borel measurable function). *A function $f : X \to Y$ is piecewise constant (PWC) Borel measurable if there exists a BFCP $\Phi$ of $X$ such that $f : \phi \to Y$ is constant for all $\phi \in \Phi$ and is called a constant-BFCP of $X$ for $f$.*

DEFINITION 2.3 (PWL Borel measurable function). *A function $f : X \to Y$ is piecewise linear (PWL) Borel measurable if there exists a BFCP $\Phi$ of $X$ such that $f : \phi \to Y$ is bounded and linear for all $\phi \in \Phi$.*

DEFINITION 2.4 (BFCP invertible function). *A function $f : X \to Y$ is BFCP invertible if, for any BFCP $\Phi$ of $Y$, there exists a BFCP $\Phi'$ of $X$, called a consistent-*

*BFCP of $X$ for $f$ to $\Phi$, such that for each $\phi' \in \Phi'$ there exists $\phi \in \Phi$ such that* $\{f(x) \mid x \in \phi'\} \subseteq \phi$.

For BFCPs $\Phi_1$ and $\Phi_2$ of $X$, we denote by $\Phi_1 + \Phi_2$ the smallest BFCP of $X$ such that $\Phi_1 + \Phi_2$ is a refinement of both $\Phi_1$ and $\Phi_2$, which can be obtained by taking all the intersections between regions of $\Phi_1$ and regions of $\Phi_2$.

**2.2. Probability measures.** Let $X$ be a Borel space. A function $f : \mathcal{B}(X) \rightarrow [0,1]$ is a probability measure on $X$ if $f(X) = 1$ and $\sum_{i \in I} f(B_i) = f(\cup_{i \in I} B_i)$ for any countable disjoint family of Borel sets $(B_i)_{i \in I}$. We denote the space of all probability measures on a Borel space $X$ by $\mathbb{P}(X)$. For Borel spaces $X$ and $Y$, a Borel measurable function $\sigma : Y \rightarrow \mathbb{P}(X)$ is called a *stochastic kernel* on $X$ given $Y$ (also known as a transition probability function from $Y$ to $X$), and we denote by $\mathbb{P}(X \mid Y)$ the set of all stochastic kernels on $X$ given $Y$. If $\sigma \in \mathbb{P}(X \mid Y)$, $y \in Y$ and $B \in \mathcal{B}(X)$, then we write $\sigma(B \mid y)$ for $\sigma(y)(B)$. It follows that $\sigma \in \mathbb{P}(X \mid Y)$ if and only if $\sigma(\cdot \mid y) \in \mathbb{P}(X)$ for all $y \in Y$ and $\sigma(B \mid \cdot)$ is Borel measurable for all $B \in \mathcal{B}(X)$.

**2.3. Neural networks.** A *neural network (NN)* is a real vector-valued function $f : \mathbb{R}^m \rightarrow \mathbb{R}^c$, where $m, c \in \mathbb{N}$, composed of a sequence of *layers* $h_1, \ldots, h_k$, where $h_i : \mathbb{R}^{m_i} \rightarrow \mathbb{R}^{c_i}$ for $1 \leq i \leq k$, $m_1 = m$, $c_i = m_{i+1}$ for $1 \leq i \leq k-1$ and $c_k = c$. Each layer $h_i$ is a data-processing module explicitly formulated as $h_i(x_i) = act_i(W_i x_i + b_i)$, where $x_i$ is the input to the $i$th layer given by the output $h_{i-1}(x_{i-1})$ of the $(i-1)$th layer, $act_i$ is an activation function, and $W_i x_i + b_i$ is a weighted sum of $x_i$ for a weight matrix $W_i$ and a bias vector $b_i$. An NN $f$ is continuous for all popular activation functions, e.g., Rectified Linear Unit (ReLU), Sigmoid and Softmax. An NN $f$ is said to be a *classifier* for a set of classes $C$ of size $c$ if, for any input $x \in \mathbb{R}^m$, the output $f(x) \in \mathbb{R}^c$ is a probability vector where, for any $1 \leq i \leq c$, the $i$th element of $f(x)$ represents the probability of the $i$th class of $C$, i.e., a classifier is a function $f : \mathbb{R}^m \rightarrow \mathbb{P}(C)$.

**2.4. Concurrent stochastic games.** We recall the model of two-player *concurrent stochastic games*.

DEFINITION 2.5. *A (two-player) concurrent stochastic game (CSG) is a tuple* $\mathsf{G} = (N, S, A, \Delta, \delta)$ *where:*
- $N = \{1, 2\}$ *is the set of two players;*
- $S$ *is a finite set of states;*
- $A = (A_1 \cup \{\perp\}) \times (A_2 \cup \{\perp\})$ *where $A_i$ is a finite set of actions available to player $i \in N$ and $\perp$ is an idle action disjoint from the set $A_1 \cup A_2$;*
- $\Delta : S \rightarrow 2^{(A_1 \cup A_2)}$ *is an action assignment function;*
- $\delta : (S \times A) \rightarrow \mathbb{P}(S)$ *is a probabilistic transition function.*

In a state $s$ of a CSG $\mathsf{G}$, each player $i \in N$ selects an action from its available actions, i.e.. from the set $\Delta(s) \cap A_i$, if this set is non-empty, and selects the idle action $\perp$ otherwise. We denote the action choices for each player $i$ in state $s$ by $A_i(s)$, i.e. $A_i(s)$ equals $\Delta(s) \cap A_i$ if $\Delta(s) \cap A_i \neq \varnothing$ and equals $\{\perp\}$ otherwise and by $A(s)$ the possible joint actions in a state, i.e., $A(s) = A_1(s) \times A_2(s)$. Supposing each player $i$ chooses action $a_i$, then with probability $\delta(s, (a_1, a_2))(s')$ there is a transition to state $s' \in S$. A *path* $\pi$ of $\mathsf{G}$ is a sequence $\pi = s_0 \xrightarrow{\alpha_0} s_1 \xrightarrow{\alpha_1} \cdots$ such that $s_k \in S$, $\alpha_k \in A(s_k)$ and $\delta(s_k, \alpha_k)(s_{k+1}) > 0$ for all $k \geq 0$. For a path $\pi$, we denote by $\pi(k)$ the $(k+1)$th state, and $\pi[k]$ the action for the transition from $\pi(k)$ to $\pi(k+1)$.

**3. Zero-sum neuro-symbolic concurrent stochastic games.** In this section we introduce our model of *neuro-symbolic concurrent stochastic games* (NS-CSGs).

We restrict attention to two-agent (player) games as we are concerned with zero-sum games in which there are two agents with directly opposing objectives. However, the approach extends to multi-agent games, by allowing the agents to form two coalitions with directly opposing objectives. An NS-CSG comprises two interacting neuro-symbolic *agents* acting in a shared, continuous-state environment. Each agent has finitely many local states and actions, and is endowed with a perception mechanism implemented as an NN through which it can observe the state of the environment, storing the observations locally in *percepts*.

DEFINITION 3.1. *A (two-agent) neuro-symbolic concurrent stochastic game (NS-CSG)* C *comprises agents* $(\mathsf{Ag}_i)_{i \in N}$, *for* $N = \{1, 2\}$, *and an environment* $E$ *where:*

$$\mathsf{Ag}_i = (S_i, A_i, \Delta_i, obs_i, \delta_i) \text{ for } i \in N, \quad E = (S_E, \delta_E)$$

*and we have:*
- $S_i = Loc_i \times Per_i$ *is a set of states for* $\mathsf{Ag}_i$, *where* $Loc_i \subseteq \mathbb{R}^{b_i}$ *and* $Per_i \subseteq \mathbb{R}^{d_i}$ *are finite sets of local states and percepts, respectively;*
- $S_E \subseteq \mathbb{R}^e$ *is a closed infinite set of environment states;*
- $A_i$ *is a nonempty finite set of actions for* $\mathsf{Ag}_i$, *and* $A := (A_1 \cup \{\bot\}) \times (A_2 \cup \{\bot\})$ *is the set of joint actions, where* $\bot$ *is an idle action disjoint from* $A_1 \cup A_2$;
- $\Delta_i : S_i \to 2^{A_i}$ *is an available action function, defining the actions* $\mathsf{Ag}_i$ *can take in each state;*
- $obs_i : (Loc_1 \times Loc_2 \times S_E) \to Per_i$ *is an observation function for* $\mathsf{Ag}_i$, *mapping the local states of all agents and the environment to a percept of the agent, implemented via an NN classifier for the set* $Per_i$;
- $\delta_i : (S_i \times A) \to \mathbb{P}(Loc_i)$ *is a partial probabilistic transition function for* $\mathsf{Ag}_i$ *determining the distribution over the agent's local states given its current state and joint action;*
- $\delta_E : (S_E \times A) \to S_E$ *is a partial deterministic transition function for the environment determining its next state given its current state and joint action.*

In an NS-CSG C the agents and environment execute concurrently and agents move between their local states probabilistically. For simplicity, we consider deterministic environments, but all the results extend directly to probabilistic environments with finite branching.

A (global) state for an NS-CSG comprises a state $s_i = (loc_i, per_i)$ for each agent $\mathsf{Ag}_i$ (a pair of a local state and percept, finite sets of real vectors) and an environment state $s_E$ drawn from a closed set of real vectors. In state $s = (s_1, s_2, s_E)$, each $\mathsf{Ag}_i$ simultaneously chooses one of the actions available in its state $s_i$ (if no action is available, i.e., $\Delta_i(s_i) = \varnothing$, then $\mathsf{Ag}_i$ chooses the idle action $\bot$). This results in a joint action $\alpha = (a_1, a_2) \in A$. Next, each $\mathsf{Ag}_i$ updates its local state to $loc_i' \in Loc_i$, according to its probabilistic local transition function $\delta_i$, applied to the current agent state $s_i$ and the joint action $\alpha$. The environment updates its state to $s_E' \in S_E$ according to its deterministic transition function $\delta_E$, based on its state $s_E$ and the actions $\alpha$. Finally, each $\mathsf{Ag}_i$ observes the new state of the other agents and the environment using the observation function $obs_i$ to generate a new percept $per_i' = obs_i(loc_1', loc_2', s_E')$. Thus, the game reaches the state $s' = (s_1', s_2', s_E')$, where $s_i' = (loc_i', per_i')$ for $1 \le i \le 2$. A state $s = (s_1, s_2, s_E)$ with $s_i = (loc_i, per_i)$ is *percept compatible* if $per_i = obs_i(loc_1, loc_2, s_E)$ for $1 \le i \le 2$.

We assume that each observation function $obs_i$ is implemented via an NN $f_i : \mathbb{R}^m \to \mathbb{P}(Per_i)$, where $m = b_1 + b_2 + e$, yielding a belief (normalised score) over the percepts of the agent, to which a rule is then applied that selects the percept

with the maximum probability. However, observation functions can be defined by any function, including other types of machine learning models. For convenience, in our formalisation the percepts (and the local states) of agents are drawn from finite sets of real-valued vectors, but any finite set could be used instead.

**3.1. Semantics of an NS-CSG.** Formally, the semantics of an NS-CSG $\mathsf{C}$ is a concurrent stochastic game $[\![\mathsf{C}]\!]$ over the product of the states of the agents and the environment.

DEFINITION 3.2 (Semantics of an NS-CSG). *Given an NS-CSG $\mathsf{C}$ consisting of two agents and an environment, the semantics of $\mathsf{C}$ is the CSG $[\![\mathsf{C}]\!] = (N, S, A, \Delta, \delta)$ where:*

- $S \subseteq S_1 \times S_2 \times S_E$ *is the set of global (percept compatible) states;*
- $A = (A_1 \cup \{\bot\}) \times (A_2 \cup \{\bot\})$;
- $\Delta(s_1, s_2, s_E) = \Delta_1(s_i) \cup \Delta_2(s_2)$;
- $\delta : (S \times ((A_1 \cup \{\bot\}) \times (A_2 \cup \{\bot\}))) \to \mathbb{P}(S)$ *is the partial probabilistic transition function, where for states $s = (s_1, s_2, s_E), s' = (s_1', s_2', s_E') \in S$ and joint action $\alpha = (a_1, a_2) \in A$, if $a_i \in \Delta_i(s_i)$ when $\Delta_i(s_i) \neq \varnothing$ and $a_i = \bot$ otherwise for $1 \leq i \leq 2$, then $\delta(s, \alpha)$ is defined and, if $s_i' = (loc_i', per_i')$, $per_i' = obs_i(loc_1', loc_2', s_E')$ for $1 \leq i \leq 2$ and $s_E' = \delta_E(s_E, \alpha)$, then*

$$\delta(s, \alpha)(s') = \delta_1(s_1, \alpha)(loc_1')\delta_2(s_2, \alpha)(loc_2')$$

*and otherwise $\delta(s, \alpha)(s') = 0$.*

Notice that the CSG $[\![\mathsf{C}]\!]$ is over percept compatible states and that the underlying transition relation $\delta$ is closed with respect to percept compatible states, by definition of $\delta_i$ for each agent $\mathsf{Ag}_i$. Since $\delta_E$ is deterministic and $Loc_i$ is a finite set, the branching set $\Theta(s, \alpha) = \{s' \mid \delta(s, \alpha)(s') > 0\}$ is finite for all $s \in S$ and $\alpha \in A(s)$. While the semantics of an NS-CSG is an instance of the general class of uncountable state space CSGs, its particular structure (see Definition 3.1) will be important in order to establish measurability and define our algorithms.

**3.2. Strategies.** The semantics of an NS-CSG is a CSG where the players are the agents, and therefore the behaviour of an NS-CSG is controlled by *strategies* for the agents. Since the state space $S$ is uncountable due to the continuous environment state space, we follow the approach of [15] and require Borel measurable conditions on the choices that the strategies can make to ensure the measurability of the induced sets of paths.

The semantics of any NS-CSG will turn out to be an instance of the class of CSGs from [15], for which *stationary strategies* achieve optimal values [15, Theorem 2(ii), Theorem 3], and therefore, to simplify the presentation, we restrict our attention to stationary strategies and refer to them simply as strategies. Before we give the formal definition, since we work with real vector spaces we require the following lemma.

LEMMA 3.3 (Borel spaces). *The sets $S$, $S_i$, $S_E$ and $A_i$ for $1 \leq i \leq 2$ are Borel spaces.*

*Proof.* By Theorem 27 [27, Chapter 9.6], $S_i$ and $S_E$ are separable metric spaces. By Theorem 12 [27, Chapter 9.4], $S_i$ and $S_E$ are complete metric spaces. Thus, $S_i$ and $S_E$ are Borel spaces. According to Theorem 1.10 [21, Chapter 1], the state space $S_1 \times S_2 \times S_E$, their Cartesian product space, is a Borel space. Since $obs_i$ is Borel measurable (stated at the end of this section) for $1 \leq i \leq 2$, then for $(loc_i, per_i) \in S_i$

and $1 \le i \le 2$, the set

$$\{((loc_1, per_1), (loc_2, per_2), s_E) \in S \mid obs_i(loc_1, loc_2, s_E) = per_i \text{ for } 1 \le i \le 2\}$$

is a Borel set of $S_1 \times S_2 \times S_E$. Then, $S$ is a Borel space as $S_1$ and $S_2$ are finite. Similarly, $A_i$ are Borel spaces, as they are finite sets. □

DEFINITION 3.4 (Strategy). *A (stationary) strategy for* $\mathsf{Ag}_i$ *of an NS-CSG* $\mathsf{C}$ *is a stochastic kernel* $\sigma_i : S \to \mathbb{P}(A_i)$, *i.e.,* $\sigma_i \in \mathbb{P}(A_i \mid S)$, *such that* $\sigma_i(A_i(s) \mid s) = 1$ *for all* $s \in S$. *A (strategy) profile* $\sigma = (\sigma_1, \sigma_2)$ *is a tuple of strategies for each agent. We denote by* $\Sigma_i$ *the set of all strategies of* $\mathsf{Ag}_i$ *and by* $\Sigma = \Sigma_1 \times \Sigma_2$ *the set of profiles.*

For a state $s \in S$, we use $\mathbb{P}(A_i(s))$ to denote the largest subset of $\mathbb{P}(A_i)$ such that $u_i(A_i(s)) = 1$ for all probability measures $u_i \in \mathbb{P}(A_i(s))$.

**3.3. Observability.** NS-CSGs are designed to model neuro-symbolic agents, whose operation depends on particular perception functions, which may result in imperfect information. However, in this paper we consider *full observability*, i.e., where agents' decisions can depend on the full state space. It is straightforward to extend the semantics above to *partially observable* CSGs (POSGs) [25, 26] where, for any state, each agent's observation function returns the agent's percept component of the state, and by restricting to observationally-equivalent strategies. However, our initial results on NS-CSGs focus on the simpler (but still challenging) case of full observability for the following reasons. First, the fully observable case represents an important baseline for this new model class, against which partially observable scenarios can later be evaluated. Second, the solution presented here can be directly used when converting imperfect-information games to perfect information, such as the mechanism described in [7, 14].

**3.4. Zero-sum NS-CSGs.** We are concerned with *zero-sum* games, namely, two-agent games in which the agents have directly opposing objectives. The objectives we consider are *discounted accumulated rewards* and we assume the first agent tries to maximise the value of this objective and the second tries to minimise it. More precisely, for a reward structure of the form $r = (r_A, r_S)$, where $r_A : S \times A \to \mathbb{R}$ is an action reward function and $r_S : S \to \mathbb{R}$ is a state reward function, the accumulated discounted reward for a path $\pi$ over the infinite-horizon is defined by:

$$(3.1) \qquad Y(\pi) = \sum_{k=0}^{\infty} \beta^k \big(r_A(\pi(k), \pi[k]) + r_S(\pi(k))\big)$$

where $\beta \in (0, 1)$ is the discount factor.

**3.5. Assumptions.** We conclude our introduction to NS-CSGs by clarifying the assumptions that are required for the results presented in this paper.

ASSUMPTION 3.5. *For any NS-CSG* $\mathsf{C}$ *and reward structure* $r = (r_A, r_S)$:
(i) $\delta_E(\cdot, \alpha) : S_E \to S_E$ *is a bimeasurable, BFCP invertible function for all* $\alpha \in A$;
(ii) $obs_i(loc_1, loc_2, \cdot) : S_E \to Per_i$ *is a PWC Borel measurable function for all* $loc_i \in Loc_i$ *and* $1 \le i \le 2$;
(iii) $r_A(\cdot, \alpha), r_S : S \to \mathbb{R}$ *are bounded, PWC Borel measurable real-valued functions for all* $\alpha \in A$.

Our assumptions for NS-CSGs differ from existing stochastic games with Borel state spaces [15, 19, 11] in that the states have both discrete and continuous elements, while the observation and reward functions are required to be PWC Borel measurable. The PWC requirements in Assumption 3.5(ii) and (iii) and BFCP invertibility

in Assumption 3.5(i) are needed to achieve PWC consistency, and hence ensure finitely many abstract state regions (and are used in Lemmas 4.2, 4.3, and 6.1 below). Bimeasurability in Assumption 3.5(ii) ensures the existence of the value of a NS-CSG with respect to a reward structure (and is used in Proposition 5.3). For the remainder of this paper, unless otherwise stated we will assume Assumption 3.5 holds.

The observation function $obs_i$ for $1 \leq i \leq 2$ is implemented through an NN classifier, i.e., a continuous mapping from the finite real vector space $\mathbb{R}^m$ to the probability space over the finite percept classes $Per_i$, with the rule of selecting the class with the highest probability. Since $f_i$ is continuous it follows that it is also Borel measurable. To ensure that $obs_i$ is Borel measurable, we also need to consider the situation where the class with the highest probability is not unique. We assume such cases are resolved using a *tie-breaking rule* defined by a function $\kappa_i : 2^{Per_i} \rightarrow Per_i$ which, given a set of percepts, i.e., those with the highest probability, returns the selected percept. It follows that, if $\kappa_i$ is a Borel measurable function, then Assumption 3.5(ii) holds. The case when observation functions are implemented using ReLU networks are discussed in Section 6, but we remark that Assumption 3.5(ii) allows a wider range of functions than just NNs for implementing observations.

**4. Game structures.** NS-CSGs are a class of stochastic games with Borel state spaces in which the state space, observation function and reward function have structured form. Before designing our algorithms, we present several characteristics of these games, which distinguish them from the general case and which will be exploited later. For the remainder of this section we fix an NS-CSG C and reward structure $r$.

Since by Assumption 3.5 the perception function $obs_i$ and the reward functions $r_A$ and $r_S$ are all PWC, there exist two FCPs of the state space, respectively denoting the abstract state spaces for perception and reward computations.

DEFINITION 4.1 (Model-induced FCPs of the state space). *There exists:*
- *(i) a smallest FCP of $S$, called the* perception FCP *of $S$ and denoted $\Phi_P$, such that all states in any $\phi \in \Phi_P$ have the same agents' states, i.e., if $(s_1, s_2, s_E), (s'_1, s'_2, s'_E) \in \phi$, then $s_i = s'_i$ for $1 \leq i \leq 2$;*
- *(ii) for each joint action $\alpha \in A$, a smallest FCP of $S$, called the* reward FCP *of $S$ under $\alpha$ and denoted $\Phi_R^\alpha$, such that all states in any $\phi \in \Phi_R^\alpha$ have the same state reward and reward when $\alpha$ is chosen, i.e., if $s, s' \in \phi$, then $r_A(s, \alpha) = r_A(s', \alpha)$ and $r_S(s) = r_S(s')$.*

The FCPs given in Definition 4.1 allow us to abstract an uncountable state space into a finite set of regions for percept and reward computations. Using Assumption 3.5, the FCPs $\Phi_P$ and $\Phi_R^\alpha$ have the following properties, which will be used for the existence of the value of C and in our algorithms.

LEMMA 4.2 (BFCPs). *For each $\alpha \in A$, $\Phi_P$ and $\Phi_R^\alpha$ are BFCPs of $S$.*

*Proof.* We consider a region $\phi \in \Phi_P$. By Definition 4.1 all states in $\phi$ have the same agents' states, say $s_1 = (loc_1, per_1)$ and $s_2 = (loc_2, per_2)$. According to Assumption 3.5, $obs_i(loc_1, loc_2, \cdot) : S_E \rightarrow Per_i$ for $1 \leq i \leq 2$ is a PWC Borel measurable function. The pre-image of $(per_1, per_2)$ under $obs_1$ and $obs_2$ over $S$ given $s_1 = (loc_1, per_1)$ and $s_2 = (loc_2, per_2)$, denoted $obs^{-1}(per_1, per_2 \mid s_1, s_2)$, equals:

$$\{(s_1, s_2, s_E) \in S \mid obs_1(loc_1, loc_2, s_E) = per_1 \wedge obs_2(loc_1, loc_2, s_E) = per_2\}$$

and therefore is a Borel set of $S$. Since $\Phi_P$ is the smallest partition of $S$ satisfying Definition 4.1(i), the regions in $\Phi_P$, which lead to the percept $(per_1, per_2)$ given $s_1$

and $s_2$, have no common boundary. Thus, $obs^{-1}(per_1, per_2 \mid s_1, s_2)$ is a finite union of disjoint regions in $\Phi_P$ which include the agents' states $s_1$ and $s_2$. Thus, each such region is a Borel set of $S$, meaning that $\phi \in \mathcal{B}(S)$. Thus, $\Phi_P$ is a BFCP of $S$. Since $r_A(\,\cdot\,, \alpha) + r_S(\,\cdot\,)$ is a PWC Borel measurable function on $S$ by Assumption 3.5, we can similarly show that $\Phi_R^\alpha$ is a BFCP of $S$ for all $\alpha \in A$. $\qquad\square$

Recall that $\Theta(s, \alpha) = \{s' \mid \delta(s, \alpha)(s') > 0\}$ is the finite branching set from $s \in S$ under $\alpha \in A(s)$.

LEMMA 4.3 (Reachability consistency).   *For each $\alpha \in A$, there exists a refinement BFCP $\Phi$ of $\Phi_P$ such that, for each $\phi \in \Phi$ and $\phi' \in \Phi_P$, if $\delta(s, \alpha)$ is defined for $s \in \phi$, then there exists $p_\alpha(\phi, \phi') \in [0, 1]$ such that:*
  1. *either $\delta(s, \alpha)(s') = p_\alpha(\phi, \phi') = 0$ for all $s \in \phi$ and $s' \in \phi'$;*
  2. *or (i) if $s, \tilde{s} \in \phi$, then there exist unique states $s', \tilde{s}' \in S$ such that $s' = \Theta(s, \alpha) \cap \phi'$, $\tilde{s}' = \Theta(\tilde{s}, \alpha) \cap \phi'$ and $\delta(s, \alpha)(s') = \delta(\tilde{s}, \alpha)(\tilde{s}') = p_\alpha(\phi, \phi') > 0$, and (ii) there exists a bimeasurable, BFCP invertible function $q_\alpha : \phi \to \phi'$ such that $q_\alpha(s) = \Theta(s, \alpha) \cap \phi'$ for all $s \in \phi$.*

*Proof.* We compute the refinement $\Phi$ of $\Phi_P$ by dividing each $\phi$ of $\Phi_P$ such that the required reachability consistency holds. Now for any $\alpha \in A$ and $\phi \in \Phi_P$. By Definition 4.1, all states in $\phi$ have the same agents' states, say $s_1$ and $s_2$. To aid the proof, for each $\phi' \in \Phi_P$, we will construct a BFCP of $\phi$ based on $\phi'$, denoted $\Phi'(\phi, \phi')$, such that the reachability consistency to the region $\phi'$ holds in each region of $\Phi'(\phi, \phi')$. If $\delta(s, \alpha)$ is not defined for $s \in \phi$, we do not divide $\phi$ and let $\Phi'(\phi, \phi') = \{\phi\}$ for all $\phi' \in \Phi_P$ and the reachability consistency to $\phi'$ is preserved.

It remains to consider the case when $\delta(s, \alpha)$ is defined. Considering any $\phi' \in \Phi_P$, by Definition 4.1 there exists agent states $s_1' = (loc_1', per_1')$ and $s_2' = (loc_2', per_2')$ such that if $(s_1'', s_2'', s_E'') \in \phi'$ then $s_1 = s_1'$ and $s_2 = s_2'$. We have the following two cases to consider.
  - If $\{(s_1', s_2', \delta_E(s_E, \alpha)) \in S \mid (s_1, s_2, s_E) \in \phi\} \cap \phi' = \varnothing$, $\delta_1(s_1, \alpha)(loc_1') = 0$ or $\delta_2(s_2, \alpha)(loc_2') = 0$, then we do not divide $\phi$ and let $\Phi'(\phi, \phi') = \{\phi\}$ and we have $\delta(s, \alpha)(s') = p_\alpha(\phi, \phi') = 0$ for all $s \in \phi$ and $s' \in \phi'$.
  - If $(\cup_{s \in \phi} \Theta(s, \alpha)) \cap \phi'$ is non-empty, then since $\delta_E(\,\cdot\,, \alpha) : S_E \to S_E$ is BFCP invertible using Assumption 3.5 and $\phi'$ is a Borel measurable region there exists a BFCP $\Phi'(\phi, \phi')$ of $\phi$ such that for each $\phi_1 \in \Phi'(\phi, \phi')$:
    - either $\delta(s, \alpha)(s') = p_\alpha(\phi_1, \phi') = 0$ for all $s \in \phi_1$ and $s' \in \phi'$;
    - or for $s, \tilde{s} \in \phi_1$ there exist unique states $s', \tilde{s}' \in S$ such that $s' = \Theta(s, \alpha) \cap \phi'$, $\tilde{s}' = \Theta(\tilde{s}, \alpha) \cap \phi'$ and $\delta(s, \alpha)(s') = \delta(\tilde{s}, \alpha)(\tilde{s}') = p_\alpha(\phi_1, \phi') > 0$.
    It remains to show that the bimeasurable, BFCP invertible function $q_\alpha$ of (ii) exists, which follows from the the fact that $\delta_E(\,\cdot\,, \alpha) : S_E \to S_E$ is BFCP invertible.

Finally, we divide $\phi$ into a BFCP $\sum_{\phi' \in \Phi_P} \Phi'(\phi, \phi')$, and therefore each region of this BFCP has the required reachability consistency. $\qquad\square$

**5. Value of zero-sum NS-CSGs.** We now proceed by establishing the value of an NS-CSG C for an objective $Y$, i.e. for a reward structure $r$ and discount factor $\beta$. We prove the existence of this value, which is a fixed point of a minimax operator. Using Banach's fixed-point theorem, a sequence of bounded, Borel measurable functions converging to this value is constructed.

Given a state $s$ and profile $\sigma = (\sigma_1, \sigma_2)$, we denote by $\mathbb{E}_s^\sigma[Y]$ the expected value of the objective $Y$ when starting from state $s$, defined in (3.1). The functions, for any

$s \in S$, given by:

$$\underline{V}(s) := \sup_{\sigma_1 \in \Sigma_1} \inf_{\sigma_2 \in \Sigma_2} \mathbb{E}_s^{\sigma_1, \sigma_2}[Y] \quad \text{and} \quad \overline{V}(s) := \inf_{\sigma_2 \in \Sigma_2} \sup_{\sigma_1 \in \Sigma_1} \mathbb{E}_s^{\sigma_1, \sigma_2}[Y]$$

are called the *lower value* and the *upper value*, respectively, of $Y$.

DEFINITION 5.1 (Value function). *If $\underline{V}(s) = \overline{V}(s)$ for all $s \in S$, then $\mathsf{C}$ is determined with respect to the objective $Y$ and the common function is called the* value *of $\mathsf{C}$, denoted by $V^\star(\,\cdot\,)$, with respect to $Y$.*

We first introduce the spaces of feasible state-action and state-action-distribution pairs and properties of these spaces. More precisely, for $1 \leq i \leq 2$, we define:

$$\Xi_i := \{(s, a_i) \in (S \times A_i) \mid a_i \in A_i(s)\}$$
$$\Lambda_i := \{(s, u_i) \in (S \times \mathbb{P}(A_i)) \mid u_i \in \mathbb{P}(A_i(s))\}$$
$$\Xi_{12} := \{(s, (a_1, a_2)) \in (S \times (A_1 \times A_2)) \mid a_1 \in A_1(s) \wedge a_2 \in A_2(s)\}$$
$$\Lambda_{12} := \{(s, (u_1, u_2)) \in (S \times (\mathbb{P}(A_1) \times \mathbb{P}(A_2))) \mid u_1 \in \mathbb{P}(A_1(s)) \wedge u_2 \in \mathbb{P}(A_2(s))\}.$$

LEMMA 5.2 (Borel spaces). *For $1 \leq i \leq 2$, the spaces $\Xi_i$ and $\Lambda_i$ are Borel sets of $S \times A_i$ and $S \times \mathbb{P}(A_i)$, respectively. Furthermore, the spaces $\Xi_{12}$ and $\Lambda_{12}$ are Borel sets of $S \times (A_1 \times A_2)$ and $S \times (\mathbb{P}(A_1) \times \mathbb{P}(A_2))$, respectively.*

*Proof.* We first consider $\Xi_i$ and $\Lambda_i$ in the case when $i = 1$, and the case for $i = 2$ follows similarly. Since $A_1$ is finite, the sets $\Xi_1$ and $\Lambda_1$ can be rearranged as:

$$\Xi_1 = \bigcup_{\hat{A}_1 \subseteq A_1} \left( \{s_1 \mid \Delta_1(s_1) = \hat{A}_1\} \times S_2 \times S_E \times \hat{A}_1 \right) \cap (S \times A_1)$$
$$\Lambda_1 = \bigcup_{\hat{A}_1 \subseteq A_1} \left( \{s_1 \mid \Delta_1(s_1) = \hat{A}_1\} \times S_2 \times S_E \times \mathbb{P}(\hat{A}_1) \right) \cap (S \times \mathbb{P}(A_1))$$

respectively. Since $\hat{A}_1$ is a subset of the finite set $A_1$, then the sets $\hat{A}_1$ and $\mathbb{P}(\hat{A}_1)$ are Borel sets of $A_1$ and $\mathbb{P}(A_1)$, respectively. Since $S_1$ is a finite set, for any $\hat{A}_1 \subseteq A_1$, the set $\{s_1 \mid \Delta_1(s_1) = \hat{A}_1\}$ is a Borel set of $S_1$. Since $S_2$ and $S_E$ are both Borel sets by Lemma 3.3, then the result follows by Theorem 1.10 [21, Chapter 1]. By the same argument, for $\Xi_{12}$ and $\Lambda_{12}$ we have:

$$\bigcup_{\hat{A}_1 \subseteq A_1, \hat{A}_2 \subseteq A_2} \left( \left( \{s_1 \mid \Delta_1(s_1) = \hat{A}_1\} \times \{s_2 \mid \Delta_2(s_2) = \hat{A}_2\} \times S_E \right) \times (\hat{A}_1 \times \hat{A}_2) \right)$$
$$\bigcup_{\hat{A}_1 \subseteq A_1, \hat{A}_2 \subseteq A_2} \left( \left( \{s_1 \mid \Delta_1(s_1) = \hat{A}_1\} \times \{s_2 \mid \Delta_2(s_2) = \hat{A}_2\} \times S_E \right) \times (\mathbb{P}(\hat{A}_1) \times \mathbb{P}(\hat{A}_2)) \right)$$

intersecting with $S \times A_1 \times A_2$ and $S \times \mathbb{P}(A_1) \times \mathbb{P}(A_2)$, and thus $\Xi_{12}$ and $\Lambda_{12}$ are also Borel sets of the respective spaces. $\square$

Using Lemmas 4.2, 4.3, and 5.2, we have the following result.

PROPOSITION 5.3 (Stochastic kernel transition function). *The probabilistic transition function $\delta$ is a stochastic kernel.*

*Proof.* For every $(s, \alpha) \in \Xi_{12}$, we have $\delta(s, \alpha)(\,\cdot\,) \in \mathbb{P}(S)$. We show that if $B \in \mathcal{B}(S)$, then $\delta(\,\cdot\,, \,\cdot\,)(B) : (S \times A) \to \mathbb{R}$ is Borel measurable on $\Xi_{12}$. More precisely we show that, for any $c \in \mathbb{R}$, pre-image of the Borel set $[c, \infty)$ of $\mathbb{R}$ under $\delta(\,\cdot\,, \,\cdot\,)(B)$ which is given by:

$$\delta^{-1}([c, \infty))(B) = \{(s, \alpha) \in \Xi_{12} \mid \delta(s, \alpha)(B) \geq c\}$$

is an element of $\mathcal{B}(\Xi_{12})$. If $c > 1$, then $\delta^{-1}([c, \infty))(B) = \varnothing \in \mathcal{B}(\Xi_{12})$, and if $c \leq 0$, then $\delta^{-1}([c, \infty))(B) = \Xi_{12} \in \mathcal{B}(\Xi_{12})$.

Therefore it remains to consider the case when $0 < c \leq 1$. For $\alpha \in A$, let $\Phi$ be the refinement of $\Phi_P$ meeting the conditions of Lemma 4.3. For each $\phi \in \Phi$ and $\phi' \in \Phi_P$ such that $p_\alpha(\phi, \phi') > 0$ using Lemma 4.3, let $q_\alpha : \phi \to \phi'$ be the associated bimeasurable, BFCP invertible function from Lemma 4.3. The image of $\phi$ under $q_\alpha$ into $\phi'$ for a fixed $\alpha$ is given by:

$$\hat{q}_\alpha(\phi, \phi') = \{s' \in \phi' \mid s' = q_\alpha(s) \wedge s \in \phi\}.$$

By Lemmas 4.2 and 4.3 both $\phi$ and $\phi'$ are Borel sets and $q_\alpha$ is bimeasurable, and therefore $\hat{q}_\alpha(\phi, \phi')$ is a Borel set. Next, the pre-image of the Borel set $\hat{q}_\alpha(\phi, \phi') \cap B$ under $q_\alpha$ over the region $\phi$ is given by:

$$\hat{q}_\alpha^{-1}(\phi, \hat{q}_\alpha(\phi, \phi') \cap B) = \{s \in \phi \mid q_\alpha(s) \in \hat{q}_\alpha(\phi, \phi') \cap B\}$$

and hence, as $q_\alpha$ is Borel measurable, is a Borel set. By combining this result with Lemma 4.3, all states in $\hat{q}_\alpha^{-1}(\phi, \hat{q}_\alpha(\phi, \phi') \cap B)$ are transferred to $B$ with the same probability $p_\alpha(\phi, \phi')$. We denote the set of all transition probabilities from $\phi$ for a fixed $\alpha$ by $P_\alpha(\phi) = \{p_\alpha(\phi, \phi') > 0 \mid \phi' \in \Phi_P\}$. Then, the collection of the subsets of $P_\alpha(\phi)$ whose all elements' sum is greater or equal to $c$ is defined as:

$$P_\alpha^{\geq c}(\phi) := \left\{ P' \subseteq P_\alpha(\phi) \mid \sum_{p' \in P'} p' \geq c \right\}$$

which is finite. For every set $P' \in P_\alpha^{\geq c}(\phi)$, all states in the set:

$$O_\alpha(\phi, P') = \bigcap_{p_\alpha(\phi, \phi') \in P'} \hat{q}_\alpha^{-1}(\phi, \hat{q}_\alpha(\phi, \phi') \cap B)$$

which reaches $B$ under $\alpha$ with probability greater or equal to $c$, where $O_\alpha(\phi, P')$ is a Borel set as $P'$ is a finite set. Thus, under a fixed $\alpha$, the states in $\phi$ reaching $B$ with probability greater or equal to $c$ are given by:

$$O_\alpha(\phi) = \bigcup_{P' \in P_\alpha^{\geq c}(\phi)} O_\alpha(\phi, P')$$

which is a Borel set since $P_\alpha^{\geq c}(\phi)$ is a finite set. Finally, we have:

$$\delta^{-1}([c, \infty))(B) = \bigcup_{\phi \in \Phi} \bigcup_{\alpha \in A} \{(s, \alpha) \in \Xi_{12} \mid s \in O_\alpha(\phi)\}$$

from which it follows that $\delta^{-1}([c, \infty))(B) \in \mathcal{B}(\Xi_{12})$ by combining Lemmas 4.3 and 5.2 as required. □

Before discussing the value function, we first introduce an operator based on the classical Bellman equation.

DEFINITION 5.4 (Minimax operator). *Given a bounded, Borel measurable real-valued function $V \in \mathbb{F}(S)$, the minimax operator $T : \mathbb{F}(S) \to \mathbb{F}(S)$ is defined, for any $s \in S$, by:*

$$[TV](s) := \max_{u_1 \in \mathbb{P}(A_1(s))} \min_{u_2 \in \mathbb{P}(A_2(s))} \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} Q(s, (a_1, a_2), V) u_1(a_1) u_2(a_2)$$

*where for any $\alpha \in A(s)$:*

$$Q(s, \alpha, V) := r_A(s, \alpha) + r_S(s) + \beta \sum_{s' \in \Theta(s, \alpha)} \delta(s, \alpha)(s') V(s').$$

The following theorem holds for the value function and minimax operator, which relies on Lemmas 3.3 and 5.2 and Proposition 5.3.

THEOREM 5.5 (Value function). *If $C$ is an NS-CSG and $Y$ is a discounted zero-sum objective, then*

(i) *$C$ is determined with respected to $Y$, i.e., $V^\star$ exists;*

(ii) *$V^\star$ is the unique fixed point of the operator $T$;*

(iii) *$V^\star$ is a bounded, Borel measurable function.*

*Proof.* The proof is through showing that the NS-CSG $C$ is an instance of a zero-sum stochastic game with Borel model as presented in [15].

From Lemma 3.3, we have that $A_1$, $A_2$ and $S$ are complete and separable metric spaces. By Lemma 5.2, the spaces $\Xi_i$ and $\Lambda_i$ are Borel sets of $S \times A_i$ and $S \times \mathbb{P}(A_i)$ for $1 \leq i \leq 2$, respectively. By Proposition 5.3, $\delta$ is a Borel stochastic kernel.

Furthermore, since $r_A + r_S$ is bounded on $S \times A$, it follows that $C$ with respect to the zero-sum objective $Y$ is an instance of a zero-sum stochastic game with Borel model and discounted payoffs introduced in [15]. Therefore (i) follows from [15, Theorems 2 and 3], and (ii) from the discounted case of [15, Theorem 1]. Finally, for (iii), since $\beta \in (0,1)$, we have that $V^\star$ is bounded, and therefore $V^\star$ is Borel measurable using [15, Lemma 3]. □

Since zero-sum NS-CSGs are determined games, as the following result demonstrates, optimal strategies exist.

LEMMA 5.6 (Optimal strategies). *There is a strategy $\sigma_1^\star$ for $\mathsf{Ag}_1$ (and a dual strategy $\sigma_2^\star$ for $\mathsf{Ag}_2$), called an optimal strategy of $\mathsf{Ag}_1$, such that under $\sigma_1^\star$ in any state $s$ the expected discounted reward is at least $V^\star(s)$ regardless of the strategy of $\mathsf{Ag}_2$; more precisely, we have $\inf_{\sigma_2 \in \Sigma_2} \mathbb{E}_s^{\sigma_1^\star, \sigma_2}[Y] = V^\star(s)$ for all $s \in S$.*

*Proof.* The result follows from [15, Theorems 2 and 3]. □

Furthermore, the following lemma guarantees that value iteration (VI) can be used to find the value function.

LEMMA 5.7 (Convergence sequence). *For any $V^0 \in \mathbb{F}(S)$, the sequence $\langle V^t \rangle_{t \in \mathbb{N}}$, where $V^{t+1} = TV^t$, converges to $V^\star$. Moreover, each $V^t$ is bounded, Borel measurable.*

*Proof.* Since $r_A + r_S$ is bounded, using [15, Lemma 2] we have that, if $V^t$ is bounded, Borel measurable, then so is $TV^t$. The result then follows from the fact that $V^\star(s) = \lim_{t \to \infty} V^t(s)$ for all $s \in S$ if $V^{t+1} = TV^t$ for all $t \in \mathbb{N}$ [15]. □

**6. Value iteration.** Despite the convergence result of Lemma 5.7, in practice there may not exist effective representations of the bounded Borel measurable functions $V^t$, due to the uncountable state space. We now show how VI can be used to *approximate* the values of $C$ with respect to a discounted accumulated reward objective $Y$, and the corresponding optimal profile, based on a sequence of bounded, PWC Borel measurable functions.

LEMMA 6.1 (PWC and measurable consistency). *If $V \in \mathbb{F}(S)$ is bounded, PWC Borel measurable, then*

(i) *$Q(\,\cdot\,, \alpha, V)$ is bounded, PWC Borel measurable for $\alpha \in A$;*

(ii) *$TV$ is bounded, PWC Borel measurable.*

*Proof.* We first proof that (i) holds. Considering any bounded, PWC Borel measurable $V$ and joint action $\alpha \in A$, since $r_A(\,\cdot\,, \alpha) + r_S(\,\cdot\,)$ is bounded, PWC Borel measurable by Assumption 3.5, the fact that $Q(\,\cdot\,, \alpha, V)$ is bounded, PWC Borel measurable follows if we can show that the function $\overline{Q}(\,\cdot\,, \alpha, V)$ where:

$$\overline{Q}(\,\cdot\,, \alpha, V) := \sum_{s' \in \Theta(\cdot, \alpha)} \delta(\,\cdot\,, \alpha)(s') V(s')$$

is bounded, PWC Borel measurable. Boundedness follows because $V$ is bounded. The indicator function of a subset $S' \subseteq S$ is the function $\chi_{S'} : S \to \mathbb{R}$ such that $\chi_{S'}(s) = 1$ if $s \in S'$ and 0 otherwise. Now $\chi_{S'}$ is Borel measurable if and only if $S'$ is a Borel set of $S$ [27]. For clarity, we use $q_\alpha(s; \phi, \phi')$ to refer to $q_\alpha$ from Lemma 4.3 for $\alpha \in A$, $s \in \phi$, $\phi \in \Phi$ and $\phi' \in \Phi_P$ (where again $\Phi$ is from Lemma 4.3). For any $s \in S$ such that $\delta(s, \alpha)$ is defined, we have:

$$
\begin{aligned}
\overline{Q}(s, \alpha, V) &= \sum_{\phi \in \Phi} \chi_\phi(s) \sum_{s' \in \Theta(s,\alpha)} \delta(s, \alpha)(s') V(s') \\
&= \sum_{\phi \in \Phi} \chi_\phi(s) \sum_{\phi' \in \Phi_P} p_\alpha(\phi, \phi') V(q_\alpha(s; \phi, \phi')) \qquad \text{by Lemma 4.3} \\
&= \sum_{\phi \in \Phi} \sum_{\phi' \in \Phi_P} p_\alpha(\phi, \phi') \chi_\phi(s) V(q_\alpha(s; \phi, \phi')) \qquad \text{rearranging.}
\end{aligned}
$$

Since $\phi$ is a Borel set of $S$, we have that $\chi_\phi$ is Borel measurable. Next, we show that $V(q_\alpha(\,\cdot\,; \phi, \phi'))$ is Borel measurable on $\phi$. Let $\Phi_V$ be a constant-BFCP of $S$ for $V$. Given $c \in \mathbb{R}$, we denote by $\Phi_V^{\geq c}$ the set of regions in $\Phi_V$ on which $V \geq c$ holds. The pre-image of $[c, \infty)$ under $V(q_\alpha(\,\cdot\,; \phi, \phi'))$ defined on $\phi$ is given by:

$$
\begin{aligned}
V^{-1}(q_\alpha([c, \infty); \phi, \phi')) &= \{ s \in \phi \mid V(q_\alpha(s; \phi, \phi')) \geq c \} \\
&= \bigcup_{\phi_V \in \Phi_V^{\geq c}} \{ s \in \phi \mid q_\alpha(s; \phi, \phi') \in \phi_V \}.
\end{aligned}
$$

Since $q_\alpha(s; \phi, \phi')$ is Borel measurable in $s \in \phi$ ($q_\alpha$ in Lemma 4.3) and $\phi_V$ is a Borel set of $S$, then $\{ s \in \phi \mid q_\alpha(s; \phi, \phi') \in \phi_V \}$ is a Borel set of $\phi$. Since $V^{-1}(q_\alpha([c, \infty); \phi, \phi'))$ is also a Borel set of $\phi$ by noting that $\Phi_V^{\geq c}$ is finite, it follows that $V(q_\alpha(\,\cdot\,; \phi, \phi'))$ is Borel measurable on $\phi$. Therefore $\overline{Q}(\,\cdot\,, \alpha, V)$ is Borel measurable.

Next, since $q_\alpha(\,\cdot\,; \phi, \phi')$ is BFCP invertible on $\phi$ by Lemma 4.3, there exists a BFCP $\Phi_q$ of $\phi$ such that all states in each region of $\Phi_q$ are mapped into the same region of $\Phi_V$ under $q_\alpha(\,\cdot\,; \phi, \phi')$. Following this, $V(q_\alpha(\,\cdot\,; \phi, \phi'))$ is constant on each region of $\Phi_q$. Therefore, using the fact that $\chi_\phi$ is PWC, it follows that $\overline{Q}(\,\cdot\,, \alpha, V)$ is PWC, which completes the proof of (i).

Regarding (ii), from Lemma 5.7 we have that $TV$ is bounded, Borel measurable. Since $Q(\,\cdot\,, \alpha, V)$ is PWC for any joint action $\alpha \in A$, $A(s)$ is PWC and $A$ is finite, it follows that $TV$ is PWC using the fact that the value of a zero-sum normal-formal game induced at every $s \in S$ is unique. □

We use the above PWC and measurable consistencies to iteratively construct a sequence of representation-friendly functions $\langle V^t \rangle_{t \in \mathbf{N}}$, where $V^0$ is an arbitrary PWC, Borel measurable function and $V^{t+1} = TV^t$. The following theorem is then a direct consequence of Lemmas 5.7 and 6.1.

THEOREM 6.2 (PWC convergence sequence). *From an arbitrary bounded, PWC Borel measurable function $V^0 \in \mathbb{F}(S)$, the sequence $\langle V^t \rangle_{t \in \mathbb{N}}$ such that $V^{t+1} = TV^t$, converges to $V^\star$. Moreover, each $V^t$ is bounded, PWC Borel measurable.*

Theorem 6.2 ensures convergence of bounded, PWC measurable functions. Using this, Algorithm 6.1 presents our VI scheme, in which the initial value function $V^0$ is a 0-valued PWC Borel measurable function defined over the BFCP $\Phi_P + \sum_{\alpha \in A} \Phi_R^\alpha$ of $S$. The steps of a single iteration of our VI algorithm are illustrated in Figure 1. These steps make use of the function $Preimage\_BFCP(\Phi_{V^t}, \Phi_P, \langle \Phi_R^\alpha \rangle_{\alpha \in A})$ in Algorithm 6.2, to compute a refinement of $\Phi_P + \sum_{\alpha \in A} \Phi_R^\alpha$ that is a consistent-BFCP of $S$ for $\delta$ to $\Phi_{V^t}$. Then, in order to compute the value $V'_\phi$ over each region $\phi \in \Phi$, we take one

---

ALGORITHM 6.1 PWC-function VI

---

1: **Input:** NS-CSG $\mathsf{C}$, discounted zero-sum objective $Y$, perception FCP $\Phi_P$, reward
   FCPs $\langle\Phi_R^\alpha\rangle_{\alpha\in A}$, threshold $\varepsilon$
2: **Output:** Approximate value function $V$
3: $\Phi_{V^0} \leftarrow \Phi_P + \sum_{\alpha\in A}\Phi_R^\alpha$,  $V_\phi^0 \leftarrow 0$ for all $\phi\in\Phi_{V^0}$,  $V^0 \leftarrow \langle V_\phi^0\rangle_{\phi\in\Phi_{V^0}}$,
4: $\bar\varepsilon \leftarrow 2\varepsilon$,  $t \leftarrow 0$
5: **while** $\bar\varepsilon > \varepsilon$ **do**
6:     $\Phi \leftarrow Preimage\_BFCP(\Phi_{V^t}, \Phi_P, \langle\Phi_R^\alpha\rangle_{\alpha\in A})$  (Algorithm 6.2)
7:     **for** $\phi\in\Phi$ **do**
8:         Take one state $s\in\phi$,  $V_\phi' \leftarrow (TV^t)(s)$
9:     $(\Phi_{V^{t+1}}, V^{t+1}) \leftarrow (\Phi, \langle V_\phi'\rangle_{\phi\in\Phi})$,
10:    $\bar\varepsilon \leftarrow Dist(V^{t+1}, V^t)$,
11:    $t \leftarrow t+1$
12: **return** $V \leftarrow V^t$

---

ALGORITHM 6.2 BFCP iteration for PWC-function VI

---

1: **Input:** A BFCP $\Phi$ of $S$, perception FCP $\Phi_P$, reward FCPs $\langle\Phi_R^\alpha\rangle_{\alpha\in A}$
2: **Output:** A BFCP of $S$
3: **procedure** $Preimage\_BFCP(\Phi, \Phi_P, \langle\Phi_R^\alpha\rangle_{\alpha\in A})$
4:     $\Phi_{\mathrm{pre}} \leftarrow \varnothing$
5:     **for** $\phi\in\Phi_P + \sum_{\alpha\in A}\Phi_R^\alpha$ **do**
6:         $B \leftarrow \varnothing$
7:         **for** $\alpha\in A$, $\phi'\in\{\phi'\in\Phi \mid (\cup_{s\in\phi}\Theta(s,\alpha))\cap\phi'\neq\varnothing\}$ **do**
8:             $B \leftarrow B \cup \big\{\{s\in\phi \mid \Theta(s,\alpha)\cap\phi'\neq\varnothing\}\big\}$
9:         $\Phi_{\mathrm{pre}} \leftarrow \Phi_{\mathrm{pre}} \cup \{\phi_1\in Intersect(B)\}$
10:    **return** $\Phi_{\mathrm{pre}}$

---

state $s\in\phi$ and then find the value of a zero-sum normal form game [33] at $s$ induced
by Definition 5.4.

As a convergence criterion for VI, we detect when the difference between successive
value approximations falls below a pre-specified threshold $\varepsilon$ (as usual for VI, this does
*not* guarantee an $\varepsilon$-optimal approximation). The function $Dist(V^{t+1}, V^t)$ computes
the difference between $V^{t+1}$ and $V^t$, which may be of different sizes due to the possible
inconsistency of $\Phi_{V^{t+1}}$ and $\Phi_{V^t}$. An intuitive method is to first evaluate $V^{t+1}$ and
$V^t$ at a finite set of points, and then compute the maximum pointwise difference. In
the usual manner for VI, an approximately optimal strategy can be extracted from
the final step of the computation.

In Algorithm 6.2, the function $Intersect(B)$ returns a refinement of $\phi$ by comput-
ing all pairwise intersections for regions in $B$ such that the refinement is a consistent-
BFCP of $\phi$ for $\delta$ to $\Phi$. The following corollary then follows from Lemmas 4.3 and 6.1.

COROLLARY 6.3 (BFCP iteration for VI). *In Algorithm 6.2, $\Phi_{\mathrm{pre}}$ is a refinement
of $\Phi_P + \sum_{\alpha\in A}\Phi_R^\alpha$ and a consistent-BFCP of $S$ for $\delta$ to $\Phi$.*

**6.1. Polytope regions.** The VI algorithm above assumes that each region in
a BFCP is finitely representable. We now briefly discuss the use of BFCPs defined
by *polytopes*, which work for the ReLU NNs. A polytope $\phi\subseteq\mathbb{R}^m$ is an intersection
of $\ell$ halfspaces $\{x\in\mathbb{R}^m \mid g_k(x)\geq 0$ for $1\leq k\leq\ell\}$, where $g_k(x) = W_k^\top x + b_k$ is a
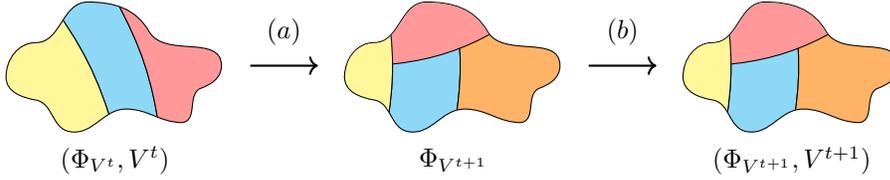
FIG. 1.  *PWC-function VI.* (a) *Find new BFCP* $\Phi_{V^{t+1}}$: *refine* $\Phi_P + \sum_{\alpha \in A} \Phi_R^\alpha$ *to be a consistent-BFCP of S for $\delta$ to* $\Phi_{V^t}$; (b) *compute a value for each* $\phi \in \Phi_{V^{t+1}}$: *take one state* $s \in \phi$ *and compute* $V^{t+1}$ *by assigning each region* $\phi$ *with value* $(TV^t)(s)$.

linear function, i.e. $W_k \in \mathbb{R}^m$ and $b_k \in \mathbb{R}$, for $1 \le k \le \ell$. If $\phi_1$ and $\phi_2$ are polytopes, represented by $\{(W_k, b_k) \mid 1 \le k \le \ell'\}$ and $\{(W_k, b_k) \mid \ell' + 1 \le k \le \ell\}$, respectively, then the *sum* (refinement polytope) of $\phi_1$ and $\phi_2$, denoted $\phi_1 + \phi_2$, is the intersection of $\ell$ halfspaces $\{x \in \mathbb{R}^m \mid g_k(x) \ge 0 \text{ for } 1 \le k \le \ell\}$ and can be represented as $\{(W_k, b_k) \mid 1 \le k \le \ell\}$. Therefore the sum $\Phi_1 + \Phi_2$ of two BFCPs $\Phi_1$ and $\Phi_2$ can be computed by considering the sum $\phi_1 + \phi_2$ of all combinations of regions $\phi_1 \in \Phi_1$ and $\phi_2 \in \Phi_2$.

The *image* of a polytope $\phi = \{x \in \mathbb{R}^m \mid g_k(x) \ge 0 \text{ for } 1 \le k \le \ell\}$ under a linear function $f : \mathbb{R}^m \to \mathbb{R}^m$, where $f(x) = Dx + b$, $D \in \mathbb{R}^{m \times m}$ is non-singular and $b \in \mathbb{R}^m$, is the polytope $f(\phi) = \{x \in \mathbb{R}^m \mid W_k^\top D^{-1}x + b_k - W_k^\top D^{-1}b \ge 0, 1 \le k \le \ell\}$ and can be represented as $\{(D^{-\top}W_k, b_k - W_k^\top D^{-1}b) \mid 1 \le k \le \ell\}$. The *preimage* of $\phi$ under $f$ is the polytope $f^{-1}(\phi) = \{x \in \mathbb{R}^m \mid W_k^\top Dx + b_k + W_k^\top b \ge 0, 1 \le k \le \ell\}$ and can be represented as $\{(D^\top W_k, b_k + W_k^\top b) \mid 1 \le k \le \ell\}$. Checking the feasibility of a set constrained by a set of linear inequalities can be solved by a linear program solver.

**6.2. ReLU networks.** If each observation function $obs_i$ is implemented via a ReLU NN and there exist polytope constant-BFCPs for $r_A(\,\cdot\,, \alpha)$ and $r_S$ for all $\alpha \in A$, then all regions in $\Phi_P$ and $\Phi_R^\alpha$ for $\alpha \in A$ are polytopes. If the inverse function of $\delta_E(\,\cdot\,, \alpha)$ is linear and $\phi'$ is a polytope (line 8 in Algorithm 6.2), then $\{s \in \phi \mid \Theta(s, \alpha) \cap \phi' \ne \varnothing\}$ is a polytope. Therefore each region in $\Phi_{\text{pre}}$ is a polytope after every iteration and the operations over polytopes, including intersections, image and preimgae computations directly follow from the computation above.

**7. Policy iteration.** In this section, we show how policy iteration (PI) can be used to approximate the values and optimal strategies of an NS-CSGs $\mathsf{C}$ with respect to a discounted accumulated reward objective $Y$. Our algorithm takes ideas from recent work [5], which proposed a new PI method to solve zero-sum stochastic games with finite state spaces. Our algorithm is the first PI algorithm for CSGs with Borel state spaces and with a convergence guarantee. The authors of [37] presented a mixed value and policy iteration method for stochastic optimal control models with Borel state space, but not for games.

**7.1. Minimax-action-free PI.** Let $\gamma \in \mathbb{R}$ be a constant such that $\gamma > 1$ and $\gamma\beta < 1$, which will be used to distribute the discount factor $\beta$ between policy evaluation and policy improvement of the two agents. Before introducing the two operators, we require the notion of a stationary Stackelberg (follower) strategy for $\mathsf{Ag}_2$, which is a stochastic kernel $\overline{\sigma}_2 : \Lambda_1 \to \mathbb{P}(A_2)$, i.e., $\overline{\sigma}_2 \in \mathbb{P}(A_2 \mid \Lambda_1)$ such that $\overline{\sigma}_2(A_2(s) \mid (s, u_1)) = 1$ for all $(s, u_1) \in \Lambda_1$. This strategy is introduced only for the PI algorithm and implies that $\mathsf{Ag}_2$ makes decisions conditioned on the current state $s$ and the current choice, i.e. action distribution $u_1$, of $\mathsf{Ag}_1$. We denote by $\overline{\Sigma}_2$ the set of all stationary Stackelberg strategies for $\mathsf{Ag}_2$.

DEFINITION 7.1 (Operator for the Max-Min value).   *For each strategy $\sigma_1 \in \Sigma_1$ of $\mathsf{Ag}_1$ and bounded, Borel measurable real-valued function $V_2 \in \mathbb{F}(\Lambda_1)$, we define the operator $H^1_{\sigma_1,V_2} : \mathbb{F}(\Lambda_1) \to \mathbb{F}(S)$ such that for $J_2 \in \mathbb{F}(\Lambda_1)$ and $s \in S$:*

$$H^1_{\sigma_1,V_2}(J_2)(s) = \frac{1}{\gamma}\min\{J_2(s,\sigma_1(s)), V_2(s,\sigma_1(s))\} = \frac{1}{\gamma}\min\{J_2(s,u_1), V_2(s,u_1)\}$$

*where $\sigma_1(s) = u_1 \in \mathbb{P}(A_1(s))$.*

DEFINITION 7.2 (Operator for the Min-Max value).   *For each Stackelberg (follower) strategy $\overline{\sigma}_2 \in \overline{\Sigma}_2$ of $\mathsf{Ag}_2$ and each $V_1 \in \mathbb{F}(S)$, we define the operator $H^2_{\overline{\sigma}_2,V_1} : \mathbb{F}(S) \to \mathbb{F}(\Lambda_1)$ such that for $J_1 \in \mathbb{F}(S)$ and $(s,u_1) \in \Lambda_1$:*

$$H^2_{\overline{\sigma}_2,V_1}(J_1)(s,u_1) = \ = \sum_{a_1 \in A_1(s)}\sum_{a_2 \in A_2(s)} Q(s,(a_1,a_2),\gamma\max[J_1,V_1])u_1(a_1)u_2(a_2)$$

*where $\overline{\sigma}_2(\cdot \mid (s,u_1)) = u_2 \in \mathbb{P}(A_2(s))$.*

Unlike the classical PI algorithms by Hoffman and Karp [12] and Pollatschek and Avi-Itzhak [23], following [5], our PI algorithm separates the policy evaluation and policy improvement of the maximiser ($\mathsf{Ag}_1$) and the minimiser ($\mathsf{Ag}_2$) through the use of the operators from Definition 7.1 and Definition 7.2, respectively.  To track the value functions after performing policy evaluation of $\mathsf{Ag}_1$ and $\mathsf{Ag}_2$, our PI algorithm introduces value functions $J_1$ and $J_2$. In addition, the value functions $V_1$ and $V_2$ are introduced to avoid the oscillatory behavior of Pollatschek and Avi-Itzhak PI algorithm [23], thus ensuring convergence, and are updated only during policy improvement. The role of $\gamma$ is to split the discount factor $\beta$ such that all the operators corresponding to policy evaluation and policy improvement of the two agents are contraction mappings, which will ensure convergence.

Before presenting our PI algorithm, we define two classes of functions, which play a key role in characterizing the functions and strategies generated during each iteration of our PI algorithm.

DEFINITION 7.3 (Bounded, CON-PWL Borel measurable function).   *A function $f \in \mathbb{F}(\Lambda_1)$ is a bounded, constant-piecewise-linear (CON-PWL) Borel measurable function if there exists a BFCP $\Phi$ of $S$ such that for each $\phi \in \Phi$, $A_1(s) = A_1(s')$ for all $s, s' \in \phi$ and it generates $\Theta = \{\theta(\phi) \mid \phi \in \Phi\}$ where $\theta(\phi) = \{(s,u_1) \in \Lambda_1 \mid s \in \phi\}$, a BFCP of $\Lambda_1$, such that for each region $\theta(\phi) \in \Theta$:*
   *(i) $f(\,\cdot\,,u_1) : \phi \to \mathbb{R}$ is constant for all $u_1 \in \mathbb{P}(A_1(s))$ where $s \in \phi$;*
   *(ii) $f(s,\,\cdot\,) : \mathbb{P}(A_1(s)) \to \mathbb{R}$ is bounded, PWL for all $s \in \phi$.*

DEFINITION 7.4 (CON-PWC stochastic kernel).   *A function $f \in \overline{\Sigma}_2$ is a constant, piecewise-constant (CON-PWC) stochastic kernel if there exists a BFCP $\Phi$ of $S$ such that for each $\phi \in \Phi$, $A(s) = A(s')$ for all $s, s' \in \phi$ and it generates $\Theta = \{\theta(\phi) \mid \phi \in \Phi\}$ where $\theta(\phi) = \{(s,u_1) \in \Lambda_1 \mid s \in \phi\}$, a BFCP of $\Lambda_1$, such that for each region $\theta(\phi) \in \Theta$:*
   *(i) $f(\,\cdot\,,u_1) : \phi \to \mathbb{P}(A_2(s))$ is constant for all $u_1 \in \mathbb{P}(A_1(s))$ where $s \in \phi$;*
   *(ii) $f(s,\,\cdot\,) : \mathbb{P}(A_1(s)) \to \mathbb{P}(A_2(s))$ is PWC for all $s \in \phi$.*

Figure 2 presents an examples of a bounded, CON-PWL Borel measurable function and CON-PWC stochastic kernel over a region. Each bounded, CON-PWL Borel measurable function $f$ can be represented by a finite set of vectors $\{(D_{\phi,\phi'}, b_{\phi,\phi'}) \in \mathbb{R}^{|A_1|} \times \mathbb{R} \mid \phi \in \Phi \wedge \phi' \in \Phi'(\phi)\}$ such that $f(s,u_1) = D_{\phi,\phi'}^\top u_1 + b_{\phi,\phi'}$ for all $s \in \phi$ and $u_1 \in \phi'$, where $\Phi$ is a BFCP of $S$ for $f$ using Definition 7.3 and $\Phi'(\phi)$ is a BFCP
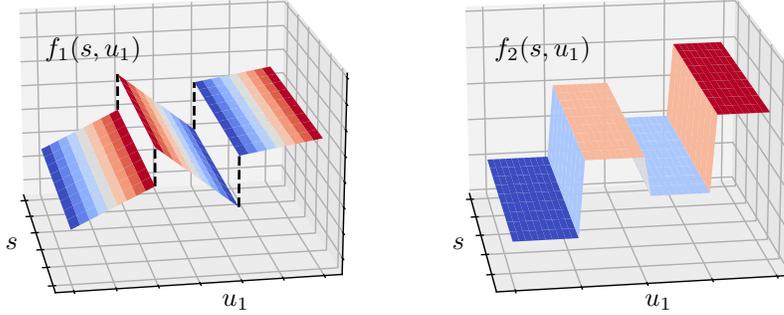
FIG. 2. *Two functions over one region in a BFCP of* $\Lambda_1$. *Bounded, constant-piecewise-linear (CON-PWL) Borel measurable function (left): given* $u_1$, $f_1(s, u_1)$ *is constant in* $s$, *and given* $s$, $f_1(s, u_1)$ *is PWL in* $u_1$. *Constant-piecewise-constant (CON-PWC) stochastic kernel (right): given* $u_1$, $f_2(s, u_1)$ *is constant in* $s$, *and given* $s$, $f_2(s, u_1)$ *is PWC in* $u_1$.

of $\{u_1 \in \mathbb{P}(A_1) \mid (s, u_1) \in \theta(\phi)\}$, and $\theta(\phi) \in \Theta$ again using Definition 7.3 such that, over each region $\phi' \in \Phi'(\phi)$, $f(s, u_1)$ is linear in $u_1$ given $s \in \phi$. Similarly using Definition 7.4, each CON-PWC stochastic kernel $f$ can be represented by a finite set of vectors $\{D_{\phi, \phi'} \in \mathbb{P}(A_2) \mid \phi \in \Phi \wedge \phi' \in \Phi(\phi)\}$ such that $f(s, u_1) = D_{\phi, \phi'}$ for all $s \in \phi$ and $u_1 \in \phi'$, where over each region $\phi' \in \Phi'(\phi)$, $f(s, u_1)$ is constant in $u_1$ given $s \in \phi$.

We introduce a criterion for selecting the maximum or minimum solution over a region, by which the strategies from policy improvement are finitely representable.

DEFINITION 7.5 (CON-1 solution). *Let* $f \in \mathbb{F}(\Lambda_1)$ *be a bounded, CON-PWL Borel measurable function. Using Definition 7.3 there exists a BFCP* $\Phi$ *of $S$ for $f$. Now, for each* $\phi \in \Phi$, *if there exists* $u_1^\phi \in \mathbb{P}(A_1(s))$ *such that:*

$$f(s, u_1^\phi) = \max_{u_1 \in \mathbb{P}(A_1(s))} f(s, u_1)$$

*for all* $s \in \phi$ *and if* $\sigma_1$ *is a strategy of* $\mathsf{Ag}_1$ *such that* $\sigma_1(s) = u_1^\phi$ *for all* $s \in \phi$, *then* $\sigma_1$ *is a* constant-1 (CON-1) solution *of $f$ over $\phi$.*

DEFINITION 7.6 (CON-2 solution). *Let* $f \in \mathbb{F}(\Lambda_{12})$ *be a bounded, Borel measurable function. If there exists a BFCP* $\Theta$ *of* $\Lambda_1$ *where, for each* $\theta \in \Theta$, $A_2(s)$ *is constant for all* $(s, u_1) \in \theta$ *and there exists* $u_2^\theta \in \mathbb{P}(A_2(s))$ *such that:*

$$f(s, u_1, u_2^\theta) = \min_{u_2 \in \mathbb{P}(A_2(s))} f(s, u_1, u_2)$$

*for all* $(s, u_1) \in \theta$, *and if* $\overline{\sigma}_2$ *is a Stackelberg strategy for* $\mathsf{Ag}_2$ *such that* $\overline{\sigma}_2(s, u_1) = u_2^\theta$ *for all* $(s, u_1) \in \theta$, *then* $\overline{\sigma}_2$ *is a* constant-2 (CON-2) solution *of $f$ over $\theta$.*

We now use the operators of Definitions 7.1 and 7.2, together with the functions and solutions from Definitions 7.3 to 7.6 specialised to Borel state spaces with game structures described in Sections 3 and 4, to derive a PI algorithm called *Minimax-action-free PI* (Algorithm 7.1) for strategy synthesis of NS-CSGs with Borel state spaces. Our algorithm closely follows the steps of the PI method of [5] for finite state spaces, but has to resolve a number of issues due to the uncountability of the underlying state space and the need to ensure Borel measurability at each iteration. To overcome these issues we (i) introduce bounded, CON-PWL Borel measurable functions and CON-PWC Borel measurable strategies to ensure measurability and finite representability; (ii) work with CON-1 and CON-2 solutions for policy improvement to ensure that the strategies generated are finitely representable and consistent; and

---

ALGORITHM 7.1 Iteration $t$ of Minimax-action-free PI

---

1: **Input:** NS-CSG C, PWC $\sigma_1^t \in \Sigma_1$, CON-PWC $\overline{\sigma}_2^t \in \overline{\Sigma}_2$, PWC $J_1^t, V_1^t \in \mathbb{F}(S)$,
   CON-PWL $J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$

2: **Perform one of the following four iterations.**

3:     Policy evaluation of $\mathsf{Ag}_1$:

4:         $J_1^{t+1} \leftarrow H^1_{\sigma_1^t, V_2^t}(J_2^t)$ via $PE1$, $\sigma_1^{t+1} \leftarrow \sigma_1^t$,

5:         $V_1^{t+1} \leftarrow V_1^t, \overline{\sigma}_2^{t+1} \leftarrow \overline{\sigma}_2^t, J_2^{t+1} \leftarrow J_2^t, V_2^{t+1} \leftarrow V_2^t$

6:     Policy improvement of $\mathsf{Ag}_1$ by CON-1 solution:

7:         $\sigma_1^{t+1}(s) \in \mathrm{argmax}_{u_1 \in \mathbb{P}(A_1(s))} H^1_{u_1, V_2^t}(J_2^t)(s)$,

8:         $V_1^{t+1} \leftarrow H^1_{\sigma_1^{t+1}, V_2^t}(J_2^t)$ via $PI1$,

9:         $J_1^{t+1} \leftarrow J_1^t, \overline{\sigma}_2^{t+1} \leftarrow \overline{\sigma}_2^t, J_2^{t+1} \leftarrow J_2^t, V_2^{t+1} \leftarrow V_2^t$

10:     Policy evaluation of $\mathsf{Ag}_2$:

11:         $J_2^{t+1} \leftarrow H^2_{\overline{\sigma}_2^t, V_1^t}(J_1^t)$ via $PE2$, $\sigma_1^{t+1} \leftarrow \sigma_1^t$,

12:         $J_1^{t+1} \leftarrow J_1^t, V_1^{t+1} \leftarrow V_1^t, \overline{\sigma}_2^{t+1} \leftarrow \overline{\sigma}_2^t, V_2^{t+1} \leftarrow V_2^t$

13:     Policy improvement of $\mathsf{Ag}_2$ by CON-2 solution:

14:         $\overline{\sigma}_2^{t+1}(s, u_1) \in \mathrm{argmin}_{u_2 \in \mathbb{P}(A_2(s))} H^2_{u_2, V_1^t}(J_1^t)(s, u_1)$,

15:         $V_2^{t+1} \leftarrow H^2_{\overline{\sigma}_2^{t+1}, V_1^t}(J_1^t)$ via $PI2$,

16:         $\sigma_1^{t+1} \leftarrow \sigma_1^t, J_1^{t+1} \leftarrow J_1^t, V_1^{t+1} \leftarrow V_1^t, J_2^{t+1} \leftarrow J_2^t$

17: $t \leftarrow t + 1$

---

(iii) propose a BFCP iteration algorithm (Algorithm 7.2) and a BFCP-based computation algorithm (Algorithm 7.3) to compute a new BFCP of the state space and the values or strategies over this BFCP. We also provide a simpler proof than that presented in [5], which does not require the introduction of any new concepts except those used in the algorithm.

The Minimax-action-free PI algorithm is initialized with a strategy for each agent, $\sigma_1^0$ and $\overline{\sigma}_2^0$, and four functions, $J_1^0, V_1^0, J_2^0, V_2^0$, where Algorithm 7.2 gives one BFCP for each initial strategy and function. An iteration of the Minimax-action-free PI is given in Algorithm 7.1. As shown later, the order and frequency by which the possible four iterations of Algorithm 7.1 are run do not affect the convergence, as long as each of them is performed infinitely often. Although we have not discussed the pros and cons of an asynchronous implementation, the Minimax-action-free PI algorithm indeed allows such an implementation, as discussed in [5] and for its single-agent counterparts in [6].

For each of the four iterations, Algorithm 7.2 provides a way to compute new BFCPs and the results below demonstrate that, over each region of these BFCPs, the corresponding computed strategies and value functions are either constant, PWC or PWL. Therefore, we can follow similar steps to our VI algorithm (see Algorithm 6.1) to compute the value functions of these new strategies and value functions (see Algorithm 7.3). The idea is to first compute the BFCPs $\Phi_{J_1^{t+1}}$, $\Phi_{V_1^{t+1}}$, $\Phi_{\sigma_1^{t+1}}$, $\Theta_{J_2^{t+1}}$, $\Theta_{V_2^{t+1}}$ and $\Theta_{\overline{\sigma}_2^{t+1}}$ via Algorithm 7.2 and then use these BFCPs to compute strategies and value functions using Algorithm 7.3. For instance, if the policy improvement of $\mathsf{Ag}_2$ is chosen at iteration $t$ then we proceed as follows. First, new BFCPs are computed via Algorithm 7.2. Second, $PI2$ of Algorithm 7.3 is performed. In this second step we take each region $\theta \in \Theta_{\overline{\sigma}_2^{t+1}}$, let $\phi = \{s \mid (s, u_1) \in \theta\}$, then take one state $s' \in \phi$, and compute a BFCP $\Phi_u$ of $\mathbb{P}(A_1(s'))$ such that $\min_{u_2 \in \mathbb{P}(A_2(s'))} H^2_{u_2, V_1^t}(J_1^t)(s', u_1)$ is con-

ALGORITHM 7.2 BFCP iteration $t$ for PI

1: **Input:** Perception FCP $\Phi_P$, reward FCPs $\langle \Phi_R^\alpha \rangle_{\alpha \in A}$
2: **Output:** BFCPs $\langle \Phi_{J_1^t}, \Phi_{V_1^t}, \Phi_{\sigma_1^t}, \Theta_{J_2^t}, \Theta_{V_2^t}, \Theta_{\overline{\sigma}_2^t} \rangle_{t \in \mathbb{N}}$ for $\langle J_1^t, V_1^t, \sigma_1^t, J_2^t, V_2^t, \overline{\sigma}_2^t \rangle_{t \in \mathbb{N}}$
3: $\Phi_{J_1^0}, \Phi_{V_1^0}, \Phi_{\sigma_1^0} \leftarrow \{S\}, \ \Theta_{J_2^0}, \Theta_{V_2^0}, \Theta_{\overline{\sigma}_2^0} \leftarrow \{\Lambda_1\}$
4: **while** Algorithm 7.1 performs iteration $t$ **do**
5:     **if** policy evaluation of $\mathsf{Ag}_1$ is chosen **then**
6:         $Preprocess\_maximiser()$,
7:         $\Phi_{J_1^{t+1}} \leftarrow \Phi_{\sigma_1^t} + \Phi_{J_2^t} + \Phi_{V_2^t}, \ \ \Phi_{V_1^{t+1}} \leftarrow \Phi_{V_1^t}, \ \ \Phi_{\sigma_1^{t+1}} \leftarrow \Phi_{\sigma_1^t}$
8:     **if** policy improvement of $\mathsf{Ag}_1$ is chosen **then**
9:         $Preprocess\_maximiser()$,
10:         $\Phi_{\sigma_1^{t+1}} \leftarrow \Phi_{J_2^t} + \Phi_{V_2^t}, \ \ \Phi_{V_1^{t+1}} \leftarrow \Phi_{\sigma_1^{t+1}}, \ \ \Phi_{J_1^{t+1}} \leftarrow \Phi_{J_1^t}$
11:     **if** policy evaluation of $\mathsf{Ag}_2$ is chosen **then**
12:         $Preprocess\_minimiser()$,
13:         $\Phi_{\overline{\sigma}_2^t} \leftarrow \big\{ \{s \mid (s, u_1) \in \theta\} \mid \theta \in \Theta_{\overline{\sigma}_2^t} \big\}$,
14:         $\Theta_{J_2^{t+1}} \leftarrow \big\{ \{(s, u_1) \in \Lambda_1 \mid s \in \phi\} \mid \phi \in \Phi_{\hat{Q}^{t+1}} + \Phi_{\overline{\sigma}_2^t} \big\}$,
15:         $\Theta_{V_2^{t+1}} \leftarrow \Theta_{V_2^t}, \ \ \Theta_{\overline{\sigma}_2^{t+1}} \leftarrow \Theta_{\overline{\sigma}_2^t}$
16:     **if** policy improvement of $\mathsf{Ag}_2$ is chosen **then**
17:         $Preprocess\_minimiser()$,
18:         $\Theta_{\overline{\sigma}_2^{t+1}} \leftarrow \big\{ \{(s, u_1) \in \Lambda_1 \mid s \in \phi\} \mid \phi \in \Phi_{\hat{Q}^{t+1}} \big\}$,
19:         $\Theta_{V_2^{t+1}} \leftarrow \Theta_{\overline{\sigma}_2^{t+1}}, \ \ \Theta_{J_2^{t+1}} \leftarrow \Theta_{J_2^t}$
20: **return** $\langle \Phi_{J_1^t}, \Phi_{V_1^t}, \Phi_{\sigma_1^t}, \Theta_{J_2^t}, \Theta_{V_2^t}, \Theta_{\overline{\sigma}_2^t} \rangle_{t \in \mathbb{N}}$
21:
22: **procedure** $Preprocess\_maximiser()$
23:     $\Theta_{J_2^{t+1}} \leftarrow \Theta_{J_2^t}, \ \ \Theta_{V_2^{t+1}} \leftarrow \Theta_{V_2^t}, \ \ \Theta_{\overline{\sigma}_2^{t+1}} \leftarrow \Theta_{\overline{\sigma}_2^t}$,
24:     $\Phi_{J_2^t} \leftarrow \big\{ \{s \mid (s, u_1) \in \theta\} \mid \theta \in \Theta_{J_2^t} \big\}, \ \ \Phi_{V_2^t} \leftarrow \big\{ \{s \mid (s, u_1) \in \theta\} \mid \theta \in \Theta_{V_2^t} \big\}$
25: **procedure** $Preprocess\_minimiser()$
26:     $\Phi_{J_1^{t+1}} \leftarrow \Phi_{J_1^t}, \ \ \Phi_{V_1^{t+1}} \leftarrow \Phi_{V_1^t}, \ \ \Phi_{\sigma_1^{t+1}} \leftarrow \Phi_{\sigma_1^t}$,
27:     $\Phi_{\hat{Q}^{t+1}} \leftarrow Preimage\_BFCP(\Phi_{J_1^t} + \Phi_{V_1^t}, \Phi_P, \langle \Phi_R^\alpha \rangle_{\alpha \in A})$

stant over each region $\phi_u \in \Phi_u$ and for all $u_1 \in \phi_u$. Third, take one $u_1' \in \phi_u$ and find $u_2' \in \mathbb{P}(A_2(s'))$ that minimises $H^2_{u_2, V_1^t}(J_1^t)(s', u_1')$. Fourth, we let $\overline{\sigma}_2^{t+1}(s, u_1) = u_2'$ for all $s \in \phi$ and $u_1 \in \phi_u$, which is a CON-2 solution of $H^2_{u_2, V_1^t}(J_1^t)(s, u_1)$ over $\{(s, u_1) \mid s \in \phi, u_1 \in \phi_u\}$ by Lemma 7.10 and $V_2^{t+1}(s, u_1)$ is CON-linear in $s \in \phi$ and $u_1 \in \phi_u$. Finally, we copy the other strategies and value functions for the next iteration.

We next introduce the following lemmas in order to prove that the strategies and value functions generated during each iteration of the Minimax-action-free PI algorithm are finitely representable.

LEMMA 7.7 (Evaluation consistency for $\mathsf{Ag}_1$). *If $\sigma_1^t \in \Sigma_1$ is a PWC stochastic kernel and $J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$ are bounded, CON-PWL Borel measurable and an iteration of policy evaluation of $\mathsf{Ag}_1$ is performed (procedure PE1), then $J_1^{t+1} = H^1_{\sigma_1^t, V_2^t}(J_2^t)$ is bounded, PWC Borel measurable.*

*Proof.* Suppose $\sigma_1^t \in \Sigma_1$ is a PWC stochastic kernel and $J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$ are bounded, CON-PWL Borel measurable. Since $\sigma_1^t$ is a PWC stochastic kernel, there

---

ALGORITHM 7.3 BFCP based computation for PI

---

1: **Input:** $J_1^t, V_1^t, \sigma_1^t, J_2^t, V_2^t, \overline{\sigma}_2^t, \Phi_{J_1^{t+1}}, \Phi_{\sigma_1^{t+1}}, \Theta_{J_2^{t+1}}, \Theta_{\overline{\sigma}_2^{t+1}}$

2: **procedure** *PE1*

3:     **for** $\phi \in \Phi_{J_1^{t+1}}$ **do**

4:         Take one state $s \in \phi$, and then $J_{1,\phi}^{t+1} \leftarrow H_{\sigma_1^t, V_2^t}^1(J_2^t)(s)$

5:     **return** $J_1^{t+1} \leftarrow \langle J_{1,\phi}^{t+1} \rangle_{\phi \in \Phi_{J_1^{t+1}}}$

6: **procedure** *PI1*

7:     **for** $\phi \in \Phi_{\sigma_1^{t+1}}$ **do**

8:         Take one state $s \in \phi$, and then $u_1 \in \operatorname{argmax}_{u_1 \in \mathbb{P}(A_1(s))} H_{u_1, V_2^t}^1(J_2^t)(s)$

9:         $\sigma_{1,\phi}^{t+1} \leftarrow u_1, \ \ V_{1,\phi}^{t+1} \leftarrow \max_{u_1 \in \mathbb{P}(A_1(s))} H_{u_1, V_2^t}^1(J_2^t)(s)$

10:     **return** $\sigma_1^{t+1} \leftarrow \langle \sigma_{1,\phi}^{t+1} \rangle_{\phi \in \Phi_{\sigma_1^{t+1}}}, \ \ V_1^{t+1} \leftarrow \langle V_{1,\phi}^{t+1} \rangle_{\phi \in \Phi_{\sigma_1^{t+1}}}$

11: **procedure** *PE2*

12:     **for** $\theta \in \Theta_{J_2^{t+1}}$ **do**

13:         $\phi \leftarrow \{s \mid (s, u_1) \in \theta\}$

14:         Take one state $s \in \phi$, and then compute a BFCP $\Phi_u$ of $\mathbb{P}(A_1(s))$ such that over $\phi_u \in \Phi_u$, $H_{\overline{\sigma}_2^t, V_1^t}^2(J_1^t)(s, u_1)$ is linear in $u_1$

15:         $J_{2,\phi,\phi_u}^{t+1} \leftarrow H_{\overline{\sigma}_2^t, V_1^t}^2(J_1^t)(s, u_1)$ is linear in $u_1$

16:     **return** $J_2^{t+1} \leftarrow \langle J_{2,\phi,\phi_u}^{t+1} \rangle_{\theta \in \Theta_{J_2^{t+1}}}$

17: **procedure** *PI2*

18:     **for** $\theta \in \Theta_{\overline{\sigma}_2^{t+1}}$ **do**

19:         $\phi \leftarrow \{s \mid (s, u_1) \in \theta\}$

20:         Take one state $s' \in \phi$, and then compute a BFCP $\Phi_u$ of $\mathbb{P}(A_1(s'))$ such that over $\phi_u \in \Phi_u$, $\min_{u_2 \in \mathbb{P}(A_2(s'))} H_{u_2, V_1^t}^2(J_1^t)(s', u_1)$ is constant for all $u_1 \in \phi_u$

21:         Take $u_1' \in \phi_u$, and $u_2' \in \operatorname{argmin}_{u_2 \in \mathbb{P}(A_2(s'))} H_{u_2, V_1^t}^2(J_1^t)(s', u_1')$ for $\phi_u \in \Phi_u$

22:         $\overline{\sigma}_{2,\phi,\phi_u}^{t+1} \leftarrow u_2', \ \ V_{2,\phi,\phi_u}^{t+1} \leftarrow H_{u_2', V_1^t}^2(J_1^t)(s', u_1)$ is linear in $u_1$

23:     **return** $\overline{\sigma}_2^{t+1} \leftarrow \langle \overline{\sigma}_{2,\phi,\phi_u}^{t+1} \rangle_{\theta \in \Theta_{\overline{\sigma}_2^{t+1}}}, \ \ V_2^{t+1} \leftarrow \langle V_{2,\phi,\phi_u}^{t+1} \rangle_{\theta \in \Theta_{\overline{\sigma}_2^{t+1}}}$

---

exists a constant-BFCP $\Phi_{\sigma_1^t}$ of $S$ for $\sigma_1^t$. Since $J_2^t$ is a CON-PWL Borel measurable function, there exists a BFCP $\Phi_{J_2^t}$ of $S$ satisfying the properties of Definition 7.3 for $J_2^t$. Therefore $J_2^t(s, \sigma_1^t(s))$ is constant on each region in the BFCP $\Phi_{\sigma_1^t} + \Phi_{J_2^t}$. We can similarly show that $V_2^t(s, \sigma_1^t(s))$ is constant on each region in the BFCP $\Phi_{\sigma_1^t} + \Phi_{V_2^t}$, where $\Phi_{V_2^t}$ is a BFCP of $S$ from Definition 7.3 for $V_2^t$. Consider the policy evaluation of $\mathsf{Ag}_1$ (procedure *PE1*). Using Definition 7.1 we have that $J_1^{t+1} = H_{\sigma_1^t, V_2^t}^1(J_2^t)$ is constant on each region in the BFCP $\Phi_{\sigma_1^t} + \Phi_{J_2^t} + \Phi_{V_2^t}$, which also implies that $J_1^{t+1}$ is Borel measurable. Since $J_2^t$ and $V_2^t$ are bounded, then $J_1^{t+1}$ is also bounded.     □

LEMMA 7.8 (Improvement consistency for $\mathsf{Ag}_1$). *If $J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$ are bounded, CON-PWL Borel measurable and an iteration of policy improvement of $\mathsf{Ag}_1$ is performed (procedure PI1), then $\sigma_1^{t+1}(s) \in \operatorname{argmax}_{u_1 \in \mathbb{P}(A_1(s))} H_{u_1, V_2^t}^1(J_2^t)(s)$ is a PWC stochastic kernel, and $V_1^{t+1} = H_{\sigma_1^{t+1}, V_2^t}^1(J_2^t)$ is bounded, PWC Borel measurable.*

*Proof.* Suppose $J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$ are bounded, CON-PWL Borel measurable functions. Using [27, Chapter 18.1] and Definition 7.3 it follows that the function $K^t \in$

$\mathbb{F}(\Lambda_1)$ where $K^t(s, u_1) \coloneqq \min\{J_2^t(s, u_1), V_2^t(s, u_1)\}$ is also a bounded, Borel measurable function for $(s, u_1) \in \Lambda_1$. Note that, over each region in $\Phi_{J_2^t} + \Phi_{V_2^t}$, $K^t(s, u_1)$ is constant in $s$ given $u_1$, and PWL in $u_1$ given $s$ ($\Phi_{J_2^t}$ and $\Phi_{V_2^t}$ from Lemma 7.7), and therefore $K^t$ is CON-PWL.

Let $\Phi_{K^t} = \Phi_{J_2^t} + \Phi_{V_2^t}$ be a BFCP of $S$ and $\Theta_{K^t}$ a BFCP of $\Lambda_1$ satisfying the properties of Definition 7.3 for $K^t$. Every state in each region of the BFCP $\Phi_{K^t}$ has the same set of available actions for $\mathsf{Ag}_1$ and same strategy $u_1$ that maximises $K^t(s, u_1)$ on a region of $\Theta_{K^t}$. Therefore, using the CON-1 solution in Definition 7.5, the strategy of $\mathsf{Ag}_1$:

$$\sigma_1^{t+1}(s) \in \mathrm{argmax}_{u_1 \in \mathbb{P}(A_1(s))}\, H_{u_1, V_2^t}^1(J_2^t)(s)$$

is constant on each region in $\Phi_{K^t}$, which also implies that $\sigma_1^{t+1}$ is Borel measurable.

Since $\sigma_1^{t+1}$ is a PWC stochastic kernel, then Lemma 7.7 implies that $V_1^{t+1}$ is bounded, PWC Borel measurable. $\qquad\square$

LEMMA 7.9 (Evaluation consistency for $\mathsf{Ag}_2$). *If $J_1^t, V_1^t \in \mathbb{F}(S)$ are bounded, PWC Borel measurable and $\overline{\sigma}_2^t \in \overline{\Sigma}_2$ is a CON-PWC stochastic kernel and an iteration of policy evaluation of $\mathsf{Ag}_2$ is performed (procedure PE1 ), then $J_2^{t+1} = H_{\overline{\sigma}_2^t, V_1^t}^2(J_1^t)$ is bounded, CON-PWL Borel measurable.*

*Proof.* The proof follows similarly to Lemma 7.7. $\qquad\square$

LEMMA 7.10 (Improvement consistency for $\mathsf{Ag}_2$). *If $J_1^t, V_1^t \in \mathbb{F}(S)$ are bounded, PWC Borel measurable and an iteration of policy improvement of $\mathsf{Ag}_2$ is performed (procedure PI2 ), then $\overline{\sigma}_2^{t+1}(s, u_1) \in \mathrm{argmin}_{u_2 \in \mathbb{P}(A_2(s))} H_{u_2, V_1^t}^2(J_1^t)(s, u_1)$ is a CON-PWC stochastic kernel, and $V_2^{t+1} = H_{\overline{\sigma}_2^{t+1}, V_1^t}^2(J_1^t)$ is bounded, CON-PWL Borel measurable.*

*Proof.* The proof follows similarly to Lemma 7.8. $\qquad\square$

By fusing Lemmas 7.7 to 7.10 we can prove that the strategies and value functions generated during each iteration of Algorithm 7.1 never leave a finitely representable class of functions. Formally, we have the following result.

THEOREM 7.11 (Representation consistency). *In any iteration of the Minimax-action-free PI algorithm (see Algorithm 7.1), if*
  (i) *$J_1^t, V_1^t \in \mathbb{F}(S)$ are bounded, PWC Borel measurable and $\sigma_1^t \in \Sigma_1$ is a PWC stochastic kernel;*
  (ii) *$J_2^t, V_2^t \in \mathbb{F}(\Lambda_1)$ are bounded, CON-PWL Borel measurable and $\overline{\sigma}_2^t \in \overline{\Sigma}_2$ is a CON-PWC stochastic kernel;*
*then so are $J_1^{t+1}$, $V_1^{t+1}$, $\sigma_1^{t+1}$, $J_2^{t+1}$, $V_2^{t+1}$ and $\overline{\sigma}_2^{t+1}$ respectively, regardless of which one of the four iterations is performed.*

*Proof.* The theorem directly follows from Lemmas 7.7 to 7.10. $\qquad\square$

We next show that the optimization problem in the policy improvement of $\mathsf{Ag}_2$ can be further simplified.

COROLLARY 7.12 (Pure strategy for $\mathsf{Ag}_2$). *For the policy improvement of $\mathsf{Ag}_2$ (procedure PI2 ), there exists a CON-PWC stochastic kernel:*

$$\overline{\sigma}_2^{t+1}(s, u_1) \in \mathrm{argmin}_{u_2 \in \mathbb{P}(A_2(s))} H_{u_2, V_1^t}^2(J_1^t)(s, u_1)$$

*such that $\overline{\sigma}_2^{t+1}(a_2 \mid (s, u_1)) = 1$ for some $a_2 \in A_2$ over each region $\theta \in \Theta_{\overline{\sigma}_2^{t+1}}$, where $\Theta_{\overline{\sigma}_2^{t+1}}$ is a BFCP of $\Lambda_1$ for $\overline{\sigma}_2^{t+1}$ using Definition 7.4.*

*Proof.* The result follows directly follows from the proof of Lemma 7.10. ☐

We next demonstrate that Algorithm 7.2 presents a way of constructing new BFCPs such that the strategies and value functions after one iteration of the Minimax-action-free PI algorithm remain constant, PWC, or PWL on each region of the constructed BFCPs.

COROLLARY 7.13 (BFCP iteration for PI). *After performing Algorithm 7.2 we have:*

(i) $\Phi_{J_1^{t+1}}$, $\Phi_{V_1^{t+1}}$ *and* $\Phi_{\sigma_1^{t+1}}$ *are constant-BFCPs of $S$ for $J_1^{t+1} = H_{\sigma_1^t, V_2^t}^1(J_2^t)$, $V_1^{t+1} = H_{\sigma_1^{t+1}, V_2^t}^1(J_2^t)$ and $\sigma_1^{t+1}(s) \in \mathrm{argmax}_{u_1 \in \mathbb{P}(A_1(s))} H_{u_1, V_2^t}^1(J_2^t)(s)$, respectively;*

(ii) $\Theta_{J_2^{t+1}}$ *and* $\Theta_{V_2^{t+1}}$ *are BFCPs of $\Lambda_1$ for $J_2^{t+1} = H_{\overline{\sigma}_2^t, V_1^t}^2(J_1^t)$ and $V_2^{t+1} = H_{\overline{\sigma}_2^{t+1}, V_1^t}^2(J_1^t)$ meeting the conditions of Definition 7.3 respectively, and $\Theta_{\overline{\sigma}_2^{t+1}}$ is a BFCP of $\Lambda_1$ for $\overline{\sigma}_2^{t+1}(s, u_1) \in \mathrm{argmin}_{u_2 \in \mathbb{P}(A_2(s))} H_{u_2, V_1^t}^2(J_1^t)(s, u_1)$ meeting the conditions of Definition 7.4.*

*Proof.* The results follow directly from Lemmas 6.1 and 7.7 to 7.10. ☐

**7.2. Convergence analysis and strategy computation.** We next prove the convergence of the Minimax-action-free PI algorithm by showing that an operator related to the algorithm is a contraction mapping with a unique fixed point, one of whose components is the value function multiplied by a known constant. The proof closely follows the steps for finite state spaces given in [5], but is more complex due to the underlying infinite state space and the need to deal with the requirement of Borel measurable consistency and finite representation of strategies and value functions.

Given PWC $\sigma_1 \in \Sigma_1$ and CON-PWC $\overline{\sigma}_2 \in \overline{\Sigma}_2$ stochastic kernels, we define the operator $G_{\sigma_1, \overline{\sigma}_2} : (\mathbb{F}(S) \times \mathbb{F}(S) \times \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1)) \to (\mathbb{F}(S) \times \mathbb{F}(S) \times \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1))$ such that:

$$(7.1)\quad G_{\sigma_1, \overline{\sigma}_2}(J_1, V_1, J_2, V_2) \coloneqq (M_{\sigma_1}^1(J_2, V_2), K^1(J_2, V_2), M_{\overline{\sigma}_2}^2(J_1, V_1), K^2(J_1, V_1))$$

where $J_1, V_1 \in \mathbb{F}(S)$ are assumed to be PWC, $J_2, V_2 \in \mathbb{F}(\Lambda_1)$ are assumed to be CON-PWL, and the four operators $M_{\sigma_1}^1$, $K^1$, $M_{\overline{\sigma}_2}^2$ and $K^2$ represent the four iterations of the Minimax-action-free PI algorithm from lines 2 to 14, and are defined as follows.

- $M_{\sigma_1}^1 : \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1) \to \mathbb{F}(S)$ corresponds to the policy evaluation of $\mathsf{Ag}_1$ (procedure *PE1*) where for any $s \in S$:

$$(7.2)\qquad\qquad M_{\sigma_1}^1(J_2, V_2)(s) \coloneqq H_{\sigma_1, V_2}^1(J_2)(s)$$

  and is bounded, PWC Borel measurable using Lemma 7.7.

- $K^1 : \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1) \to \mathbb{F}(S)$ corresponds to the policy improvement of $\mathsf{Ag}_1$ (procedure *PI1*) and for any $s \in S$:

$$(7.3)\qquad\quad K^1(J_2, V_2)(s) \coloneqq \max_{u_1 \in \mathbb{P}(A_1(s))} H_{u_1, V_2}^1(J_2)(s)$$

  and is bounded, PWC Borel measurable using Lemma 7.8.

- $M_{\overline{\sigma}_2}^2 : \mathbb{F}(S) \times \mathbb{F}(S) \to \mathbb{F}(\Lambda_1)$ corresponds to the policy evaluation of $\mathsf{Ag}_2$ (procedure *PE2*) and for any $(s, u_1) \in \Lambda_1$:

$$(7.4)\qquad\qquad M_{\overline{\sigma}_2}^2(J_1, V_1)(s, u_1) \coloneqq H_{\overline{\sigma}_2, V_1}^2(J_1)(s, u_1)$$

  and is bounded, CON-PWL Borel measurable using Lemma 7.9.

- $K^2 : \mathbb{F}(S) \times \mathbb{F}(S) \to \mathbb{F}(\Lambda_1)$ corresponds to the policy improvement of $\mathsf{Ag}_2$ (procedure *PI2*) and any $(s, u_1) \in \Lambda_1$:

(7.5)
$$K^2(J_1, V_1)(s, u_1) := \min_{u_2 \in \mathbb{P}(A_2(s))} H^2_{u_2, V_1}(J_1)(s, u_1)$$

and is bounded, CON-PWL Borel measurable using Lemma 7.10.

For both the spaces $\mathbb{F}(S) \times \mathbb{F}(S)$ and $\mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1)$, we consider the norm:

$$\|(J, V)\| = \max\{\|J\|, \|V\|\}$$

and for the space $\mathbb{F}(S) \times \mathbb{F}(S) \times \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1)$ the norm:

(7.6)
$$\|(J_1, V_1, J_2, V_2)\| = \max\{\|J_1\|, \|V_1\|, \|J_2\|, \|V_2\|\}.$$

We next require the following properties of theses norms.

LEMMA 7.14. *For any $J_1, V_1, J_1', V_1' \in \mathbb{F}(S)$ and $J_2, V_2, J_2', V_2' \in \mathbb{F}(\Lambda_1)$, we have:*

$$\|\max[J_1, V_1] - \max[J_1', V_1']\| \leq \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}$$
$$\|\min[J_2, V_2] - \min[J_2', V_2']\| \leq \max\{\|J_2 - J_2'\|, \|V_2 - V_2'\|\}.$$

*Proof.* Consider any $J_1, V_1, J_1', V_1' \in \mathbb{F}(S)$. The sup-norm for the space $\mathbb{F}(S)$ implies that for every $s \in S$:

(7.7)
$$J_1(s) \leq J_1'(s) + \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}$$
(7.8)
$$V_1(s) \leq V_1'(s) + \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}$$

from which we have:

(7.9)
$$\max\{J_1(s), V_1(s)\} \leq \max\{J_1'(s), V_1'(s)\} + \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}.$$

Exchanging $(J_1, V_1)$ with $(J_1', V_1')$ in (7.7) and (7.8) derives an inequality similar to (7.9), and combining it with (7.9) leads to the inequality:

$$|\max\{J_1(s), V_1(s)\} - \max\{J_1'(s), V_1'(s)\}| \leq \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}$$

implying that:

(7.10) $$\sup_{s \in S} |\max\{J_1(s), V_1(s)\} - \max\{J_1'(s), V_1'(s)\}| \leq \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}.$$

Since $J_1, V_1, J_1'$ and $V_1'$ are bounded, Borel measurable, so is $\max[J_1, V_1] - \max[J_1', V_1']$ by [27, Chapter 18.1], i.e., $\max[J_1, V_1] - \max[J_1', V_1'] \in \mathbb{F}(S)$. Thus, by (7.10):

$$\|\max[J_1, V_1] - \max[J_1', V_1']\| \leq \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\}.$$

The second inequality of the lemma can be proved following the same steps for $J_2, V_2, J_2', V_2' \in \mathbb{F}(\Lambda_1)$. □

Using the above operators and results, we are now in position to prove the convergence of the Minimax-action-free PI algorithm.

THEOREM 7.15 (Convergence guarantee). *Suppose that each of the four iterations of the Minimax-action-free PI algorithm from lines 2 to 14 is performed infinitely often. Then the sequence $\langle \gamma V_1^t \rangle_{t=0}^{\infty}$ generated by the algorithm converges to $V^\star$.*

*Proof.* We prove that each component of the operator $G_{\sigma_1,\overline{\sigma}_2}$ satisfies a contraction property. Suppose that $J_1, V_1, J_1', V_1' \in \mathbb{F}(S)$ are PWC and $J_2, V_2, J_2', V_2' \in \mathbb{F}(\Lambda_1)$ are CON-PWL.

- For $M_{\sigma_1}^1$, since $M_{\sigma_1}^1(J_2, V_2) - M_{\sigma_1}^1(J_2', V_2') \in \mathbb{F}(S)$ by [27, Chapter 18.1] and by Definition 7.1 and the definition of the sup-norm for $\mathbb{F}(S)$ we have:

$$(7.11) \quad \|M_{\sigma_1}^1(J_2, V_2) - M_{\sigma_1}^1(J_2', V_2')\|$$

$$= \sup_{s \in S} \left| \frac{1}{\gamma} \min\{J_2(s, \sigma_1(s)), V_2(s, \sigma_1(s))\} - \frac{1}{\gamma} \min\{J_2'(s, \sigma_1(s)), V_2'(s, \sigma_1(s))\} \right|$$

$$\leq \frac{1}{\gamma} \sup_{(s, u_1) \in \Lambda_1} \left| \min\{J_2(s, u_1), V_2(s, u_1)\} - \min\{J_2'(s, u_1), V_2'(s, u_1)\} \right|$$

$$\text{rearranging and since } \{(s, \sigma_1(s)) \mid s \in S\} \subseteq \Lambda_1$$

$$= \frac{1}{\gamma} \left\| \min[J_2, V_2] - \min[J_2', V_2'] \right\|$$

$$\text{since } \min[J_2, V_2] - \min[J_2', V_2'] \in \mathbb{F}(\Lambda_1) \text{ using } [27, \text{ Chapter 18.1}]$$

$$\leq \frac{1}{\gamma} \max\{\|J_2 - J_2'\|, \|V_2 - V_2'\|\} \qquad\qquad \text{by Lemma 7.14}$$

$$\leq \frac{1}{\gamma} \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|, \|J_2 - J_2'\|, \|V_2 - V_2'\|\}.$$

- For $K^1$, since $K^1(J_2, V_2) - K^1(J_2', V_2') \in \mathbb{F}(S)$ by Definition 7.1 and definition of the sup-norm for $\mathbb{F}(S)$:

$$(7.12) \quad \|K^1(J_2, V_2) - K^1(J_2', V_2')\|$$

$$= \sup_{s \in S} \left| \max_{u_1 \in \mathbb{P}(A_1(s))} \frac{1}{\gamma} \min\{J_2(s, u_1), V_2(s, u_1)\} \right.$$

$$\left. - \max_{u_1 \in \mathbb{P}(A_1(s))} \frac{1}{\gamma} \min\{J_2'(s, u_1), V_2'(s, u_1)\} \right|$$

$$\leq \frac{1}{\gamma} \sup_{(s, u_1) \in \Lambda_1} \left| \min\{J_2(s, u_1), V_2(s, u_1)\} - \min\{J_2'(s, u_1), V_2'(s, u_1)\} \right|$$

$$\text{rearranging and since } \{(s, u_1) \mid u_1 \in \mathbb{P}(A_1(s))\} \subseteq \Lambda_1$$

$$\leq \frac{1}{\gamma} \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|, \|J_2 - J_2'\|, \|V_2 - V_2'\|\}$$

where the final inequality follows from similar arguments used in (7.11).

- For $M_{\overline{\sigma}_2}^2$, since $M_{\overline{\sigma}_2}^2(J_1, V_1) - M_{\overline{\sigma}_2}^2(J_1', V_1') \in \mathbb{F}(\Lambda_1)$ by Definition 7.2 and the sup-norm for $\mathbb{F}(\Lambda_1)$ we have:

$$(7.13) \quad \|M_{\overline{\sigma}_2}^2(J_1, V_1) - M_{\overline{\sigma}_2}^2(J_1', V_1')\|$$

$$= \sup_{(s, u_1) \in \Lambda_1} \left| \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} \left( Q(s, (a_1, a_2), \right.\right.$$

$$\left.\left. \gamma \max[J_1, V_1]) - Q(s, (a_1, a_2), \gamma \max[J_1', V_1']) \right) u_1(a_1) \overline{\sigma}_2(a_2 \mid (s, u_1)) \right|$$

$$= \sup_{(s, u_1) \in \Lambda_1} \left| \sum_{(a_1, a_2) \in A(s)} \gamma\beta \sum_{s' \in \Theta(s, (a_1, a_2))} \delta(s, (a_1, a_2))(s') \right.$$

$$\left. \left(\max\{J_1(s'), V_1(s')\} - \max\{J_1'(s'), V_1'(s')\}\right)u_1(a_1)\overline{\sigma}_2(a_2 \mid (s, u_1))\right|$$

rearranging, by Definition 5.4 and the sup-norm for $\mathbb{F}(\Lambda_1)$

$$\leq \gamma\beta \sup_{(s,u_1)\in\Lambda_1} \sum_{(a_1,a_2)\in A(s)} \sum_{s'\in\Theta(s,(a_1,a_2))} \delta(s,(a_1,a_2))(s')$$

$$\left|\max\{J_1(s'), V_1(s')\} - \max\{J_1'(s'), V_1'(s')\}\right| u_1(a_1)\overline{\sigma}_2(a_2 \mid (s, u_1))$$

rearranging and since $\delta$, $u_1$ and $\overline{\sigma}_2$ are non-negative

$$\leq \gamma\beta \sup_{(s,u_1)\in\Lambda_1} \sum_{(a_1,a_2)\in A(s)} \sum_{s'\in\Theta(s,(a_1,a_2)} \delta(s,(a_1,a_2))(s')$$

$$\sup_{s''\in S}\left|\max\{J_1(s''), V_1(s'')\} - \max\{J_1'(s''), V_1'(s'')\}\right| u_1(a_1)\overline{\sigma}_2(a_2 \mid (s, u_1))$$

since $f(s') \leq \sup_{s''\in S} f(s'')$ for any $f \in \mathbb{F}(S)$

$$= \gamma\beta \sup_{s''\in S}\left|\max\{J_1(s''), V_1(s'')\} - \max\{J_1'(s''), V_1'(s'')\}\right|$$

$$\sup_{(s,u_1)\in\Lambda_1} \sum_{(a_1,a_2)\in A(s)} \sum_{s'\in\Theta(s,(a_1,a_2))} \delta(s,(a_1,a_2))(s')u_1(a_1)\overline{\sigma}_2(a_2 \mid (s, u_1))$$

rearranging

$$= \gamma\beta \sup_{s''\in S}\left|\max\{J_1(s''), V_1(s'')\} - \max\{J_1'(s''), V_1'(s'')\}\right|$$

since $\delta \in \mathbb{P}(S \times A)$, $u_1 \in \mathbb{P}(A_1)$ and $\overline{\sigma}_2 \in \mathbb{P}(A_2 \mid \Lambda_1)$

$$= \gamma\beta \left\|\max[J_1, V_1] - \max[J_1', V_1']\right\|$$

since $\max[J_1, V_1] - \max[J_1', V_1'] \in \mathbb{F}(S)$

$$\leq \gamma\beta \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|\} \qquad \text{by Lemma 7.14}$$

$$\leq \gamma\beta \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|, \|J_2 - J_2'\|, \|V_2 - V_2'\|\}$$

- For $K^2$, since $K^2(J_1, V_1) - K^2(J_1', V_1') \in \mathbb{F}(\Lambda_1)$, by definition of the sup-norm for $\mathbb{F}(\Lambda_1)$ we have:

$$(7.14) \quad \|K^2(J_1, V_1) - K^2(J_1', V_1')\|$$

$$= \sup_{(s,u_1)\in\Lambda_1}\left| \min_{u_2\in\mathbb{P}(A_2(s))} \sum_{a_1\in A_1(s)}\sum_{a_2\in A_2(s)} \gamma\beta \sum_{s'\in\Theta(s,(a_1,a_2))} \delta(s,(a_1,a_2))(s') \right.$$

$$\left. \left(\max\{J_1(s'), V_1(s')\} - \max\{J_1'(s'), V_1'(s')\}\right) u_1(a_1)u_2(a_2)\right|$$

$$\leq \gamma\beta \sup_{(s,u_1)\in\Lambda_1} \min_{u_2\in\mathbb{P}(A_2(s))} \sum_{(a_1,a_2)\in A(s)}\sum_{s'\in\Theta(s,(a_1,a_2))} \delta(s,(a_1,a_2))(s')$$

$$\left|\max\{J_1(s'), V_1(s')\} - \max\{J_1'(s'), V_1'(s')\}\right| u_1(a_1)u_2(a_2) \quad \text{rearranging}$$

$$\leq \gamma\beta \max\{\|J_1 - J_1'\|, \|V_1 - V_1'\|, \|J_2 - J_2'\|, \|V_2 - V_2'\|\}$$

where the final inequality follows from similar arguments used in (7.13).
Next we prove that $G_{\sigma_1,\overline{\sigma}_2}$ is a contraction mapping using the above inequalities. More precisely, by definition, see (7.1), we have:

$$\left\|G_{\sigma_1,\overline{\sigma}_2}(J_1, V_1, J_2, V_2) - G_{\sigma_1,\overline{\sigma}_2}(J_1', V_1', J_2', V_2')\right\|$$
$$= \left\|(M_{\sigma_1}^1(J_2, V_2) - M_{\sigma_1}^1(J_2', V_2'), K^1(J_2, V_2) - K^1(J_2', V_2'),\right.$$
$$\left. M_{\overline{\sigma}_2}^2(J_1, V_1) - M_{\overline{\sigma}_2}^2(J_1', V_1'), K^2(J_1, V_1) - K^2(J_1', V_1'))\right\|$$

$$= \max\left\{ \left\| M^1_{\sigma_1}(J_2, V_2) - M^1_{\sigma_1}(J'_2, V'_2) \right\|, \left\| K^1(J_2, V_2) - K^1(J'_2, V'_2) \right\|, \right.$$

$$\left. \left\| M^2_{\overline{\sigma}_2}(J_1, V_1) - M^2_{\overline{\sigma}_2}(J'_1, V'_1) \right\|, \left\| K^2(J_1, V_1) - K^2(J'_1, V'_1) \right\| \right\} \qquad \text{rearranging}$$

$$\le \max\left\{ \frac{1}{\gamma}, \gamma\beta \right\} \max\{ \|J_1 - J'_1\|, \|V_1 - V'_1\|, \|J_2 - J'_2\|, \|V_2 - V'_2\| \}$$

where the final inequality follows from (7.11)–(7.14).

Therefore, since $\max\{\frac{1}{\gamma}, \gamma\beta\} < 1$ and assuming $\sigma_1$ is PWC and $\overline{\sigma}_2$ is CON-PWC, we have that $G_{\sigma_1, \overline{\sigma}_2}$ is a contraction mapping for all $(\sigma_1, \overline{\sigma}_2) \in \Sigma_1 \times \overline{\Sigma}_2$. Now since $\mathbb{F}(S) \times \mathbb{F}(S) \times \mathbb{F}(\Lambda_1) \times \mathbb{F}(\Lambda_1)$ is a complete metric space with respect to the norm (7.6), we conclude that $G_{\sigma_1, \overline{\sigma}_2}$ has a unique fixed point $(J_1^\star, V_1^\star, J_2^\star, V_2^\star)$. In view of (7.2)–(7.5), this fixed point satisfies for each $(s, u_1) \in \Lambda_1$:

$$(7.15) \qquad J_1^\star(s) = \frac{1}{\gamma} \min\{J_2^\star(s, \sigma_1(s)), V_2^\star(s, \sigma_1(s))\}$$

$$(7.16) \qquad V_1^\star(s) = \max_{u_1 \in \mathbb{P}(A_1(s))} \frac{1}{\gamma} \min\{J_2^\star(s, u_1), V_2^\star(s, u_1)\}$$

$$(7.17) \quad J_2^\star(s, u_1) = \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} Q(s, (a_1, a_2)\gamma \max[J_1^\star, V_1^\star]) u_1(a_1) \overline{\sigma}_2(a_2 \mid (s, u_1))$$

$$(7.18) \quad V_2^\star(s, u_1) = \min_{u_2 \in \mathbb{P}(A_2(s))} \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)}$$
$$Q(s, (a_1, a_2), \gamma \max[J_1^\star, V_1^\star]) u_1(a_1) u_2(a_2).$$

By combining (7.15) and (7.16), and combining (7.17) and (7.18), we have for each $(s, u_1) \in \Lambda_1$:

$$J_1^\star(s) \le V_1^\star(s),$$
$$J_2^\star(s, u_1) \ge V_2^\star(s, u_1)$$

from which (7.16) and (7.18) can be simplified to:

$$V_1^\star(s) = \max_{u_1 \in \mathbb{P}(A_1(s))} \frac{1}{\gamma} V_2^\star(s, u_1)$$
$$V_2^\star(s, u_1) = \min_{u_2 \in \mathbb{P}(A_2(s))} \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} Q(s, (a_1, a_2), \gamma V_1^\star) u_1(a_1) u_2(a_2)$$

implying that:

$$\gamma V_1^\star(s) = \max_{u_1 \in \mathbb{P}(A_1(s))} \min_{u_2 \in \mathbb{P}(A_2(s))} \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} Q(s, (a_1, a_2), \gamma V_1^\star) u_1(a_1) u_2(a_2)$$
$$= T(\gamma V_1^\star).$$

Thus, we have $\gamma V_1^\star = V^\star$, which completes the proof. □

Next we compute the strategies for the agents based on the function returned by the Minimax-action-free PI algorithm.

DEFINITION 7.16 (CON-3 solution). *Let $f \in \mathbb{F}(\Lambda_{12})$. If there exists a BFCP $\Phi$ of $S$ where, for each region $\phi \in \Phi$, $A(s) = A(s')$ for all $s, s' \in \phi$, and there exists a pair of probability measures $u_1^\phi \in \mathbb{P}(A_1(s))$ and $u_2^\phi \in \mathbb{P}(A_2(s))$ for all $s \in \phi$ such that $f(s, u_1^\phi, u_2^\phi) = \max_{u_1 \in \mathbb{P}(A_1(s))} \min_{u_2 \in \mathbb{P}(A_2(s))} f(s, u_1, u_2)$ for all $s \in \phi$, and if $\sigma_1 \in \Sigma_1$, $\sigma_2 \in \Sigma_2$ are such that $\sigma_1(s) = u_1^\phi$ and $\sigma_2(s) = u_2^\phi$ for all $s \in \phi$, then $(\sigma_1, \sigma_2)$ is a constant-3 (CON-3) solution of $f$ over $\phi$.*
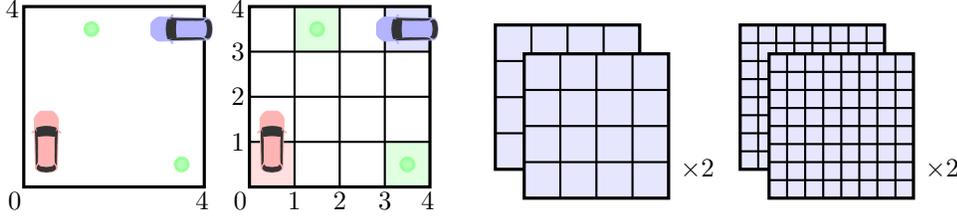
FIG. 3. *Two-vehicle dynamic parking. Left: continuous environment and $4 \times 4$ discrete percepts of the agents. Right: sum of perception and reward FCPs: $\Phi_P + \sum_{\alpha \in A} \Phi_R^\alpha$ (4-by-4 grid) and converged BFCP (8-by-8 grid).*

LEMMA 7.17 (PWC strategies). *Let $V$ be a function returned by Minimax-action-free PI. If $(\sigma_1, \sigma_2) \in \Sigma$ achieves the maximum and the minimum in Definition 5.4 for $V$ and all $s \in S$ via a CON-3 solution, then $\sigma_1$ and $\sigma_2$ are PWC stochastic kernels.*

*Proof.* By Theorems 7.11 and 7.15, $V$ is bounded, PWC Borel measurable. Now for any $\alpha \in A$, $Q(\cdot, \alpha, V) : S \to \mathbb{R}$ is bounded, PWC Borel measurable by Lemma 6.1. Let $\Phi_Q$ be a BFCP of $S$ such that $Q(\cdot, \alpha, V)$ is constant on each region in $\Phi_Q$ for all $\alpha \in A$, and $\Phi_A$ be a BFCP of $S$ such that $A(s)$ is constant on each region in $\Phi_A$. Then, given $u_1 \in \mathbb{P}(A_1(s))$ and $u_2 \in \mathbb{P}(A_2(s))$, the function $Q'(\cdot, u_1, u_2) : S \to \mathbb{R}$ where for any $s \in S$:

$$Q'(s, u_1, u_2) = \sum_{a_1 \in A_1(s)} \sum_{a_2 \in A_2(s)} Q(s, (a_1, a_2), V) u_1(a_1) u_2(a_2)$$

is constant in each region of $\Phi_Q + \Phi_A$. Therefore there exists a CON-3 solution $(\sigma_1, \sigma_2)$ of $Q'(s, u_1, u_2)$ and since $\Phi_Q + \Phi_A$ is a BFCP, the result follows.  □

**8. Case study.** To illustrate our approach, we build a small NS-CSG model of a dynamic vehicle parking problem (a static version is considered in [3]). We then synthesise strategies for it using an implementation of the VI algorithm of Section 6.

**8.1. NS-CSG modelling.** Figure 3 (left) shows two agents (vehicles) $N = \{1, 2\}$ and two parking spots $M = \{m_1, m_2\}$ in a (continuous) rectangular region $\mathcal{R} = \{(x, y) \in \mathbb{R}^2 \mid 0 \le x, y \le 4\}$ representing the states of the environment. The coordinates of these parking spots are $w_{m_1} = (1.5, 3.5)$ and $w_{m_2} = (3.5, 0.5)$, which are assumed to be known to agents. The zero-sum objective is for $\mathsf{Ag}_1$ to park at a chosen parking spot without crashing into $\mathsf{Ag}_2$, while $\mathsf{Ag}_2$ seeks to prevent $\mathsf{Ag}_1$ from parking. Suppose that the two agents can start from any positions in $\mathcal{R}$ and have the same speed. The agents' percepts are implemented via the same linear regression model for multi-class classification, see Figure 3 (left). The actions of the agents are to move either up, down, left or right, or park.

Formally, the agents and the environment are defined as follows.

- The local state spaces of the agents are given by $Loc_1 = \{w_{m_1}, w_{m_2}\}$ and $Loc_2 = \{loc_2\}$, where the local state of $\mathsf{Ag}_1$ is the coordinate of the parking spot it aims to approach and the local state of $\mathsf{Ag}_2$ is a dummy state. Here we work with integer percepts and the set of percepts of $\mathsf{Ag}_i$ for $1 \le i \le 2$ is given by $Per_i = \{1, 2, 3, 4\}^2 \times \{1, 2, 3, 4\}^2$ and the discrete coordinates the agent perceives are the positions of the two vehicles.
- The set of environment states is $S_E = \mathcal{R} \times \mathcal{R}$ and the environment is in state $s_E = (w_1, w_2)$ if $w_i \in \mathcal{R}$ is the continuous coordinate vector of $\mathsf{Ag}_i$.

- The action set of $\mathsf{Ag}_i$ for $1 \le i \le 2$ is given by $A_i = \{u_i, d_i, l_i, r_i, p_i\}$ representing the four possible directions an agent can choose to move and park.
- The action availability function of $\mathsf{Ag}_i$ for $1 \le i \le 2$ is given by:

$$\Delta_i(loc_i, (\bar{w}_1, \bar{w}_2)) = \left\{ \begin{array}{ll} \{u_i, d_i, l_i, r_i, p_i\} & \text{if } \bar{w}_i \in \{\bar{w}_{m_1}, \bar{w}_{m_2}\} \\ \{u_i, d_i, l_i, r_i\} & \text{otherwise.} \end{array} \right.$$

- In line with the percepts, the observation functions of the agents are both implemented via the same linear regression model for multi-class classification with classifier boundaries given by

$$\bigcup\nolimits_{\ell \in \{1,2,3\}} \left( \{(x, y) \in \mathcal{R} \mid x = \ell\} \cup \{(x, y) \in \mathcal{R} \mid y = \ell\} \right).$$

The Borel measurable tie-breaking rule used here is assigning boundary points to the left and lower discrete coordinate, e.g., the class of environment state $(2, 3, 3.1, 1.7)$ is $(2, 3, 4, 2)$. We denote by $\bar{w}_i$ the coordinate of the class of $w_i$ under $obs_i$.
- For any $s_1 = (z_1, (\bar{w}_1, \bar{w}_2)) \in S_1$ and $\alpha = (a_1, a_2)$ to define $\delta_1$ we have the following two cases to consider:
    - if $\|\bar{z}_1 - \bar{w}_1\|_2 > \|\bar{z}_1 - \bar{w}_2\|_2$, i.e. $\mathsf{Ag}_2$ is closer to the chosen parking spot of $\mathsf{Ag}_1$, and the joint action $(a_1, a_2)$ indicate both agents are trying to approach $z_1$, then $\delta_1(s_1, \alpha)(w_{m_j}) = 0.5$ for $1 \le j \le 2$, i.e., $\mathsf{Ag}_1$ changes its chosen parking spot with probability 0.5;
    - otherwise $\delta_1(s_1, \alpha)(z_1) = 1$, i.e. $\mathsf{Ag}_1$ sticks with its chosen parking spot.
- Since $Loc_2 = \{loc_2\}$, for any $s_2 = (loc_2, (\bar{w}_1, \bar{w}_2)) \in S_2$ and $\alpha = (a_1, a_2)$ we let $\delta_2(s_2, \alpha) = loc_2$.
- For any $(w_1, w_2) \in S_E$ and $(a_1, a_2) \in A$, we let $\delta_E((w_1, w_2), (a_1, a_2)) = (w_1', w_2')$ where for $1 \le i \le 2$

$$w_i' = \left\{ \begin{array}{ll} w_i + d_{a_i} \Delta t & \text{if } (w_i + d_{a_i} \Delta t) \in \mathcal{R} \\ w_i & \text{otherwise} \end{array} \right.$$

$d_{a_i}$ is the direction of movement of the action $a_i$, e.g., $d_{u_i} = (0, 1)$ and $d_{l_i} = (-1, 0)$ and $\Delta t = 0.5$ is the time step.

We consider the reward structure where all action rewards are 0, there is a negative reward in states where the $\mathsf{Ag}_1$ has yet to reach its current parking spot and the agents crash (that is, $r_S(s) = -1000$ if $\bar{w}_1 = \bar{w}_2 \ne \bar{z}_1$), a positive reward in states where the $\mathsf{Ag}_1$ has reached its parking spot which increases if the agents have not crashed (that is, $r_S(s) = 500$ if $\bar{w}_1 = \bar{w}_2 = \bar{z}_1$ and $r_S(s) = 1000$ if $\bar{w}_1 = \bar{z}_1$ and $\bar{w}_1 \ne \bar{w}_2$), and otherwise the state reward is 0. We fix the discount factor $\beta = 0.6$.

**8.2. Strategy synthesis.** We implement our VI algorithm (see Algorithm 6.1) using a polyhedra representation of regions. We have partitioned the state space of the game into two sets corresponding to the two possible local states of $\mathsf{Ag}_1$.

The VI algorithm converges after 46 iterations when $\varepsilon = 10^{-6}$ and takes $3,825$s to complete. For each set in the partition of the state space, the BFCP of this set converges to the product of two $8 \times 8$ grids, see Figure 3 (right). For the current chosen parking spot of $\mathsf{Ag}_1$ (red square) and coordinate of $\mathsf{Ag}_2$ (purple triangle), the value function with respect to the coordinate of $\mathsf{Ag}_1$ is presented in Figure 4 (left) and shows that, the closer $\mathsf{Ag}_1$ is to its chosen parking spot, the higher the optimal value. The lightest-colour class is caused by an immediate crash, and its position follows from the observation function.
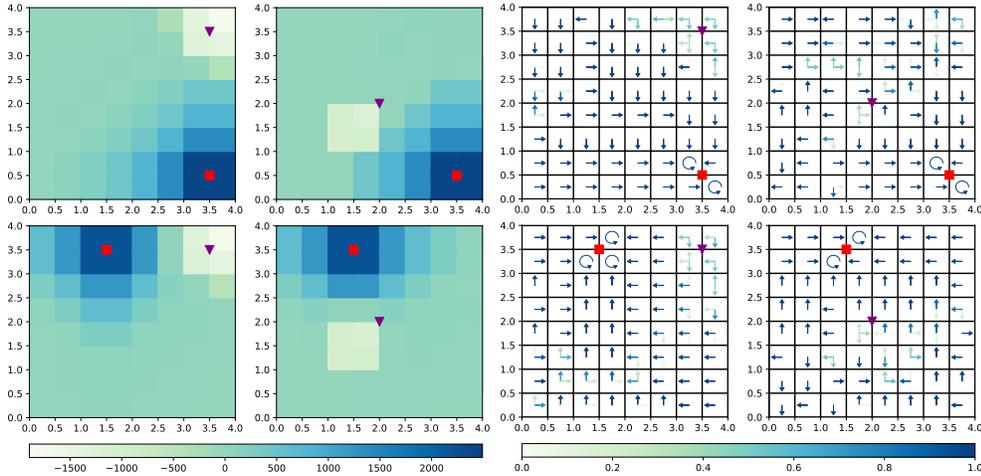
FIG. 4. *Value function (left) and optimal strategy (right) over different coordinates of* $Ag_1$ *for a fixed local state of* $Ag_1$ *(red square) and a fixed coordinate of* $Ag_2$ *(purple triangle).*

An (approximate) optimal strategy for $Ag_1$ is presented in Figure 4 (right), where the colour of an arrow is proportional to the probability of moving in that direction and the rotating arrow represents the parking action. There are several interesting choices which are not intuitive. For example, although a crash cannot be avoided before reaching its current parking spot, $Ag_1$ moves left when in $[1.0, 1.5] \times [3.5, 4.0]$ (top left) as it is better to crash later in a discounted setting and $Ag_1$ moves right when in $[1.5, 2.0] \times [0.5, 1.0]$ (down right) since by moving in this direction it will meet the conditions to (randomly) update its chosen parking spot required by $Ag_1$'s local transition function.

**9. Conclusions.** We have proposed a novel modelling formalism called neuro-symbolic concurrent stochastic games (NS-CSGs) for representing probabilistic finite-state agents with neural network perception mechanisms interacting in a shared, continuous-state environment. We proved the existence of the value of NS-CSGs under Borel measurability and piecewise constant restrictions, and presented the first implementable policy iteration and value iteration algorithms for computing the values and strategies of NS-CSGs with respect to zero-sum discounted cumulative rewards, assuming a fully observable setting. We illustrated our approach by modelling a dynamic vehicle parking problem as an NS-CSG and synthesising optimal values and strategies using value iteration.

A number of improvements and challenges remain as future work. For example, in value iteration, the number of regions increases exponentially as the number of dimensions in the environment's state space increases, which directly influences computation costs in terms of both time and space. In addition, generalising to other observation functions, e.g., using Sigmoid NNs, will necessitate working with abstractions in order to represent regions with a more general shape. There is also work to be done in devising practical stopping criteria for iterative methods over BFCPs. Next, we plan to investigate extending the algorithms to allow for partial observability, as discussed in Section 3, and nonzero-sum objectives through equilibria-based properties [16]. For the latter, preliminary progress has been made in [35], building upon the NS-CSG model presented in this paper.

## REFERENCES

[1] M. E. Akintunde, E. Botoeva, P. Kouvaros, and A. Lomuscio, *Verifying Strategic Abilities of Neural-symbolic Multi-agent Systems*, in Proc. 17th Int. Conf. Principles of Knowledge Representation and Reasoning (KR'20), IJCAI Organization, 9 2020, pp. 22–32.

[2] G. Anderson, A. Verma, I. Dillig, and S. Chaudhuri, *Neurosymbolic reinforcement learning with formally verified exploration*, in Proc. 33rd Int. Conf. Advances in Neural Information Processing Systems (NeurIPS'20), Curran Associates, Inc., 2020.

[3] D. Ayala, O. Wolfson, B. Xu, B. Dasgupta, and J. Lin, *Parking slot assignment games*, in Proc. 19th ACM SIGSPATIAL Int. Conf. Advances in Geographic Information Systems (GIS'11), ACM, 2011, p. 299–308.

[4] A. Basu and Ł. Stettner, *Zero-sum Markov games with impulse controls*, SIAM Journal on Control and Optimization, 58 (2020), pp. 580–604.

[5] D. P. Bertsekas, *Distributed asynchronous policy iteration for sequential zero-sum games and minimax control*, arXiv:2107.10406, (2021).

[6] D. P. Bertsekas and H. Yu, *Q-learning and enhanced policy iteration in discounted dynamic programming*, Mathematics of Operations Research, 37 (2012), pp. 66–94.

[7] N. Brown, A. Bakhtin, A. Lerer, and Q. Gong, *Combining deep reinforcement learning and search for imperfect-information games*, in Proc. 33rd Int. Conf. Advances in Neural Information Processing Systems (NeurIPS'20), Curran Associates, Inc., 2020.

[8] A. Cosso, *Stochastic differential games involving impulse controls and double-obstacle quasi-variational inequalities*, SIAM Journal on Control and Optimization, 51 (2013), pp. 2102–2131.

[9] J. Filar and K. Vrieze, *Competitive Markov decision processes*, Springer, 1997.

[10] J. Gupta, M. Egorov, and M. Kochenderfer, *Cooperative multi-agent control using deep reinforcement learning*, in Proc. 16th Int. Conf. Autonomous Agents and Multiagent Systems (AAMAS'17), Springer, 2017, pp. 66–83.

[11] O. Hernández-Lerma and J. Lasserre, *Zero-sum stochastic games in borel spaces: average payoff criteria*, SIAM Journal on Control and Optimization, 39 (2000), pp. 1520–1539.

[12] A. J. Hoffman and R. M. Karp, *On non-terminating stochastic games*, Management Science, 12 (1966), pp. 359–370.

[13] I. Hogeboom-Burr and S. Yuksel, *Comparison of information structures for zero-sum games and a partial converse to Blackwell ordering in standard borel spaces*, SIAM Journal on Control and Optimization, 59 (2021), pp. 1781–1803.

[14] V. Kovařík, M. Schmid, N. Burch, M. Bowling, and V. Lisý, *Rethinking formal models of partially observable multiagent decision making*, Artificial Intelligence, 303 (2022), p. 103645.

[15] P. R. Kumar and T.-H. Shiau, *Existence of value and randomized strategies in zero-sum discrete-time stochastic dynamic games*, SIAM Journal on Control and Optimization, 19 (1981), pp. 617–634.

[16] M. Kwiatkowska, G. Norman, D. Parker, and G. Santos, *Automatic verification of concurrent stochastic systems*, Formal Methods in System Design, (2021), pp. 1–63.

[17] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, *Multi-agent actor-critic for mixed cooperative-competitive environments*, in Proc. 31st Int. Conf. Neural Information Processing Systems (NIPS'17), Curran Associates Inc., 2017, p. 6382–6393.

[18] A. Maitra and T. Parthasarathy, *On stochastic games*, Journal of Optimization Theory and Applications, 5 (1970), pp. 289–300.

[19] A. S. Nowak, *Universally measurable strategies in zero-sum stochastic games*, The Annals of Probability, 13 (1985), pp. 269–287.

[20] A. S. Nowak, *Optimal strategies in a class of zero-sum ergodic stochastic games*, Mathematical methods of operations research, 50 (1999), pp. 399–419.

[21] K. R. Parthasarathy, *Probability measures on metric spaces*, American Mathematical Soc., 1967.

[22] J. Perolat, B. Scherrer, B. Piot, and O. Pietquin, *Approximate dynamic programming for two-player zero-sum Markov games*, in Proc. Int. Conf. Machine Learning, PMLR, 2015, pp. 1321–1329.

[23] M. A. Pollatschek and B. Avi-Itzhak, *Algorithms for stochastic games with geometrical interpretation*, Management Science, 15 (1969), pp. 399–415.

[24] L. D. Raedt, S. Dumancic, R. Manhaeve, and G. Marra, *From statistical relational to neural-symbolic artificial intelligence*, in Proc. 29th Int. Joint Conf. Artificial Intelligence (IJCAI'20), IJCAI Organization, 07 2020, pp. 4943–4950.

[25] J. H. Reif, *Universal games of incomplete information*, in Proc. ACM Symp. Theory of Com-

puting (STOC'79), ACM, 1979, pp. 288–308.

[26] J. H. Reif, *The complexity of two-player games of incomplete information*, Journal of Computer and System Sciences, 29 (1984), pp. 274–301.

[27] H. L. Royden and P. Fitzpatrick, *Real analysis (fourth edition)*, Macmillan New York, 2010.

[28] S. Shalev-Shwartz, S. Shammah, and A. Shashua, *Safe, multi-agent, reinforcement learning for autonomous driving*, arXiv:1610.03295, (2016).

[29] L. S. Shapley, *Stochastic games*, Proc. National Academy of Sciences, 39 (1953), pp. 1095–1100.

[30] D. Silver, A. Huang, C. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al., *Mastering the game of go with deep neural networks and tree search*, Nature, 529 (2016), pp. 484–489.

[31] B. Tolwinski, *Newton-type methods for stochastic games*, in Differential games and applications, Springer, 1989, pp. 128–144.

[32] J. Van Der Wal, *Discounted Markov games: Generalized policy iteration method*, Journal of Optimization Theory and Applications, 25 (1978), pp. 125–138.

[33] J. von Neumann, O. Morgenstern, H. Kuhn, and A. Rubinstein, *Theory of Games and Economic Behavior*, Princeton University Press, 1944.

[34] R. Yan, X. Duan, Z. Shi, Y. Zhong, J. Marden, and F. Bullo, *Policy evaluation and seeking for multi-agent reinforcement learning via best response*, IEEE Transactions on Automatic Control, 67 (2022), pp. 1898–1913.

[35] R. Yan, G. Santos, X. Duan, D. Parker, and M. Kwiatkowska, *Finite-horizon equilibria for neuro-symbolic concurrent stochastic games*, in Proc. 38th Conf. Uncertainty in Artificial Intelligence (UAI'22), AUAI Press, 2022.

[36] H. Yu, *On convergence of value iteration for a class of total cost Markov decision processes*, SIAM Journal on Control and Optimization, 53 (2015), pp. 1982–2016.

[37] H. Yu and D. P. Bertsekas, *A mixed value and policy iteration method for stochastic control with universally measurable policies*, Mathematics of Operations Research, 40 (2015), pp. 926–968.