

# Partially Observable Game-Theoretic Agent Programming in Golog

Alberto Finzi<sup>a</sup>, Thomas Lukasiewicz<sup>b</sup>

<sup>a</sup>*DIETI, Università di Napoli “Federico II”, Italy*

<sup>b</sup>*Department of Computer Science, University of Oxford, UK*

---

## Abstract

In this paper, we present the agent programming language POGTGolog (Partially Observable Game-Theoretic Golog), which integrates explicit agent programming in Golog with game-theoretic multi-agent planning in partially observable stochastic games. In this framework, we assume one team of cooperative agents acting under partial observability, where the agents may also have different initial belief states and not necessarily the same rewards. POGTGolog allows for specifying a partial control program in a high-level logical language, which is then completed by an interpreter in an optimal way. To this end, we define a formal semantics of POGTGolog programs in terms of Nash equilibria, and we then specify a POGTGolog interpreter that computes one of these Nash equilibria.

*Keywords:* Agent programming, cognitive robotics, multi-agent systems, reasoning about actions and change.

---

## 1. Introduction

High-level agent programming is an approach to autonomous agent control design that combines explicit agent programming and planning [59, 25]. While an agent programming language provides high-level primitives and constructs that enable a programmer to easily specify the agent behavior, in the planning approach, goal or rewards functions are specified, and the agent is given a planning ability to achieve a goal or to maximize a reward function. An integration of these approaches in a decision-theoretic setting has been proposed through the

---

*Email addresses:* alberto.finzi@unina.it (Alberto Finzi),  
thomas.lukasiewicz@cs.ox.ac.uk (Thomas Lukasiewicz)

seminal language DTGolog [8], which integrates explicit agent programming in Golog [59] with decision-theoretic planning in (fully observable) Markov decision processes (MDPs) [55]. It allows for partially specifying a control program in a high-level language as well as for optimally filling in missing details through decision-theoretic planning. It can thus be seen as a decision-theoretic extension to Golog, where choices left to the agent are made by maximizing expected utility. From a different perspective, it can also be seen as a formalism that gives advice to a decision-theoretic planner, since it naturally constrains the search space by providing fragments of control knowledge (partially specified control programs).

The agent programming language DTGolog has several other nice features, as it is closely related to first-order extensions of decision-theoretic planning (see especially [7, 70, 32, 63]), which allow for (i) compactly representing decision-theoretic planning problems without explicitly referring to atomic states and state transitions, (ii) exploiting such compact representations for efficiently solving large-scale problems, and (iii) nice properties such as *modularity* (parts of the specification can be easily added, removed, or modified) and *elaboration tolerance* (solutions can be easily reused for similar problems with few or no extra cost).

However, DTGolog is designed only for the single-agent framework. That is, the model of the world essentially consists of a single agent that we control by a DTGolog program and the environment summarized in “nature”. But there are many applications where we encounter multiple agents, which may compete against each other, or which may also cooperate with each other. For example, in *robotic soccer*, we have two competing teams of agents, where each team consists of cooperating agents. The optimal actions of one agent generally depend on the actions of all the other (“enemy” and “friend”) agents. In particular, there is a bidirectional dependence between the actions of two different agents, which generally makes it inappropriate to model enemies and friends of the agent that we control simply as a part of “nature”. As an example for an important cooperative domain, in *robotic rescue*, mobile agents may be used in the emergency area to acquire new detailed information (e.g., the locations of injured people in the emergency area) or to perform certain rescue operations. In general, acquiring information as well as performing rescue operations involves several and different rescue elements (agents and/or teams of agents), which cannot effectively handle the rescue situation on their own. Only the cooperative work among all the rescue elements may solve it. As most of the rescue tasks involve a certain level of risk for humans (depending on the type of rescue situation), mobile agents can play a major role in rescue situations, especially teams of cooperative heterogeneous mobile agents.

This is the motivation behind GTGolog [18], which is a generalization of DT-

Golog that integrates agent programming in Golog with game-theoretic multi-agent planning in (fully observable) stochastic games [49], also called Markov games [64, 41]. The following example shows a program in GTGolog.

**Example 1.1.** (*Rugby Domain*) Consider a rugby player  $a_1$ , who is deciding his next  $n > 0$  moves and wants to cooperate with a teammate  $a_2$ . He has to deliberate about if and when it is worth to pass the ball. His options can be encoded by the following GTGolog program:

```

proc step( $n$ )
if (haveBall( $a_1$ )  $\wedge$   $n > 0$ ) then
   $\pi x (\pi y (\mathbf{choice}(a_1 : moveTo(x) \mid passTo(a_2)) \parallel$ 
     $\mathbf{choice}(a_2 : moveTo(y) \mid receive(a_1))))$ ;
  step( $n-1$ )
end.

```

This program encodes that while agent  $a_1$  is the ball owner and  $n > 0$ , the two agents  $a_1$  and  $a_2$  perform a parallel action choice ( $\mathbf{choice}(a_1 : \cdot) \parallel \mathbf{choice}(a_2 : \cdot)$ ) in which  $a_1$  (resp.,  $a_2$ ) can either go somewhere or pass (resp., receive) the ball. In this choice, the two agents should simultaneously select one of the two possible actions (*moveTo* or *passTo* for  $a_1$ ; *moveTo* or *receive* for  $a_2$ ) according to some strategical game-theoretic criterion (as explained in Section 3.3). The preconditions and effects of the actions are to be formally specified in a suitable action theory. Given this high-level program and the action theory for  $a_1$  and  $a_2$ , the program interpreter should fill in the best moves for  $a_1$  and  $a_2$ , reasoning about all the possible interactions between the two agents.

However, another crucial aspect of real-world environments is that they are typically only partially observable, due to noisy and inaccurate sensors, or because some relevant parts of the environment simply cannot be sensed. For example, in the robotic rescue domain, every agent has generally only a very partial view on the environment and the other agents. But both DT- and GTGolog assume full observability, and have not been generalized to the partially observable case so far.

In this paper, we try to fill this gap. We present the agent programming language POGTGolog, which extends GTGolog and thus also DTGolog by partial observability. Its main contributions are summarized as follows:

- We present the agent programming language POGTGolog, which combines agent programming in Golog with game-theoretic multi-agent planning in

partially observable stochastic games. POGTGolog allows for modeling one team of cooperative agents under partial observability, where the agents may also have different initial belief states and not necessarily the same rewards (and thus, in some sense, the team does not necessarily have to be homogeneous). We assume multiple agents with free communication (as in [60, 56]), that is, each agent can freely communicate beliefs, actions, and observations to the other agents. We do not assume any external authority that enforces agreements among the agents, hence the agents' cooperative behavior should rely on self-enforcing strategies. Since the underlying model is Markovian, the agents' actions in a team run synchronously.

- POGTGolog allows for specifying a partial control program in a high-level language, which is then completed in an optimal way. To this end, we associate with every POGTGolog program a set of (finite-horizon) policies, which are possible (finite-horizon) program instantiations where missing details are filled in. We then define a semantics of a POGTGolog program in terms of Nash equilibria, which are optimal policies (that is, optimal instantiations) of the program.
- We define a POGTGolog interpreter and show that it associates with each program one of its Nash equilibria. We also prove that POGTGolog programs faithfully extend partially observable stochastic games. Furthermore, we illustrate the usefulness of the POGTGolog approach along several examples.

The rest of this paper is organized as follows. In Section 2, we recall the basic concepts of the situation calculus, Golog, normal form games, and partially observable stochastic games. Section 3 defines the domain theory, syntax, and semantics of POGTGolog programs. In Section 4, we give a formal specification of the POGTGolog interpreter and provide an optimality result. Section 5 illustrates the overall framework through another example. In Section 6, we discuss related work. Finally, Section 7 summarizes our main results and gives an outlook on future research. Note that detailed proofs of all results are given in Appendix A.

## 2. Preliminaries

In this section, we first recall the main concepts of the situation calculus (in its standard and concurrent version) and of the agent programming language Golog;

for further details and background, see especially [59]. We then recall the basics of normal form games and partially observable stochastic games.

### 2.1. The Situation Calculus

The situation calculus [45, 59] provides a first-order language for representing dynamically changing worlds. Its main ingredients are *actions*, *situations*, and *fluents*. An *action* is a first-order term of the form  $a(u_1, \dots, u_n)$ , where the function symbol  $a$  is its *name*, and the  $u_i$ 's are its *arguments*. All changes to the world are the result of actions. A *situation* is a first-order term encoding a sequence of actions. It is either a constant symbol or of the form  $do(a, s)$ , where  $do$  is a distinguished binary function symbol,  $a$  is an action, and  $s$  is a situation. The constant symbol  $S_0$  is the *initial situation* and represents the empty sequence, while  $do(a, s)$  encodes the sequence obtained from executing  $a$  after the sequence of  $s$ . We write  $Poss(a, s)$ , where  $Poss$  is a distinguished binary predicate symbol, to denote that it is possible to execute the action  $a$  in the situation  $s$ . A (*relational*) *fluent* represents a world or agent property that may change when executing an action. It is a predicate symbol whose most right argument is a situation. A situation calculus formula is *uniform* in a situation  $s$  iff (i) it does not mention the predicates  $Poss$  and  $\sqsubset$  (which denotes the proper subsequence relationship on situations), (ii) it does not quantify over situation variables, (iii) it does not mention equality on situations, and (iv) every situation in the situation argument of a fluent coincides with  $s$  [59].

**Example 2.1.** (*Rugby Domain cont'd*) The action  $moveTo(r, x, y)$  may stand for moving the agent  $r$  to the position  $(x, y)$ , while the situation  $do(moveTo(r, 1, 2), do(moveTo(r, 3, 4), S_0))$  stands for executing  $moveTo(r, 1, 2)$  after the execution of  $moveTo(r, 3, 4)$  in the initial situation  $S_0$ . The (relational) fluent  $at(r, x, y, s)$  may express that the agent  $r$  is at the position  $(x, y)$  in the situation  $s$ .

In the situation calculus, a dynamic domain is represented by a *basic action theory (BAT)*  $AT = (\Sigma, \mathcal{D}_{una}, \mathcal{D}_{S_0}, \mathcal{D}_{ssa}, \mathcal{D}_{ap})$ , where:

- $\Sigma$  is the set of (domain-independent) foundational axioms for situations [59].
- $\mathcal{D}_{una}$  is the set of unique name axioms for actions, which express that different actions are interpreted in a different way (that is, any two actions with different names or different arguments have a different meaning).

- $\mathcal{D}_{S_0}$  is a set of first-order formulas that are uniform in  $S_0$ , which describe the initial state of the domain (represented by  $S_0$ ).
- $\mathcal{D}_{ssa}$  is the set of *successor state axioms* [58, 59]. For each fluent  $F(\mathbf{x}, s)$ , it contains an axiom of the form  $F(\mathbf{x}, do(a, s)) \equiv \Phi_F(\mathbf{x}, a, s)$ , where  $\Phi_F(\mathbf{x}, a, s)$  is a formula uniform in  $s$  with free variables among  $\mathbf{x}, a, s$ . These axioms specify the truth of the fluent  $F$  in the next situation  $do(a, s)$  in terms of the current situation  $s$ , and are a solution to the frame problem (for deterministic actions).
- $\mathcal{D}_{ap}$  is the set of *action precondition axioms*. For each action  $a$ , it contains an axiom of the form  $Poss(a(\mathbf{x}), s) \equiv \Pi(\mathbf{x}, s)$ , where  $\Pi$  is a formula uniform in  $s$  with free variables among  $\mathbf{x}, s$ . This axiom specifies the preconditions of  $a$ .

**Example 2.2.** (*Rugby Domain cont'd*) The formula  $at(r, 1, 2, S_0) \wedge at(r', 3, 4, S_0)$  in  $\mathcal{D}_{S_0}$  may express that the agents  $r$  and  $r'$  are initially at the positions  $(1, 2)$  and  $(3, 4)$ , respectively. Notice that first-order logic formulas enable an expressive denotation of the initial state, for instance,  $\exists x, y (at(r, x, 2, S_0)) \wedge at(r', y, 4, S_0)$  provides an incomplete specification of the initial state (see [25]). The axiom  $at(r, x, y, do(a, s)) \equiv a = moveTo(r, x, y) \vee (at(r, x, y, s) \wedge \neg \exists x', y' (a = moveTo(r, x', y')))$  in  $\mathcal{D}_{ssa}$  may express that the agent  $r$  is at the position  $(x, y)$  in the situation  $do(a, s)$  iff either  $r$  moves to  $(x, y)$  in the situation  $s$ , or  $r$  is already at the position  $(x, y)$  and does not move away in  $s$ . The axiom  $Poss(moveTo(r, x, y), s) \equiv \neg \exists r' at(r', x, y, s)$  in  $\mathcal{D}_{ap}$  may express that it is possible to move the agent  $r$  to  $(x, y)$  in  $s$  iff no agent  $r'$  is at  $(x, y)$  in  $s$ .

The BAT specification enables a first-order logic representation of the domain theory and provides a reasoning machinery based on regression and first-order inference in the initial state. Additional details can be found in [59, 53, 25].

In this paper, we use the concurrent version of the situation calculus [59], which is an extension of the above standard situation calculus by concurrent actions. A *concurrent action*  $c$  is a set of standard actions, which are concurrently executed when  $c$  is executed. A situation is then a sequence of concurrent actions of the form  $do(c_m, \dots, do(c_0, S_0))$ , where  $do(c, s)$  states that all the simple actions  $a$  in  $c$  are executed at the same time in the situation  $s$ .

In order to encode concurrent actions, some slight modifications to standard basic action theories are necessary. In particular, the successor state axioms in

$\mathcal{D}_{ssa}$  are now defined relative to concurrent actions. Furthermore, the action preconditions in  $\mathcal{D}_{ap}$  are extended by further axioms expressing (i) that a singleton concurrent action  $c = \{a\}$  is executable if its standard action  $a$  is executable, (ii) that if a concurrent action is executable, then it is nonempty and all its standard actions are executable, and (iii) preconditions for concurrent actions. Note that precondition axioms for standard actions are in general not sufficient, since two standard actions may each be executable, but their concurrent execution may not be permitted. This *precondition interaction problem* [59] (see also [52] for a discussion) requires some domain-dependent extra precondition axioms.

**Example 2.3.** (*Rugby Domain cont'd*) The axiom  $at(r, x, y, do(a, s)) \equiv a = moveTo(r, x, y) \vee (at(r, x, y, s) \wedge \neg \exists x', y' (a = moveTo(r, x', y')))$  in  $\mathcal{D}_{ssa}$  in the standard situation calculus is replaced by the axiom  $at(r, x, y, do(c, s)) \equiv moveTo(r, x, y) \in c \vee (at(r, x, y, s) \wedge \neg \exists x', y' (moveTo(r, x', y') \in c))$  in the concurrent one.

## 2.2. Golog

Golog is an agent programming language that is based on the situation calculus. It allows for constructing complex actions (also called *programs*) from (standard or concurrent) primitive actions that are defined in a basic action theory  $AT$ , where standard (and not so standard) Algol-like control constructs can be used. *Programs*  $p$  in Golog have one of the following forms (where  $c$  is a (standard or concurrent) primitive action,  $\phi$  is a *condition*, which is obtained from a situation calculus formula that is uniform in  $s$  by suppressing the situation argument,  $p, p_1, p_2, \dots, p_n$  are programs,  $P_1, \dots, P_n$  are procedure names, and  $x, x_1, \dots, x_n$  are arguments):

- (1) *Primitive action*:  $c$ . Do  $c$ .
- (2) *Test action*:  $\phi?$ . Test the truth of  $\phi$  in the current situation.
- (3) *Sequence*:  $[p_1; p_2]$ . Do  $p_1$  followed by  $p_2$ .
- (4) *Nondeterministic choice of two programs*:  $(p_1 \mid p_2)$ . Do either  $p_1$  or  $p_2$ .
- (5) *Nondeterministic choice of program argument*:  $\pi x (p(x))$ . Do any  $p(x)$ .
- (6) *Nondeterministic iteration*:  $p^*$ . Do  $p$  zero or more times.
- (7) *Conditional*: **if**  $\phi$  **then**  $p_1$  **else**  $p_2$ . If  $\phi$  is true in the current situation, then do  $p_1$  else do  $p_2$ .
- (8) *While-loop*: **while**  $\phi$  **do**  $p$ . While  $\phi$  is true in the current situation, do  $p$ .

(9) *Procedures*: **proc**  $P_1(x_1)$   $p_1$  **end**; ... ; **proc**  $P_n(x_n)$   $p_n$  **end**;  $p$ .

**Example 2.4.** (*Rugby Domain cont'd*) The Golog program **while**  $\neg at(r, 1, 2)$  **do**  $\pi x, y (moveTo(r, x, y))$  stands for “while the agent  $r$  is not at the position  $(1, 2)$ , move  $r$  to a nondeterministically chosen position  $(x, y)$ ”.

Golog has a declarative formal semantics in the situation calculus. Given a Golog program  $p$ , its execution is represented by a situation calculus formula  $Do(p, s, s')$ , which encodes that the situation  $s'$  can be reached by executing  $p$  in the situation  $s$ .

### 2.3. Normal Form Games

Normal form games from classical game theory [66] describe the possible actions of  $n \geq 2$  agents and the rewards that the agents receive when they simultaneously execute one action each. A *normal form game*  $G = (I, (A_i)_{i \in I}, (R_i)_{i \in I})$  consists of a set of *agents*  $I = \{1, \dots, n\}$  with  $n \geq 2$ , a nonempty finite set of *actions*  $A_i$  for each agent  $i \in I$ , and a *reward function*  $R_i: A \rightarrow \mathbf{R}$  for each  $i \in I$ , which associates with every *joint action*  $a \in A = \times_{i \in I} A_i$  a *reward*  $R_i(a)$  to agent  $i$ . If  $n = 2$ , then  $G$  is a *two-player normal form game* (or *matrix game*). If also  $R_1 = -R_2$ , then  $G$  is a *zero-sum matrix game*; we then often omit  $R_2$  and abbreviate  $R_1$  by  $R$ .

The behavior of the agents in a normal form game is expressed through the notions of pure and mixed strategies, which specify which of its actions an agent should execute and which of its actions an agent should execute with which probability, respectively. A *pure strategy* for agent  $i \in I$  is any action  $a_i \in A_i$ . A *pure strategy profile* is any joint action  $a \in A$ . If the agents play  $a$ , then the *reward* to agent  $i \in I$  is given by  $R_i(a)$ . A *mixed strategy* for agent  $i \in I$  is any probability distribution  $\pi_i$  over its set of actions  $A_i$ . A *mixed strategy profile*  $\pi = (\pi_i)_{i \in I}$  consists of a mixed strategy  $\pi_i$  for each agent  $i \in I$ . If the agents play  $\pi$ , then the *expected reward* to agent  $i \in I$ , denoted  $\mathbf{E}[R_i(a) \mid \pi]$  (or  $R_i(\pi)$ ), is  $\sum_{a=(a_j)_{j \in I} \in A} R_i(a) \cdot \prod_{j \in I} \pi_j(a_j)$ .

Towards optimal behavior of the agents in a normal form game, we are especially interested in mixed strategy profiles  $\pi$  (called *Nash equilibria*) where no agent has the incentive to deviate from its part, once the other agents play their parts. Formally, given a normal form game  $G = (I, (A_i)_{i \in I}, (R_i)_{i \in I})$ , a mixed strategy profile  $\pi = (\pi_i)_{i \in I}$  is a *Nash equilibrium* of  $G$  iff for every agent  $i \in I$ , it holds that  $R_i(\pi \leftarrow \pi'_i) \leq R_i(\pi)$  for every mixed strategy  $\pi'_i$ , where  $\pi \leftarrow \pi'_i$  is



obtained from  $\pi$  by replacing  $\pi_i$  by  $\pi'_i$ . Every normal form game  $G$  has at least one Nash equilibrium among its mixed (but not necessarily pure) strategy profiles, and many have multiple Nash equilibria. In the two-player case, they can be computed by linear complementary programming and linear programming in the general and the zero-sum case, respectively. A *Nash selection function*  $f$  associates with every normal form game  $G$  a unique Nash equilibrium  $f(G)$ . The expected reward to agent  $i \in I$  under  $f(G)$  is denoted by  $v_f^i(G)$ .

**Example 2.5.** (*two-finger Morra*) In a *two-finger Morra* game, two players  $E$  and  $O$  simultaneously show one or two fingers. Let  $f$  be the total numbers of fingers shown. If  $f$  is odd, then  $O$  gets  $f$  dollars from  $E$ , and if  $f$  is even, then  $E$  gets  $f$  dollars from  $O$ . A pure strategy for player  $E$  (or  $O$ ) is to show two fingers, and a mixed strategy for  $E$  (or  $O$ ) is to show one finger with the probability  $7/12$  and two fingers with the probability  $5/12$ . The mixed strategy profile where each player shows one finger resp. two fingers with the probability  $7/12$  resp.  $5/12$  is a Nash equilibrium.

#### 2.4. Partially Observable Stochastic Games

Partially observable stochastic games [49] generalize normal form games, partially observable Markov decision processes (POMDPs) [50], and decentralized POMDPs (Dec-POMDPs) [51, 29, 48]. A partially observable stochastic game consists of a set of states  $S$ , a normal form game for each state  $s \in S$ , a set of joint observations of the agents  $O$ , and a transition function that associates with every state  $s \in S$  and joint action of the agents  $a \in A$  a probability distribution on all combinations of next states  $s' \in S$  and joint observations  $o \in O$ . Formally, a *partially observable stochastic game (POSG)*  $G = (I, S, (A_i)_{i \in I}, (O_i)_{i \in I}, P, (R_i)_{i \in I})$  consists of a set of *agents*  $I = \{1, \dots, n\}$ ,  $n \geq 2$ , a nonempty finite set of *states*  $S$ , two nonempty finite sets of *actions*  $A_i$  and *observations*  $O_i$  for each  $i \in I$ , a transition function  $P: S \times A \rightarrow PD(S \times O)$ , which associates with every state  $s \in S$  and joint action  $a \in A = \times_{i \in I} A_i$  a probability distribution over  $S \times O$ , where  $O = \times_{i \in I} O_i$ , and a *reward function*  $R_i: S \times A \rightarrow \mathbf{R}$  for each  $i \in I$ , which associates with every state  $s \in S$  and joint action  $a \in A$  a *reward*  $R_i(s, a)$  to agent  $i$ .

Since the actual state  $s \in S$  of the POSG  $G$  is not fully observable, every agent  $i \in I$  has a belief state  $b_i$  that associates with every state  $s \in S$  the belief of agent  $i$  about  $s$  being the actual state. A *belief state*  $b = (b_i)_{i \in I}$  of  $G$  consists of a probability function  $b_i$  over  $S$  for each agent  $i \in I$ . The POSG  $G$  then defines probabilistic transitions between belief states as follows. The new belief

state  $b^{a,o} = (b_i^{a,o})_{i \in I}$  after executing the joint action  $a \in A$  in  $b = (b_i)_{i \in I}$  and jointly observing  $o \in O$  is given by  $b_i^{a,o}(s') = \sum_{s \in S} P(s', o | s, a) \cdot b_i(s) / P_b(b_i^{a,o} | b_i, a)$ , where  $P_b(b_i^{a,o} | b_i, a) = \sum_{s' \in S} \sum_{s \in S} P(s', o | s, a) \cdot b_i(s)$  is the probability of observing  $o$  after executing  $a$  in  $b_i$ .

The notions of finite-horizon pure (resp., mixed) policies and their rewards (resp., expected rewards) can now be defined as usual using the above probabilistic transitions between belief states. Informally, given a finite horizon  $H \geq 0$ , a pure (resp., mixed) time-dependent policy associates with every belief state  $b$  of  $G$  and number of steps to go  $h \in \{0, \dots, H\}$  a pure (resp., mixed) normal form game strategy.

Finally, the notion of a finite-horizon Nash equilibrium for a POSG  $G$  is then defined as follows. A policy  $\pi$  is a *Nash equilibrium* of  $G$  under a belief state  $b$  iff for every agent  $i \in I$ , it holds that  $G_i(H, b, \pi \leftarrow \pi'_i) \leq G_i(H, b, \pi)$  for all policies  $\pi'_i$ , where  $G_i(H, b, \alpha)$  denotes the  $H$ -step reward to agent  $i \in I$  under an initial belief state  $b = (b_i)_{i \in I}$  and the policy  $\alpha$ . A policy  $\pi$  is a *Nash equilibrium* of  $G$  iff it is a Nash equilibrium of  $G$  under every belief state  $b$ .

### 3. Partially Observable GTGolog (POGTGolog)

In this section, we present the agent programming language POGTGolog, which is a generalization of GTGolog [18] that allows for partial observability. We first define the domain theory and belief states of POGTGolog. We then introduce the syntax and the semantics of POGTGolog programs.

We focus on the case of one team of cooperative agents acting under partial observability, where the agents may also have different initial belief states and not necessarily the same rewards. That is, the agents are collaborative in the sense that they belong to the same team, but they may also have to solve conflicts of opinions about the state of the world (expressed through different initial beliefs) and its significance (expressed through different rewards). Like in [60, 56], we assume *free communication* between the agents, which means that the agents can communicate without cost, and thus we can assume that (i) each agent is aware about the initial local belief state of every other agent, and (ii) after each action execution, each agent can observe the actions of the other agents and receives their local observations.

#### 3.1. Domain Theory

POGTGolog programs are interpreted relative to a domain theory, which extends a basic action theory by stochastic actions and reward/utility functions. In

addition to a basic action theory  $AT$ , a *domain theory*  $DT = (AT, ST, OT)$  consists of a *stochastic theory*  $ST$  and an *optimization theory*  $OT$ , which are defined below. We assume a team  $I = \{1, \dots, n\}$  of  $n \geq 2$  cooperative agents  $1, \dots, n$ . The nonempty finite set of primitive actions  $A$  is partitioned into nonempty sets of primitive actions  $A_1, \dots, A_n$  of agents  $1, \dots, n$ , respectively. A *single-agent action* of agent  $i \in I$  (resp., *multi-agent action*) is any concurrent action over  $A_i$  (resp.,  $A$ ). Every multi-agent action  $c$  is associated with a nonempty finite set of *multi-agent observations*  $O_c = \times_{i \in I} O_{c,i}$ , where every  $O_{c,i}$  is a nonempty finite set of pairwise exclusive and exhaustive conditions (called *single-agent observations*) of agent  $i \in I$ .

**Example 3.1.** (*Rugby Domain cont'd*) Let  $I = \{1, 2\}$ . Then, the concurrent actions  $\{moveTo(1, 1, 2)\} \subseteq A_1$  and  $\{moveTo(2, 2, 3)\} \subseteq A_2$  are single-agent actions of agents 1 and 2, respectively, and thus also multi-agent actions, while the concurrent action  $\{moveTo(1, 1, 2), moveTo(2, 2, 3)\}$  is a multi-agent action.

A *stochastic theory*  $ST$  is a set of axioms that define stochastic actions. We represent stochastic actions through a finite set of deterministic actions, as usual [26, 8]. When a stochastic action is executed, then with a certain probability, “nature” executes exactly one of its deterministic actions and produces exactly one possible observation. We use the predicate  $stochastic(c, s, n, o, \mu)$  to encode that when executing the stochastic action  $c$  in the situation  $s$ , “nature” chooses the deterministic action  $n$ , producing the observation  $o \in O_c$  with the probability  $\mu$ . For every stochastic action  $c$  and situation  $s$ , the set of all  $(n, o, \mu)$  such that  $stochastic(c, s, n, o, \mu)$  is a probability function on the set of all deterministic components  $n$  and observations  $o \in O_c$  of  $c$  in  $s$ . We also use the notation  $prob(c, s, n, o)$  to denote the probability  $\mu$  such that  $stochastic(c, s, n, o, \mu)$ . Notice that  $prob(c, s, n, o)$  is here the analogous of the POSGs transition function, in that it provides the probability of the pair  $s' = do(n, s)$  and  $o$ , given  $c$  executed from  $s$ . We assume that  $c$  and all its nature choices  $n$  have the same preconditions. A stochastic action  $c$  is then indirectly represented by providing a successor state axiom for every associated nature choice  $n$ . The stochastic action  $c$  is *executable* in a situation  $s$  with the observation  $o \in O_c$ , denoted  $Poss(c_o, s)$ , iff  $prob(c, s, n, o) > 0$  for some  $n$ .

**Example 3.2.** (*Rugby Domain cont'd*) The stochastic action  $moveS(\alpha, x, y)$ , which (i) is observed as being either successful or a failure, and (ii) moves the

agent  $\alpha$  to either the position  $(x, y)$  or the position  $(x, y + 1)$ , can be encoded as follows:

$$\begin{aligned} stochastic(moveS(\alpha, x, y), s, n, o, \mu) \stackrel{def}{=} & \exists y' (n = moveTo(\alpha, x, y') \wedge \\ & (o = obs(success) \wedge (y' = y + 1 \wedge \mu = 0.5 \vee y' = y \wedge \mu = 0.3) \vee \\ & o = obs(failure) \wedge y' = y + 1 \wedge \mu = 0.2)). \end{aligned}$$

The precondition and successor state axioms of  $moveS(\alpha, x, y)$  are then specified through the precondition and successor state axioms of  $moveTo(\alpha, x, y')$ .

The *optimization theory*  $OT$  specifies a reward function, a utility function, and some Nash selection functions. The reward function associates with each situation  $s$  and multi-agent action  $a$ , a reward to each agent  $i \in I$ , denoted  $reward(i, a, s)$ . Note that the reward function for stochastic actions is defined through a reward function for their deterministic components. The utility function  $utility$  maps every pair consisting of a reward  $v$  and a probability value  $pr$  (that is, a real from the unit interval  $[0, 1]$ ) to a real-valued utility  $utility(v, pr)$ . We assume that  $utility(v, 1) = v$  and  $utility(v, 0) = 0$  for all rewards  $v$ . Informally, the utility function suitably mediates between the rewards of the agents and the failure of actions due to unsatisfied preconditions: differently from actions in decision-theoretic planning, the actions here may fail due to unsatisfied preconditions. Thus, the usefulness of an action/program does not only depend on its reward, but also on the probability that it is executable. Similarly to all arithmetic operations, utility functions are pre-interpreted (that is, rigid) and thus not explicitly axiomatized in the domain theory.

**Example 3.3.** (*Rugby Domain cont'd*) The reward function  $reward(1, \{moveTo(1, x, y)\}, s) = y$  encodes that the reward to agent 1 when moving to  $(x, y)$  in the situation  $s$  is given by  $y$ . An example of a utility function is  $utility(v, pr) = v \cdot pr$ .

Notice that since  $ST$  and  $OT$  are represented by means of successor state actions and preconditions, the domain theory preserves the properties of a basic action theory, including the associated inference mechanisms [59, 53].

The following example describes a more detailed domain theory for the Rugby Domain, which is inspired by the soccer domain in [41].

**Example 3.4.** (*Rugby Domain cont'd*) The rugby field (see Fig. 1) is a  $4 \times 7$  grid of 28 squares, and it includes two designated areas representing two goals. We assume a team of two agents  $\mathbf{a} = \{a_1, a_2\}$  against a (static) team of two agents

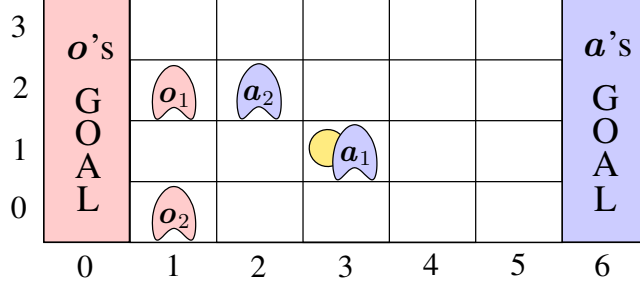


Figure 1: Rugby Domain: Two teams  $\{a_1, a_2\}$  and  $\{o_1, o_2\}$ .

$o = \{o_1, o_2\}$ , where  $a_1$  and  $o_1$  are the *captains* of  $a$  and  $o$ , respectively. Each agent occupies a square and is able to do one of the following actions on each turn:  $N$ ,  $S$ ,  $E$ ,  $W$ ,  $stand$ ,  $passTo(\alpha)$ , and  $receive(\alpha)$  (move up, move down, move right, move left, no move, pass the ball to  $\alpha$ , and receive the ball from  $\alpha$ , respectively). The ball is represented by a circle and also occupies a square. An agent is a *ball owner* iff it occupies the same square as the ball. The ball follows the moves of the ball owner, and we have a goal when the ball owner steps into the adversary goal. An agent can also pass the ball to another agent of the same team, but this is possible only if the receiving agent is not closer to the opposing end of the field than the ball, otherwise an offside fault is called by the referee, and the ball possession goes to the captain of the opposing team. When the ball owner goes into the square occupied by another agent, if the other agent stands, then the ball possession changes. Thus, a good defensive maneuver is to stand where the ball owner wants to go.

We define the domain theory  $DT = (AT, ST, OT)$  as follows. The basic action theory  $AT$  specifies the deterministic action  $move(\alpha, m)$  (encoding that agent  $\alpha$  executes  $m$ ), where  $\alpha \in a \cup o$ ,  $m \in \{N, S, E, W, stand, passTo(\alpha'), receive(\alpha')\}$ , and  $\alpha'$  is a teammate of  $\alpha$ , and the fluents  $at(\alpha, x, y, s)$  (encoding that agent  $\alpha$  is at position  $(x, y)$  in situation  $s$ ) and  $haveBall(\alpha, s)$  (encoding that agent  $\alpha$  has the ball in situation  $s$ ). They are defined by the following successor

state axioms:

$$\begin{aligned}
at(\alpha, x, y, do(c, s)) &\equiv \exists x', y' (at(\alpha, x', y', s) \wedge \exists m (move(\alpha, m) \in c \wedge \\
&((m = stand \vee m = receive \vee \exists \beta (m = passTo(\beta))) \wedge x = x' \wedge y = y') \vee \\
&(m = N \wedge x = x' \wedge y = y' + 1) \vee (m = S \wedge x = x' \wedge y = y' - 1) \vee \\
&(m = E \wedge x = x' + 1 \wedge y = y') \vee (m = W \wedge x = x' - 1 \wedge y = y'))); \\
haveBall(\alpha, do(c, s)) &\equiv \exists \beta (haveBall(\beta, s) \wedge \\
&(\alpha = \beta \wedge \neg \exists \beta' (cngBall(\beta', c, s) \vee rcvBall(\beta', c, s))) \vee \\
&(\alpha \neq \beta \wedge (cngBall(\alpha, c, s) \vee rcvBall(\alpha, c, s)))).
\end{aligned}$$

In this case,  $cngBall(\alpha, c, s)$  is true iff the ball possession changes to  $\alpha$  after an action  $c$  in  $s$  (in the case of an adversary block), that is,

$$cngBall(\alpha, c, s) \stackrel{def}{=} \exists \beta, x, y (\beta \neq \alpha \wedge \neg sameTeam(\alpha, \beta) \wedge haveBall(\beta, s) \wedge \\
at(\beta, x, y, do(c, s)) \wedge at(\alpha, x, y, s) \wedge move(\alpha, stand) \in c).$$

The predicate  $rcvBall(\alpha, c, s)$  is true iff agent  $\alpha$  receives the ball from the ball owner or because of an offside ball passing, that is,

$$\begin{aligned}
rcvBall(\alpha, c, s) &\stackrel{def}{=} \exists \beta, \alpha' (\beta \neq \alpha \wedge haveBall(\beta, s) \wedge \\
&move(\beta, passTo(\alpha')) \in c \wedge ((move(\alpha', receive(\beta)) \wedge \\
&\neg offside(\beta, \alpha', s) \wedge \alpha' = \alpha) \vee (\neg move(\alpha', receive(\beta)) \in c \vee \\
&offside(\beta, \alpha', s)) \wedge \neg sameTeam(\alpha, \alpha') \wedge captain(\alpha))).
\end{aligned}$$

In this representation,  $offside(\beta, \alpha, s)$  is true iff the players  $\alpha$  and  $\beta$  are teammates, and  $\alpha$  is closer to the goal of the adversary team than  $\beta$ . The deterministic actions  $move(\alpha, m)$  are associated with precondition axioms as follows:

$$\begin{aligned}
Poss(move(\alpha, m), s) &\equiv \neg \exists x, y (at(\alpha, x, y, s) \wedge \\
&((y = 4 \wedge m = N) \vee (y = 1 \wedge m = S) \vee \\
&(x = 6 \wedge m = E) \vee (x = 0 \wedge m = W))) \vee \\
&\exists \beta (m = passTo(\beta) \wedge haveBall(\alpha, s)) \vee \\
&\exists \beta (m = receive(\beta) \wedge haveBall(\beta, s)).
\end{aligned}$$

As for the stochastic theory  $ST$ , we assume the stochastic action  $moveS(\alpha, m)$ , which represents agent  $\alpha$ 's attempt in doing  $m \in \{N, S, E, W, stand, passTo(\beta), receive(\beta)\}$ . It can either succeed, and then the deterministic action  $move(\alpha, m)$  is executed, or it can fail, and then the deterministic action  $move(\alpha, stand)$  (that is, no change) is executed. Furthermore, after each execution of  $moveS(\alpha, m)$ ,

agent  $\alpha$  can observe the presence of a teammate  $\alpha'$  in the direction of the movement, given that agent  $\alpha'$  is visible, that is, not covered by another agent:

$$\begin{aligned}
& \text{stochastic}(\{\text{moveS}(\alpha, m)\}, s, \{a\}, \{\text{obs}(\beta, \text{res})\}, \mu) \stackrel{\text{def}}{=} \\
& \quad \exists \mu_1, \mu_2 ((a = \text{move}(\alpha, m) \wedge (\text{res} = \text{success} \wedge \mu_1 = 0.8 \vee \\
& \quad \quad \text{res} = \text{failure} \wedge \mu_1 = 0.1) \vee a = \text{move}(\alpha, \text{stand}) \wedge \\
& \quad \quad (\text{res} = \text{success} \wedge \mu_1 = 0.01 \vee \text{res} = \text{failure} \wedge \mu_1 = 0.09)) \wedge \\
& \quad (\text{visible}(\alpha, \beta, a, s) \wedge \mu_2 = 0.7 \vee \text{visible}(\alpha, \text{nil}, a, s) \wedge \mu_2 = 0.1 \vee \\
& \quad \neg \text{visible}(\alpha, \text{nil}, a, s) \wedge \mu_2 = 0.2) \wedge \mu = \mu_1 \cdot \mu_2); \\
& \text{stochastic}(\{\text{moveS}(\alpha, m), \text{moveS}(\alpha', m')\}, s, \{a_\alpha, a_{\alpha'}\}, \{o_\alpha, o_{\alpha'}\}, \mu) \stackrel{\text{def}}{=} \\
& \quad \exists \mu_1, \mu_2 (\text{stochastic}(\{\text{moveS}(\alpha, m)\}, s, \{a_\alpha\}, \{o_\alpha\}, \mu_1) \wedge \\
& \quad \quad \text{stochastic}(\{\text{moveS}(\alpha', m')\}, s, \{a_{\alpha'}\}, \{o_{\alpha'}\}, \mu_2) \wedge \mu = \mu_1 \cdot \mu_2).
\end{aligned}$$

Here,  $\text{visible}(\alpha, \alpha', a, s)$  is true if  $\alpha$  can observe  $\alpha'$  after the execution of  $a$  in  $s$ . The stochastic action  $\text{moveS}(\alpha, m)$  is associated with the observations  $\text{obs}(\beta, \text{res})$ , where  $\beta \in \{\alpha', \text{nil}\}$  and  $\text{res} \in \{\text{success}, \text{failure}\}$ . That is, after the execution of the action  $\text{move}(\alpha, m)$ , agent  $\alpha$  can observe both whether its teammate  $\alpha'$  is present or not (first argument) and the success or failure of the action (second argument). Note that we assume that  $\text{obs}(\alpha', \text{res})$  has the probability zero, if  $\alpha'$  is not visible. Notice also that in the last axiom, we assume the independence of the observations.

As for the optimization theory *OT*, the reward function for the agents is defined by:

$$\begin{aligned}
\text{reward}(\alpha, c, s) = r & \stackrel{\text{def}}{=} \exists \alpha' (\text{goal}(\alpha', \text{do}(c, s)) \wedge (\alpha' = \alpha \wedge r = M \vee \\
& \quad \text{sameTeam}(\alpha, \alpha') \wedge r = M' \vee \neg \text{sameTeam}(\alpha', \alpha) \wedge r = -M)) \vee \\
& \quad \neg \exists \alpha' (\text{goal}(\alpha', \text{do}(c, s)) \wedge \text{evalTeamPos}(\alpha, c, r, s)).
\end{aligned}$$

In this formulation, the reward of agent  $\alpha$  is very high (that is,  $M$  stands for a “big” integer) if  $\alpha$  itself scores a goal, and a bit lower (that is,  $M' < M$ ) if the goal is scored by a teammate. Otherwise, the reward depends on  $\text{evalTeamPos}(\alpha, c, r, s)$ , that is, the position of its team relative to the adversary team as well as the ball possession. The predicate  $\text{goal}(\alpha, s)$  is defined as follows:

$$\text{goal}(\alpha, s) \stackrel{\text{def}}{=} \exists x, y (\text{haveBall}(\alpha, s) \wedge \text{at}(\alpha, x, y, s) \wedge \text{posGoal}(\alpha, x, y)),$$

where  $\text{posGoal}(\alpha, x, y)$  is true iff  $(x, y)$  are the goal coordinates of the adversary team of  $\alpha$ . The predicate  $\text{evalTeamPos}(c, r, s)$  is defined as follows:

$$\begin{aligned}
\text{evalTeamPos}(\alpha, c, r, s) & \stackrel{\text{def}}{=} \exists \alpha', r' ((\text{haveBall}(\alpha', \text{do}(c, s)) \wedge \text{evalPos}(\alpha', r', s) \\
& \quad \wedge (\text{sameTeam}(\alpha, \alpha') \wedge r = r' \vee \neg \text{sameTeam}(\alpha', \alpha) \wedge r = -r')).
\end{aligned}$$

Informally, the value  $r$  in  $evalTeamPos(\alpha, c, r, s)$  depends on the position of the ball owner  $evalPos(\alpha, r, s)$  with respect to the goal, the teammates, and the adversary team.

### 3.2. Belief States

We now introduce belief states along with the executability of actions in belief states and the semantics of actions in terms of transitions between belief states. We also describe how the initial belief state of the agents can be encoded. It should be noted that belief states are not defined as new terms that extend the situation calculus language; instead, they are introduced outside the language to denote probability distributions over situations.

A *belief state (over situations)* has the form  $b = (b_i)_{i \in I}$ , where every  $b_i$  is a finite set of pairs  $(s, \mu)$  consisting of a situation  $s$  and a real  $\mu \in (0, 1]$  such that all  $\mu$  sum up to 1. Informally, every  $b_i$  represents the belief of agent  $i \in I$  expressed as a probability distribution over ordinary situations. The *probability* of a formula  $\phi(s)$  that is uniform in  $s$  in the belief state  $b = (b_i)_{i \in I}$ , denoted  $\phi(b)$ , is the probability vector  $pr = (pr_i)_{i \in I}$ , where every  $pr_i$  with  $i \in I$  is the sum of all  $\mu$  such that  $\phi(s)$  is true and  $(s, \mu) \in b_i$ . Similarly,  $reward(c, b)$  denotes the vector  $r = (r_i)_{i \in I}$ , where every  $r_i$  with  $i \in I$  is the sum of all  $reward(i, c, s) \cdot \mu$  such that  $(s, \mu) \in b_i$ . Intuitively, the expression  $\phi(b)$  provides a probability value for each agent, which is obtained by assessing the situation calculus formula  $\phi(s)$  in the situations mentioned in  $b$ , given the domain theory introduced above.

Notice that, the belief state collects a finite set of pairs composed of (ground) situations and reals; therefore, the expression  $\phi(b)$  may be associated to a formula in the situation calculus. Specifically, the probability vector  $(pr_i)_{i \in I}$  for  $\phi(b)$ , can be denoted through conjunctions and summations as follows:  $\bigwedge_{i \in I} (\bigwedge_{(s_{i,j}, \mu_{i,j}) \in b_i} \psi(s_{i,j}, \mu_{i,j}, p_{i,j}) \wedge pr_i = \sum_j p_{i,j})$  (with  $p_{i,j}$  existentially quantified), where  $\psi(s, \mu, p)$  stands for  $(\phi(s) \wedge p = \mu) \vee (\neg \phi(s) \wedge p = 0)$ . An analogous formula can be introduced to denote the reward vector  $reward(c, b)$ .

**Example 3.5.** (*Rugby Domain cont'd*) Consider the following scenario relative to the domain theory of Example 3.4 (see Fig. 2). We focus on controlling the members of the team  $\mathbf{a}$ , which cooperate to score a goal against the (static) team  $\mathbf{o}$ . The captain  $\mathbf{a}_1$  of  $\mathbf{a}$  has a complete view of the situation, and its belief state  $b_{\mathbf{a}_1}$  coincides with the state shown in Fig. 2, upper part: There is only the situation  $s_1$  with the probability 1 such that  $at(\mathbf{a}_1, 2, 1, s_1)$ ,  $at(\mathbf{a}_2, 2, 4, s_1)$ ,  $at(\mathbf{o}_2, 1, 1, s_1)$ ,  $at(\mathbf{o}_1, 5, 2, s_1)$ , and  $haveBall(\mathbf{a}_1, s_1)$  are all true. That is, the captain  $\mathbf{o}_1$  of  $\mathbf{o}$  is very close to the goal of  $\mathbf{a}$ . From the perspective of  $\mathbf{a}_1$ , its team can score



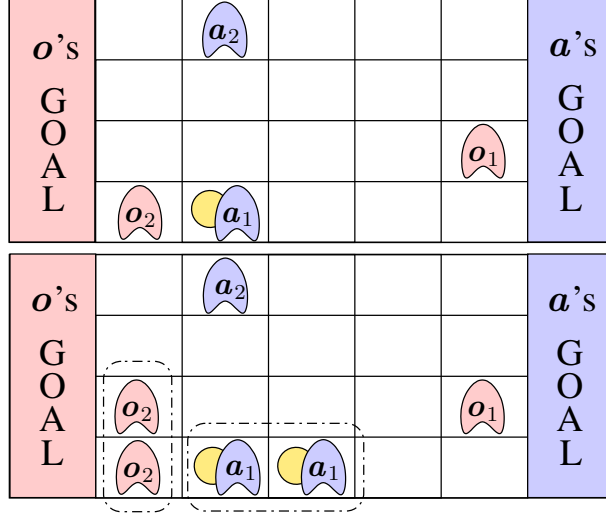


Figure 2: Rugby Domain: Belief states of  $a_1$  and  $a_2$ , respectively.

a goal as follows:  $a_1$  can pass to  $a_2$ , which has a paved way towards the goal. Unfortunately,  $a_1$  has to cooperate with  $a_2$ , whose vision of the situation is more confused (see Fig. 2, lower part): From  $a_2$ 's point of view (that is, belief state  $b_{a_2}$ ),  $o_2$  could be at either (a) (1, 1) or (b) (1, 2), and  $a_1$  could be at either (c) (2, 1) or (d) (3, 1). Hence,  $a_2$ 's belief state  $b_{a_2}$  consists of four possible states with, for example, the following probability distribution:  $\{(s_{a,c}, 0.5), (s_{a,d}, 0.3), (s_{b,c}, 0.1), (s_{b,d}, 0.1)\}$ .

A deterministic action  $c$  is *executable* in a belief state  $b = (b_i)_{i \in I}$  iff  $Poss(c, b) > 0$  (that is,  $c$  is executable in a situation  $s$  such that  $(s, \mu) \in b_i$  for some  $i \in I$ ). A stochastic action  $c$  is *executable* in a belief state  $b = (b_i)_{i \in I}$  producing the observation  $o \in O_c$  iff  $Poss(c_o, b) > 0$  (that is,  $c$  is executable in a situation  $s$  such that  $(s, \mu) \in b_i$  for some  $i \in I$  and  $prob(c, s, n, o) > 0$  for some  $n$ ).

Given a deterministic action  $c$  and a belief state  $b = (b_i)_{i \in I}$ , the *successor belief state* after executing  $c$  in  $b$ , denoted  $do(c, b)$ , is the belief state  $b' = (b'_i)_{i \in I}$ , where

$$b'_i = \{(do(c, s), \mu / Poss(c, b)) \mid (s, \mu) \in b_i, Poss(c, s)\}$$

for every  $i \in I$ . Given a stochastic action  $c$ , an observation  $o \in O_c$ , and a belief state  $b = (b_i)_{i \in I}$ , the *successor belief state* after executing  $c$  in  $b$  and observing  $o$ , denoted  $do(c_o, b)$ , is the belief state  $b' = (b'_i)_{i \in I}$ , where  $b'_i$  is obtained from all pairs

$(do(n, s), \mu \cdot \mu')$  such that  $(s, \mu) \in b_i$ ,  $Poss(c, s)$ , and  $\mu' = prob(c, s, n, o) > 0$  by normalizing the probabilities to sum up to 1. The probability of making the observation  $o \in O_c$  after the execution of the stochastic action  $c$  in  $b = (b_i)_{i \in I}$ , denoted  $prob(c, b, o)$ , is the vector  $pr = (pr_i)_{i \in I}$ , where every probability  $pr_i$  with  $i \in I$  is the sum of all  $\mu \cdot \mu'$  such that  $(s, \mu) \in b_i$  and  $\mu' = prob(c, s, n, o) > 0$ .

**Example 3.6.** (*Rugby Domain cont'd*) The successor belief state after executing the stochastic action  $c = \{moveS(\mathbf{a}_1, m), moveS(\mathbf{a}_2, m)\}$  in the belief state  $b = (\{(S_0, 1)\}, \{(S_0, 1)\})$  and jointly observing the failure resp. success of the action of agent  $\mathbf{a}_1$  resp.  $\mathbf{a}_2$  (that is,  $obs_{\mathbf{a}_1}(failure)$  resp.  $obs_{\mathbf{a}_2}(success)$ ) is given by the following  $b' = (b'_{\mathbf{a}_1}, b'_{\mathbf{a}_2})$  (where  $c_{i,j} = c_i \cup c'_j$  with  $c_0 = \{moveTo(\mathbf{a}_1, stand)\}$ ,  $c_1 = \{moveTo(\mathbf{a}_1, m)\}$ ,  $c'_0 = \{moveTo(\mathbf{a}_2, stand)\}$ , and  $c'_1 = \{moveTo(\mathbf{a}_2, m)\}$ ), and the probabilities follow from the stochastic theory  $ST$  in Example 3.4):

$$b'_{\mathbf{a}_1} = b'_{\mathbf{a}_2} = \{(do(c_{0,0}, S_0), \frac{0.09}{0.1+0.09} \cdot \frac{0.01}{0.8+0.01}), (do(c_{1,0}, S_0), \frac{0.1}{0.1+0.09} \cdot \frac{0.01}{0.8+0.01}), \\ (do(c_{0,1}, S_0), \frac{0.09}{0.1+0.09} \cdot \frac{0.8}{0.8+0.01}), (do(c_{1,1}, S_0), \frac{0.1}{0.1+0.09} \cdot \frac{0.8}{0.8+0.01})\}.$$

We next describe how the initial belief state of the agents can be encoded. Notice that the *basic action theory* (see Section 2.1) provides a unique initial state in  $S_0$ ; however, multiple initial states can be represented by exploiting initial fictitious deterministic actions used to bootstrap them. Let  $b = (b_i)_{i \in I}$  be an initial belief state with  $b_i = \{(s_{i,j}, \mu_{i,j}) \mid j \in \{1, \dots, n_i\}\}$  for all  $i \in I$ . For every  $i \in I$  and  $j \in \{1, \dots, n_i\}$ , we assume a deterministic action  $g_{i,j}$ , which performs a transition from  $S_0$  into the situation  $s_{i,j}$ . For every  $i \in I$ , we then generate  $b_i$  by a stochastic action  $g_i$ , which has every  $g_{i,j}$  (along with the probability  $\mu_{i,j}$ ) such that  $j \in \{1, \dots, n_i\}$  as a deterministic component, and we generate  $b$  by the multi-agent stochastic action  $g = \{g_i \mid i \in I\}$ . Since the deterministic actions  $g_{i,j}$  generate only the possible situations  $s_{i,j}$  in the initial belief state, we assume the precondition axiom  $Poss(g_{i,j}, s) \equiv s = S_0$ . We also want that the other deterministic actions are only executable after  $S_0$ . Hence, for each primitive action, the precondition axiom  $Poss(a, s) \equiv \Psi(a, s)$  is slightly rewritten as  $Poss(a, s) \equiv \Psi(a, s) \wedge s \neq S_0$ . Finally, instead of using  $\mathcal{D}_{S_0}$  to specify the initial situation  $S_0$ , we use it to describe each of the possible situations in the initial belief state.

**Example 3.7.** (*Rugby Domain cont'd*) Consider again the belief state  $b = (b_{\mathbf{a}_1}, b_{\mathbf{a}_2})$  described in Example 3.4, which is given by  $b_{\mathbf{a}_1} = \{(s_1, 1)\}$  and  $b_{\mathbf{a}_2} = \{(s_{a,c}, 0.5), (s_{a,d}, 0.3), (s_{b,c}, 0.1), (s_{b,d}, 0.1)\}$ . In order to encode  $b_{\mathbf{a}_2}$ , we use the deterministic actions  $g_{2;a,c}$ ,  $g_{2;a,d}$ ,  $g_{2;b,c}$ , and  $g_{2;b,d}$  to generate all possible situations

in  $b_{a_2}$ :

$$\begin{aligned} s_{a,c} &= do(g_{2;a,c}, S_0), \quad s_{a,d} = do(g_{2;a,d}, S_0), \\ s_{b,c} &= do(g_{2;b,c}, S_0), \quad s_{b,d} = do(g_{2;b,d}, S_0), \end{aligned}$$

and we use the stochastic action  $g_2$  to associate them with their probabilities in  $b_{a_2}$ :

$$\begin{aligned} stochastic(\{g_2\}, s, \{a\}, \{o\}, \mu) &\stackrel{def}{=} \\ a &= g_{2;a,c} \wedge \mu = 0.5 \vee a = g_{2;a,d} \wedge \mu = 0.3 \vee \\ a &= g_{2;b,c} \wedge \mu = 0.1 \vee a = g_{2;b,d} \wedge \mu = 0.1. \end{aligned}$$

Assuming a similar stochastic action  $g_1$  for  $b_{a_1}$  (associated with a unique deterministic action  $g_{1,1}$  such that  $s_1 = do(g_{1,1}, S_0)$ ), the belief state  $b = (b_{a_1}, b_{a_2})$  is generated by executing the multi-agent stochastic action  $g = \{g_1, g_2\}$  in  $S_0$ :

$$\begin{aligned} stochastic(\{g_1, g_2\}, s, \{a, b\}, \{o_1, o_2\}, \mu) &\stackrel{def}{=} \\ \exists \mu_1, \mu_2 (stochastic(\{g_1\}, s, \{a\}, \{o_1\}, \mu_1) \wedge & \\ stochastic(\{g_2\}, s, \{b\}, \{o_2\}, \mu_2) \wedge \mu = \mu_1 \times \mu_2). & \end{aligned}$$

Notice that the properties of the possible situations in a belief state can be defined as usual. For example, the properties of the ones in  $b_{a_2}$  can be defined by:

$$\begin{aligned} at(\mathbf{a}_1, 2, 1, s_{a,c}) \wedge at(\mathbf{o}_2, 1, 1, s_{a,c}), \quad at(\mathbf{a}_1, 3, 1, s_{a,d}) \wedge at(\mathbf{o}_2, 1, 1, s_{a,d}), \\ at(\mathbf{a}_1, 2, 1, s_{b,c}) \wedge at(\mathbf{o}_2, 1, 2, s_{b,c}), \quad at(\mathbf{a}_1, 3, 1, s_{b,d}) \wedge at(\mathbf{o}_2, 1, 2, s_{b,d}). \end{aligned}$$

### 3.3. Syntax of POGTGolog

Given a domain theory  $DT$ , we define POGTGolog by induction as follows. A *program*  $p$  in POGTGolog has one of the following forms (where  $\alpha$  is a multi-agent action or the empty action  $nop$  (which is always executable and does not change the state of the world),  $\phi$  is a condition,  $p, p_1, p_2, \dots, p_n$  are programs without procedure declarations,  $P_1, \dots, P_n$  are procedure names,  $x, \mathbf{x}_1, \dots, \mathbf{x}_n$  are arguments,  $\tau = \{\tau_1, \tau_2, \dots, \tau_n\}$  is a nonempty finite set of ground terms,  $a_{i,1}, \dots, a_{i,n_i}$  are single-agent actions of agent  $i \in I$ , and  $J \subseteq I$  with  $|J| \geq 2$ ):

1. *Deterministic or stochastic action*:  $\alpha$ . Do  $\alpha$ .
2. *Nondeterministic action choice of agent  $i \in I$* : **choice**( $i: a_{i,1} \mid \dots \mid a_{i,n_i}$ ).  
Do an optimal action (for agent  $i \in I$ ) among  $a_{i,1}, \dots, a_{i,n_i}$ .
3. *Nondeterministic joint action choice*:  $\parallel_{j \in J} \mathbf{choice}(j: a_{j,1} \mid \dots \mid a_{j,n_j})$ .  
Do any action  $\parallel_{j \in J} a_{j,i_j}$  with an optimal probability  $\pi = \prod_{j \in J} \pi_{j,i_j}$ .
4. *Test action*:  $\phi?$ . Test the truth of  $\phi$  in the current situation.

5. *Action sequence*:  $[p_1; p_2]$ . Do  $p_1$  followed by  $p_2$ .
6. *Nondeterministic choice of two programs*:  $(p_1 \mid p_2)$ . Do  $p_1$  or  $p_2$ .
7. *Nondeterministic choice of program argument*:  $\pi[x : \tau](p(x))$ . Do any  $p(\tau_i)$ .
8. *Nondeterministic iteration*:  $p^*$ . Do  $p$  zero or more times.
9. *Conditional*: **if**  $\phi$  **then**  $p_1$  **else**  $p_2$ .
10. *While-loop*: **while**  $\phi$  **do**  $p$ .
11. *Procedures*: **proc**  $P_1(x_1)$   $p_1$  **end**;  $\dots$ ; **proc**  $P_n(x_n)$   $p_n$  **end**;  $p$ .

Hence, compared to Golog, we now also have multi-agent actions and stochastic actions (instead of only primitive resp. deterministic actions). Furthermore, we now additionally have different kinds of nondeterministic action choices for the agents in (2) and (3), where one or any subset of the agents in  $I$  can choose among a finite set of single-agent actions. The formal semantics of (2) and (3) is defined in such a way that an optimal action is chosen for the agents (see Section 3.4). As usual, the sequence operator “;” is associative (that is,  $[[p_1; p_2]; p_3]$  and  $[p_1; [p_2; p_3]]$  have the same meaning), and we often use “ $p_1; p_2$ ” to abbreviate “ $[p_1; p_2]$ ”.

**Example 3.8.** (*Rugby Domain cont’d*) Consider again the scenario (and its belief states  $b_{a_1}$  and  $b_{a_2}$ ) of Example 3.5 relative to the domain theory of Example 3.4 (see Fig. 2). Both agents  $a_1$  and  $a_2$  have to decide when (and if) it is worth to pass the ball, considering that if  $a_1$  tries to pass while  $a_2$  is in offside (for example, in  $s_{a,d}$  or  $s_{b,d}$ ), then the ball goes to the captain  $o_1$  of the adversary team  $o$ , which is in a very good position to score a goal. The subsequent POGTGolog procedure, denoted *schema*, represents a way of acting of  $a_1$  and  $a_2$  in this scenario, where  $a_1$  and  $a_2$  have two possible chances to coordinate themselves to pass the ball (*passTo* and *receive* actions to be jointly executed), alternative actions (*moveS*) can be executed otherwise; thereafter, both of them have to run towards the goal (with or without the ball):

```

proc schema()
  choice( $a_1$ : moveS( $a_1, E$ ) | moveS( $a_1, stand$ ) | moveS( $a_1, passTo(a_2)$ )) ||
    choice( $a_2$ : moveS( $a_2, S$ ) | moveS( $a_2, E$ ) | moveS( $a_2, receive(a_1)$ ));
  choice( $a_1$ : moveS( $a_1, E$ ) | moveS( $a_1, stand$ ) | moveS( $a_1, passTo(a_2)$ )) ||
    choice( $a_2$ : moveS( $a_2, E$ ) | moveS( $a_2, receive(a_1)$ ));
  {moveS( $a_1, E$ ), moveS( $a_2, E$ )};
  {moveS( $a_1, E$ ), moveS( $a_2, E$ )}
end.

```

### 3.4. Policies and Nash Equilibria of POGTGolog

We now define the formal semantics of POGTGolog programs  $p$  relative to a domain theory  $DT$  in terms of Nash equilibria of  $p$ , which are optimal finite-horizon policies of  $p$ . We first associate with every POGTGolog program  $p$ , belief state  $b$ , and horizon  $H \geq 0$ , a set of  $H$ -step policies  $\pi$  along with their expected utility  $U_i$  to every agent  $i \in I$ . We then define the notion of an  $H$ -step Nash equilibrium to characterize a subset of optimal such policies, which is the natural semantics of a POGTGolog program relative to a domain theory.

Intuitively, given a horizon  $H \geq 0$ , an  $H$ -step policy  $\pi$  of a POGTGolog program  $p$  in a belief state  $b$  relative to a domain theory  $DT$  is obtained from the  $H$ -horizon part of  $p$  by replacing every single-agent choice by a single action, and every multi-agent choice by a collection of probability distributions, one over the actions of each agent. Every such  $H$ -step policy  $\pi$  is associated with an expected  $H$ -step reward to  $i \in I$ , an  $H$ -step success probability (which is the probability that  $\pi$  is executable in  $b$ ), and an expected  $H$ -step utility to  $i \in I$  (which is computed from the expected  $H$ -step reward and the  $H$ -step success probability using the utility function).

Formally, for every POGTGolog program  $p$ , the *nil-terminated variant* of  $p$ , denoted  $\hat{p}$ , is inductively defined by  $\hat{p} = [p_1; \hat{p}_2]$ , if  $p = [p_1; p_2]$ , and  $\hat{p} = [p; \text{nil}]$ , otherwise. Given a POGTGolog program  $p$  relative to a domain theory  $DT$ , a horizon  $H \geq 0$ , and a start belief state  $b$ , we say that  $\pi$  is an  $H$ -step policy of  $p$  in  $b$  relative to  $DT$  with *expected  $H$ -step reward*  $v_i$ ,  *$H$ -step success probability*  $pr_i$ , and *expected  $H$ -step utility*  $U_i(H, b, \pi) = \text{utility}(v_i, pr_i)$  to agent  $i \in I$  iff  $DT \models G(\hat{p}, b, H, \pi, (v_i)_{i \in I}, (pr_i)_{i \in I})$ , where the macro  $G(\hat{p}, b, h, \pi, v, pr)$ , for every number of steps to go  $h \in \{0, \dots, H\}$ , is defined by induction as follows ( $\hat{p}$ ,  $b$ , and  $h$  are the input values of  $G$ , while  $\pi$ ,  $v = (v_i)_{i \in I}$ , and  $pr = (pr_i)_{i \in I}$  are the output values of  $G$ ):

- *Null program or zero horizon:*

If  $\hat{p} = \text{nil}$  or  $h = 0$ , then:

$$G(\hat{p}, b, h, \pi, v, pr) \stackrel{\text{def}}{=} \pi = \text{nil} \wedge v = \mathbf{0} \wedge pr = \mathbf{1}.$$

Intuitively,  $p$  ends when it is null or at the horizon end.

- *Deterministic first program action (resp., stochastic first program action with observation):* If  $\hat{p} = [c; p']$ , where  $c$  is a deterministic action (resp.,

stochastic action with observation), and  $h > 0$ , then:

$$G([c; p'], b, h, \pi, v, pr) \stackrel{def}{=} \\ (Poss(c, b) = \mathbf{0} \wedge \pi = stop \wedge v = \mathbf{0} \wedge pr = \mathbf{1}) \vee \\ (Poss(c, b) > \mathbf{0} \wedge \exists \pi', v', pr' (G(p', do(c, b), h-1, \pi', v', pr') \wedge \\ \pi = c; \pi' \wedge v = v' + reward(c, b) \wedge pr = pr' \cdot Poss(c, b))) ,$$

where,  $(s_i)_{i \in I} op (t_i)_{i \in I} = (s_i op t_i)_{i \in I}$  for  $op \in \{+, \cdot\}$ . Informally, suppose that  $\hat{p} = [c; p']$ , where  $c$  is a deterministic action (resp., stochastic action with observation). If  $c$  is not executable in the belief state  $b$ , then  $p$  has only the policy  $\pi = stop$  along with the expected reward  $v = \mathbf{0}$  and the success probability  $pr = \mathbf{0}$ . Here,  $stop$  is a zero-cost action, which takes the agents to an absorbing state, where they stop the execution of the policy and wait for further instructions. Otherwise, the optimal execution of  $[c; p']$  in the belief state  $b$  depends on that one of  $p'$  in  $do(c, b)$ . Observe that  $c$  is executable in  $b$  with the probability  $Poss(c, b)$ , which affects the overall success probability  $pr$ .

- *Stochastic first program action (choice of nature):*

If  $\hat{p} = [c; p']$ , where  $c$  is a stochastic action, and  $h > 0$ , then:

$$G([c; p'], b, h, \pi, v, pr) \stackrel{def}{=} \\ \exists l, \pi_1, \dots, \pi_l, v_1, \dots, v_l, pr_1, \dots, pr_l (\bigwedge_{q=1}^l G([c_{o_q}; p'], b, h, c_{o_q}; \pi_q, \\ v_q, pr_q) \wedge \pi = c; \textbf{for } q = 1 \textbf{ to } l \textbf{ do if } o_q \textbf{ then } \pi_q \wedge \\ v = \sum_{q=1}^l v_q \cdot prob(c, b, o_q) \wedge pr = \sum_{q=1}^l pr_q \cdot prob(c, b, o_q)) ,$$

where  $o_1, \dots, o_l$  are the possible observations. The generated policy consists of  $c$  and a conditional plan in which every such observation  $o_q$  is considered.

- *Nondeterministic first program action (choice of agent  $i \in I$ ):*

If  $\hat{p} = [\textbf{choice}(i: a_1 | \dots | a_n); p']$  and  $h > 0$ , then:

$$G([\textbf{choice}(i: a_1 | \dots | a_n); p'], b, h, \pi, v, pr) \stackrel{def}{=} \\ \exists \pi_1, \dots, \pi_n, v_1, \dots, v_n, pr_1, \dots, pr_n, k (\bigwedge_{q=1}^n G([a_q; p'], b, h, a_q; \pi_q, \\ v_q, pr_q) \wedge k \in \{1, \dots, n\} \wedge \pi = a_k; \textbf{for } q = 1 \textbf{ to } n \textbf{ do if } \psi_q \textbf{ then } \pi_q \wedge \\ v = v_k \wedge pr = pr_k) .$$

Informally, every policy  $\pi$  of  $p$  consists of any action  $a_k$  and one policy  $\pi_q$  of  $p'$  for every possible action  $a_q$ . The expected reward and the success probability of  $\pi$  are given by the expected reward  $v_q$  and the success probability  $pr_q$  of  $\pi_q$ . For the other agents to observe which action among  $a_1, \dots, a_n$  was actually executed by agent  $i$ , we use a cascade of if-then-else statements with conditions  $\psi_q$ .

- *Nondeterministic first program action (joint choice of the agents in  $J$ ):*  
If  $\widehat{p} = [\|_{j \in J} \mathbf{choice}(j: a_{j,1} \mid \dots \mid a_{j,n_j}); p']$  and  $h > 0$ , then:

$$\begin{aligned} G([\|_{j \in J} \mathbf{choice}(j: a_{j,1} \mid \dots \mid a_{j,n_j}); p'], b, h, \pi, v, pr) &\stackrel{def}{=} \\ \exists \pi_a (a \in A), v_a (a \in A), pr_a (a \in A), \pi_j (j \in J) &(\bigwedge_{a \in A} G([\bigcup_{j \in J} a_j; p'], \\ b, h, \bigcup_{j \in J} a_j; \pi_a, v_a, pr_a) \wedge \bigwedge_{j \in J} &(\pi_j \in PD(\{a_{j,1}, \dots, a_{j,n_j}\})) \wedge \\ \pi = \prod_{j \in J} \pi_j; \mathbf{for\ each\ } a \in A \mathbf{ do\ if\ } \phi_a &\mathbf{ then\ } \pi_a \wedge \\ v = \sum_{a \in A} v_a \cdot \prod_{j \in J} \pi_j(a_j) \wedge pr = &\sum_{a \in A} pr_a \cdot \prod_{j \in J} \pi_j(a_j)), \end{aligned}$$

where  $A = \times_{j \in J} \{a_{j,1}, \dots, a_{j,n_j}\}$ . We denote by  $PD(S)$  the set of all probability distributions over  $S$ , and  $(\prod_{j \in J} \pi_j)(a) = \prod_{j \in J} \pi_j(a_j)$  for all  $a = (a_j)_{j \in J}$ . Informally, every policy  $\pi$  of  $p$  consists of one probability distribution  $\pi_j$  over  $\{a_{j,1}, \dots, a_{j,n_j}\}$  for every agent  $j \in J$ , and one policy  $\pi_a$  of  $p'$  for every possible joint action  $a \in A$ . The expected reward and the success probability of  $\pi$  are given by the expected reward and the expected success probability of the policies  $\pi_a$ . Here,  $\pi_j$  specifies the probabilities with which agent  $j \in J$  should execute the actions  $\{a_{j,1}, \dots, a_{j,n_j}\}$ . Hence, assuming the usual probabilistic independence between the distributions  $\pi_j$  with  $j \in J$  in stochastic games, every possible joint action  $a$  is executed with the probability  $(\prod_{j \in J} \pi_j)(a)$ . Note that the conditions  $\phi_a$  with  $a \in A$  are to observe what the agents have actually executed.

- *Test action:*

If  $\widehat{p} = [\phi?; p']$  and  $h > 0$ , then:

$$\begin{aligned} G([\phi?; p'], b, h, \pi, v, pr) &\stackrel{def}{=} (\phi[b] = \mathbf{0} \wedge \pi = stop \wedge v = \mathbf{0} \wedge pr = \mathbf{0}) \vee \\ \exists pr' (\phi[b] > \mathbf{0} \wedge G(p', b, h, \pi, v, pr') &\wedge pr = pr' \cdot \phi[b]). \end{aligned}$$

Informally, let  $\widehat{p} = [\phi?; p']$ . If  $\phi$  is false in  $b$ , then  $p$  has only the policy  $\pi = stop$ , the expected reward  $v = \mathbf{0}$ , and the success probability  $pr = \mathbf{0}$ . Otherwise,  $\pi$  is a policy of  $p$  with the expected reward  $v$  and success probability

$pr' \cdot \phi[b]$  iff  $\pi$  is a policy of  $p'$  with the expected reward  $v$  and success probability  $pr'$ .

- *Nondeterministic choice of two programs:*

If  $\hat{p} = [(p_1 \mid p_2); p']$  and  $h > 0$ , then:

$$\begin{aligned} G([(p_1 \mid p_2); p'], b, h, \pi, v, pr) &\stackrel{def}{=} \\ \exists \pi_1, \pi_2, v_1, v_2, pr_1, pr_2, k \ (\bigwedge_{q \in \{1,2\}} G([p_q; p'], b, h, \pi_q, v_q, pr_q) \wedge \\ k \in \{1, 2\} \wedge \pi = \pi_k \wedge v = v_k \wedge pr = pr_k) . \end{aligned}$$

- *Conditional:*

If  $\hat{p} = [\text{if } \phi \text{ then } p_1 \text{ else } p_2; p']$  and  $h > 0$ , then:

$$\begin{aligned} G([\text{if } \phi \text{ then } p_1 \text{ else } p_2; p'], b, h, \pi, v, pr) &\stackrel{def}{=} \\ G([( [\phi?; p_1] \mid [\neg\phi?; p_2] ); p'], b, h, \pi, v, pr) . \end{aligned}$$

This case is reduced to test action and nondeterministic choice of two programs.

- *While-loop:*

If  $\hat{p} = [\text{while } \phi \text{ do } p; p']$  and  $h > 0$ , then:

$$G([\text{while } \phi \text{ do } p; p'], b, h, \pi, v, pr) \stackrel{def}{=} G([( [\phi?; p]^* ; \neg\phi? ); p'], b, h, \pi, v, pr) .$$

This case is reduced to test action and nondeterministic iteration.

- *Nondeterministic choice of program argument:*

If  $\hat{p} = [\pi[x : \tau](p(x)); p']$ , where  $\tau = \{\tau_1, \tau_2, \dots, \tau_n\}$ , and  $h > 0$ , then:

$$\begin{aligned} G([\pi[x : \tau](p(x)); p'], b, h, \pi, v, pr) &\stackrel{def}{=} \\ G([( \dots (p(\tau_1) \mid p(\tau_2)) \mid \dots \mid p(\tau_n) ); p'], b, h, \pi, v, pr) . \end{aligned}$$

This case is reduced to nondeterministic choice of two programs.

- *Nondeterministic iteration:*

If  $\hat{p} = [p^*; p']$  and  $h > 0$ , then:

$$\begin{aligned} G([p^*; p'], b, h, \pi, v, pr) &\stackrel{def}{=} \\ G([\text{proc } nit \ (nop \mid [p; nit]) \ \text{end}; nit]; p'], b, h, \pi, v, pr) . \end{aligned}$$

This case is reduced to procedures and nondeterministic choice of two programs.



- *Procedures:* We consider the two cases of (1) handling procedure declarations and (2) handling procedure calls. To this end, we slightly extend the first argument of  $G$  with a store for procedure declarations, which can be safely ignored in all the above constructs of POGTGolog.

(1) If  $\hat{p} = [\mathbf{proc } P_1(\mathbf{x}_1) p_1 \mathbf{end}; \dots; \mathbf{proc } P_n(\mathbf{x}_n) p_n \mathbf{end}; p] \langle \rangle$  and  $h > 0$ , then:

$$G([\mathbf{proc } P_1(\mathbf{x}_1) p_1 \mathbf{end}; \dots; \mathbf{proc } P_n(\mathbf{x}_n) p_n \mathbf{end}; p] \langle \rangle, b, h, \pi, v, pr) \stackrel{def}{=} G([p] \langle \mathbf{proc } P_1(\mathbf{x}_1) p_1 \mathbf{end}; \dots; \mathbf{proc } P_n(\mathbf{x}_n) p_n \mathbf{end} \rangle, b, h, \pi, v, pr).$$

We store the procedure declarations at the end of the first argument of  $G$ .

(2) If  $\hat{p} = [P_i(\mathbf{x}_i); p'] \langle d \rangle$  and  $h > 0$ , then:

$$G([P_i(\mathbf{x}_i); p'] \langle d \rangle, b, h, \pi, v, pr) \stackrel{def}{=} G([p_d(P_i(\mathbf{x}_i)); p'] \langle d \rangle, b, h, \pi, v, pr).$$

We replace a procedure call  $P_i(\mathbf{x}_i)$  with its code (body)  $p_d(P_i(\mathbf{x}_i))$ , which is encoded in  $d$ .

We are now ready to define the notion of an  $H$ -step Nash equilibrium as follows. An  $H$ -step policy  $\pi$  of a POGTGolog program  $p$  in a belief state  $b$  relative to  $DT$  is an  $H$ -step Nash equilibrium of  $p$  in  $b$  relative to  $DT$  iff, for every agent  $i \in I$ , it holds that  $U_i(H, b, \pi') \leq U_i(H, b, \pi)$  for all  $H$ -step policies  $\pi'$  of  $p$  in  $b$  relative to  $DT$  obtained from  $\pi$  by modifying only actions of agent  $i$ .

**Example 3.9.** (*Rugby Domain cont'd*) Consider again the scenario (and its belief states  $b_{a_1}$  and  $b_{a_2}$ ) of Example 3.5 relative to the domain theory of Example 3.4 (see Fig. 2). Assuming the horizon  $H = 4$ , a 4-step policy  $\pi$  of the POGTGolog program *schema* of Example 3.8 is given by:

$$DT \models G([schema; nil], (b_{a_1}, b_{a_2}), 4, \pi, (v_1, v_2), (pr_1, pr_2)).$$

For agent  $a_1$ , an optimal way of acting is to pass the ball as soon as possible, which can be encoded by the following (pure) 4-step policy  $\pi_{a_1} = c; \pi_{a_1}^1$ , where  $c = \{moveS(a_1, passTo(a_2)), moveS(a_2, receive(a_1))\}$ , and  $\pi_{a_1}^1$  is an optimal 3-step policy of *schema'* in the belief state  $(do(c, b_{a_1}), do(c, b_{a_2}))$ . The procedure *schema'* is obtained from *schema* by removing the first nondeterministic joint action choice. The policy  $\pi_{a_1}^1$  gives to agent  $a_2$  three  $moveS(a_2, E)$  attempts to achieve the touch-line. From the standpoint of  $a_2$ , instead, it is worth to do a

$moveS(a_2, S)$  to observe if agent  $a_1$  is aligned, trying to minimize the likelihood of a wrong passing. In this case,  $a_1$  has to delay the passing, waiting for the move of  $a_2$ . The resulting (pure) 4-step policy  $\pi_{a_2}$  is more favorable to  $a_2$ 's belief state:

$$\begin{aligned} \pi_{a_2} = c; & \text{ if } obs(a_1, success) \text{ then } \pi_{a_2}^{1,o_1} \\ & \text{ else if } obs(a_1, failure) \text{ then } \pi_{a_2}^{1,o_2} \\ & \text{ else if } obs(nil, success) \text{ then } \pi_{a_2}^{2,o_3} \\ & \text{ else if } obs(nil, failure) \text{ then } \pi_{a_2}^{2,o_4}, \end{aligned}$$

where  $c = \{moveS(a_1, S), moveS(a_2, stand)\}$ , and  $\pi_{a_2}^{k,o_i}$  is an optimal 3-step policy of  $schema'$ , when observing  $o_i$  after executing  $c$  in  $(b_{a_1}, b_{a_2})$ . Given this conflict of opinions, an optimal compromise for both  $a_1$  and  $a_2$  is a Nash equilibrium.

#### 4. A POGTGolog Interpreter

In this section, we first define an interpreter for POGTGolog programs, and we then provide some optimality and faithfulness results for the interpreter.

##### 4.1. Formal Specification

We now define an interpreter for POGTGolog programs  $p$  relative to a domain theory  $DT$  by specifying the macro  $OptG(\hat{p}, b, H, \pi, v, pr)$ , which takes as input the *nil*-terminated variant  $\hat{p}$  of a POGTGolog program  $p$ , a belief state  $b = (b_i)_{i \in I}$ , and a finite horizon  $H \geq 0$ , and which computes as output an optimal policy, that is, a  $H$ -step policy  $\pi$  which is a  $H$ -step Nash equilibrium of  $p$  in  $b$  (see Theorem 4.1). Such policy is associated with the vectors  $v = (v_i)_{i \in I}$  and  $pr = (pr_i)_{i \in I}$ , where each  $v_i$  is the expected  $H$ -step reward of  $\pi$  in  $b$  to  $i$ , each  $pr_i \in [0, 1]$  is the  $H$ -step success probability of  $\pi$  in  $b$  for  $i$ , and  $utility(v_i, pr_i)$  is the expected  $H$ -step utility of  $\pi$  in  $b$  to  $i$ . Notice that, the macro  $G$  of Section 3.4 defines all the legal  $H$ -step policies allowed by a POGTGolog program  $p$ ; therefore, any possible action choice compatible with the program specification and the domain theory is enabled. In contrast, the macro  $OptG$  defines an optimal such legal one, that is, a  $H$ -step Nash equilibrium.

For this purpose, we define the macro  $OptG(\hat{p}, b, h, \pi, v, pr)$  in nearly the same way as the macro  $G(\hat{p}, b, h, \pi, v, pr)$  in Section 3.4, except for the following modifications:

- *Nondeterministic first program action (choice of agent  $i \in I$ ):* The definition of  $OptG$  is obtained from the one of  $G$  by replacing the condition “ $k \in \{1,$

$\dots, n\}$ ” by the condition “ $k = \operatorname{argmax}_{q \in \{1, \dots, n\}} \operatorname{utility}(v_{q,i}, pr_{q,i})$ ”, where  $v_q = (v_{q,i})_{i \in I}$  and  $pr_q = (pr_{q,i})_{i \in I}$ . Informally, given the possible actions  $a_1, \dots, a_n$  for agent  $i \in I$ , we select an optimal one for  $i$ , that is, one with greatest  $\operatorname{utility}(v_{q,i}, pr_{q,i})$ .

- *Nondeterministic first program action (joint choice of the agents in  $J$ ):* The definition of  $\operatorname{Opt}G$  is obtained from the one of  $G$  by replacing “ $\bigwedge_{j \in J} (\pi_j \in PD(\{a_{j,1}, \dots, a_{j,n_j}\}))$ ” by “ $(\pi_j)_{j \in J} = \operatorname{selectNash}(\{\operatorname{utility}(v_a, pr_a)|_J \mid a \in A\})$ ”, where  $\operatorname{utility}((s_i)_{i \in I}, (t_i)_{i \in I}) = (\operatorname{utility}(s_i, t_i))_{i \in I}$ , and  $s|_J$  is the restriction of  $s$  to  $J$ , for  $s = (s_i)_{i \in I}$  and  $J \subseteq I$ . Informally, we compute a local Nash equilibrium  $(\pi_j)_{j \in J}$  from a normal form game using the Nash selection function  $\operatorname{selectNash}$ . Note that we assume that all agents have the same Nash selection functions, and thus they automatically select a common unique Nash equilibrium. Such Nash selection function can be implemented by any Nash equilibrium solver for normal form games (see, for instance, [46]), provided that it produces the same solution for the same problem instance.
- *Nondeterministic choice of two programs:* The definition of  $\operatorname{Opt}G$  is obtained from the one of  $G$  by replacing “ $k \in \{1, 2\}$ ” by “ $k = \operatorname{argmax}_{q \in \{1, 2\}} \operatorname{utility}(v_{q,j}, pr_{q,j})$ ”. Informally, given two possible programs  $p_1$  and  $p_2$ , we select an optimal one for agent  $j$ , that is, one with greatest  $\operatorname{utility}(v_{q,j}, pr_{q,j})$ .

#### 4.2. Optimality and Faithfulness

The following theorem shows the important result that the macro  $\operatorname{Opt}G$  captures an optimal policy (that is, a Nash equilibrium) among the policies admitted by a POGTGolog program (that is, the ones defined by the macro  $G$ ). More specifically, for every horizon  $H \geq 0$ , among the set of all  $H$ -step policies  $\pi$  of a POGTGolog program  $p$  relative to a domain theory  $DT$  in a belief state  $b$ , the macro  $\operatorname{Opt}G$  defines an  $H$ -step Nash equilibrium and its expected  $H$ -step utility.

**Theorem 4.1.** *Let  $DT = (AT, ST, OT)$  be a domain theory, and let  $p$  be a POGTGolog program relative to  $DT$ . Let  $b$  be a belief state, let  $H \geq 0$  be a horizon, and let  $DT \models \operatorname{Opt}G(\hat{p}, b, H, \pi, v, pr)$ . Then,  $\pi$  is an  $H$ -step Nash equilibrium of  $p$  in  $b$  relative to  $DT$ , and  $\operatorname{utility}(v_i, pr_i)$  is its expected  $H$ -step utility to agent  $i \in I$ .*

The following theorem shows that POGTGolog programs faithfully extend POSGs, that is, any POSG can be encoded as a POGTGolog program, which is semantically interpreted as that POSG. More formally, the theorem says that given any  $H \geq 0$ , every POSG can be encoded as a program  $p$  in POGTGolog such that  $OptG$  specifies one of its  $H$ -step Nash equilibria and its expected  $H$ -step reward.

**Theorem 4.2.** *Let  $G = (I, Z, (A_i)_{i \in I}, (O_i)_{i \in I}, P, (R_i)_{i \in I})$  be a POSG, and let  $H \geq 0$  be a horizon. Then, there exists a domain theory  $DT = (AT, ST, OT)$ , and a set of POGTGolog programs  $\{p^h \mid h \in \{0, \dots, H\}\}$  relative to  $DT$  such that  $\sigma = (\sigma_i)_{i \in I}$  is an  $H$ -step Nash equilibrium of  $G$ , where every  $(\sigma_i(b, h))_{i \in I} = (\pi_i)_{i \in I}$  is given by  $DT \models OptG(\hat{p}^h, B_b, h+1, \Pi_{i \in I} \pi_i; \pi', v, pr)$ , for every belief state  $b$  of  $G$  and every  $h \in \{0, \dots, H\}$ , with  $B_b$  being a belief state of  $DT$  associated with  $b$ . Furthermore, the expected  $H$ -step reward of  $\sigma$  in  $b$  to agent  $i \in I$  is given by  $utility(v_i, pr_i)$ , where  $DT \models OptG(\hat{p}^H, B_b, H+1, \pi, v, pr)$ , for every belief state  $b$  of  $G$ .*

## 5. Example

In this section, we provide an extended example to illustrate the overall framework at work. This example is inspired by the stratagus domain in [42].

**Example 5.1. (Stratagus Domain)** The stratagus field consists of  $9 \times 9$  positions (see Fig. 3). We assume a team of two agents  $\mathbf{a} = \{\mathbf{a}_1, \mathbf{a}_2\}$ , which occupy one position each. The stratagus field has designated areas representing two *gold-mines*, one *forest*, and one *base* for each agent (see Fig. 3). The two agents can move one step in one of the directions  $N, S, E$ , and  $W$ , or remain stationary. Each of the two agents can also pick up one unit of wood (resp., gold) at the forest (resp., gold-mines), and drop these resources at the base. We assume that if  $\mathbf{a}_1$  and  $\mathbf{a}_2$  are in the same location, then they cannot pick up anything. Each action of the two agents can fail, resulting in a stationary move. Any carried object drops when the two agents collide. After each step, the agents  $\mathbf{a}_1$  and  $\mathbf{a}_2$  receive the rewards  $r_{\mathbf{a}_1}$  and  $r_{\mathbf{a}_2}$ , respectively, where  $r_k$  for  $k \in \{\mathbf{a}_1, \mathbf{a}_2\}$  is 0, 1, and 2 when  $k$  brings nothing, one unit of wood, and one unit of gold to its base, respectively.

We define the domain theory  $DT = (AT, ST, OT)$  as follows. As for the basic action theory  $AT$ , we assume the deterministic actions  $move(\alpha, m)$  (agent  $\alpha$  performs  $m$  among  $N, S, E, W$ , and  $stand$ ),  $pickUp(\alpha, o)$  (agent  $\alpha$  picks up the object  $o$ ), and  $drop(\alpha, o)$  (agent  $\alpha$  drops the object  $o$ ), as well as the relational fluents  $at(q, x, y, s)$  (agent or object  $q$  is at the position  $(x, y)$  in the situation  $s$ ),

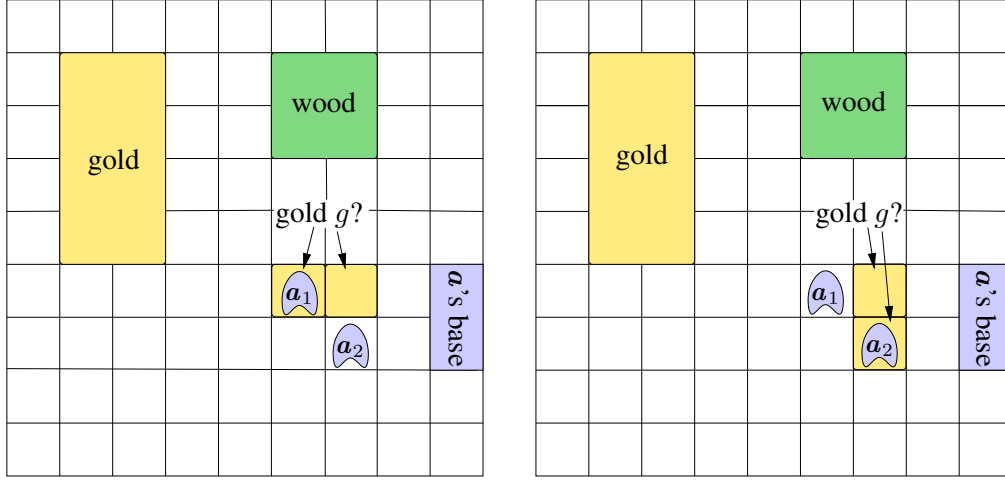


Figure 3: Stratagus Domain: Initial belief states of  $a_1$  and  $a_2$ , respectively.

and  $holds(\alpha, o, s)$  (agent  $\alpha$  holds the object  $o$  in the situation  $s$ ), which are defined through the following successor state axioms:

$$\begin{aligned}
at(q, x, y, do(c, s)) &\equiv agent(q) \wedge (at(q, x, y, s) \wedge move(q, stand) \in c \vee \\
&\quad \exists x', y' (at(q, x', y', s) \wedge \exists m (move(\alpha, m) \in c \wedge \phi(x, y, x', y', m)))) \vee \\
&\quad object(q) \wedge (at(q, x, y, s) \wedge \neg \exists \alpha (pickUp(\alpha, q) \in c) \vee \\
&\quad \exists \alpha ((drop(\alpha, q) \in c \vee collision(c, s)) \wedge at(\alpha, x, y, s) \wedge holds(\alpha, q, s))) ; \\
holds(\alpha, o, do(c, s)) &\equiv holds(\alpha, o, s) \wedge drop(\alpha, o) \notin c \wedge \\
&\quad \neg collision(c, s) \vee pickUp(\alpha, o) \in c.
\end{aligned}$$

$\phi(x, y, x', y', m)$  is true iff the coordinates change from the position  $(x', y')$  to the position  $(x, y)$  due to  $m \in \{N, S, E, W, stand\}$ , that is,

$$\begin{aligned}
\phi(x, y, x', y', m) &\stackrel{def}{=} (m \notin \{N, S, E, W\} \wedge x = x' \wedge y = y') \vee \\
&\quad (m = N \wedge x = x' \wedge y = y' + 1) \vee (m = S \wedge x = x' \wedge y = y' - 1) \vee \\
&\quad (m = E \wedge x = x' + 1 \wedge y = y') \vee (m = W \wedge x = x' - 1 \wedge y = y'),
\end{aligned}$$

and  $collision(c, s)$  encodes that executing the action  $c$  in the situation  $s$  causes a

collision between the agents  $\alpha_1$  and  $\alpha_2$  in the situation  $do(c, s)$ , that is,

$$\begin{aligned} collision(c, s) \stackrel{def}{=} & \exists \alpha, \beta, x, y (\alpha \neq \beta \wedge \\ & \exists x', y' (at(\alpha, x', y', s) \wedge \exists m (move(\alpha, m) \in c \wedge \phi(x, y, x', y', m))) \wedge \\ & \exists x'', y'' (at(\beta, x'', y'', s) \wedge \exists m (move(\beta, m) \in c \wedge \phi(x, y, x'', y'', m))) \wedge \\ & (x' \neq x \vee y' \neq y \vee x'' \neq x \vee y'' \neq y)) . \end{aligned}$$

That is, we have a collision between two agents iff (i) at least one them moves, and (ii) they are in the same location thereafter. The deterministic actions  $move(\alpha, m)$ ,  $drop(\alpha, o)$ , and  $pickUp(\alpha, o)$  have the following precondition axioms:

$$\begin{aligned} Poss(move(\alpha, m), s) & \equiv \neg \exists x, y (at(\alpha, x, y, s) \wedge ((y = 9 \wedge m = N) \vee \\ & (y = 1 \wedge m = S) \vee (x = 9 \wedge m = E) \vee (x = 1 \wedge m = W))) ; \\ Poss(drop(\alpha, o), s) & \equiv holds(\alpha, o, s) ; \\ Poss(pickUp(\alpha, o), s) & \equiv \neg \exists x holds(\alpha, x, s) \wedge \exists x, y (at(\alpha, x, y, s) \wedge at(o, x, y, s)) . \end{aligned}$$

The first axiom forbids  $\alpha$  to go out of the  $9 \times 9$  game-field, while the second axiom states that  $\alpha$  can only drop the object  $o$  if  $\alpha$  is holding  $o$ , and the third axiom permits  $\alpha$  to pick up  $o$  if  $\alpha$  is at same location and not holding anything else.

As for the stochastic theory  $ST$ , we assume the stochastic actions  $moveS(\alpha, m)$ ,  $pickUpS(\alpha, o)$ , and  $dropS(\alpha, o)$ , which are specified below. Each execution of such an action  $c$  is followed by an associated observation, either  $obs_\alpha(success)$  or  $obs_\alpha(failure)$ . Intuitively,  $\alpha$  can observe if the execution of  $\alpha$  was successful or not.

$$\begin{aligned} stochastic(\{moveS(\alpha, m)\}, s, \{a\}, \{o\}, \mu) & \stackrel{def}{=} o = obs_\alpha(success) \wedge \\ & (a = move(\alpha, m) \wedge \mu = 0.54 \vee a = move(\alpha, stand) \wedge \mu = 0.36) \vee \\ & o = obs_\alpha(failure) \wedge a = move(\alpha, stand) \wedge \mu = 0.1 ; \\ stochastic(\{pickUpS(\alpha, o)\}, s, \{a\}, \{o\}, \mu) & \stackrel{def}{=} o = obs_\alpha(success) \wedge \\ & (a = pickUp(\alpha, o) \wedge \mu = 0.72 \vee a = move(\alpha, stand) \wedge \mu = 0.18) \vee \\ & o = obs_\alpha(failure) \wedge a = move(\alpha, stand) \wedge \mu = 0.1 ; \\ stochastic(\{dropS(\alpha, o)\}, s, \{a\}, \{o\}, \mu) & \stackrel{def}{=} o = obs_\alpha(success) \wedge \\ & (a = drop(\alpha, o) \wedge \mu = 0.81 \vee a = move(\alpha, stand) \wedge \mu = 0.09) \vee \\ & o = obs_\alpha(failure) \wedge a = move(\alpha, stand) \wedge \mu = 0.1 . \end{aligned}$$

Observe that the observation  $obs_\alpha(failure)$  is reliable, that is, if agent  $\alpha$  observes  $obs_\alpha(failure)$ , then its executed action was not successful with the probability 1.

Multi-agent stochastic actions are defined as follows (assuming independence):

$$\begin{aligned} stochastic(\{moveS(\alpha, m), moveS(\alpha', m')\}, s, \{a_\alpha, a_{\alpha'}\}, \{o_\alpha, o_{\alpha'}\}, \mu) &\stackrel{def}{=} \\ \exists \mu_1, \mu_2 (stochastic(\{moveS(\alpha, m)\}, s, \{a_\alpha\}, \{o_\alpha\}, \mu_1) \wedge & \\ stochastic(\{moveS(\alpha', m')\}, s, \{a_{\alpha'}\}, \{o_{\alpha'}\}, \mu_2) \wedge \mu = \mu_1 \cdot \mu_2) . & \end{aligned}$$

As for the optimization theory *OT*, we use the product as the utility function *utility*. Furthermore, we define the reward function for agent  $\alpha$  as follows:

$$\begin{aligned} reward(\alpha, c, s) = r &\stackrel{def}{=} \exists r_\alpha, r_\beta (rewardAct(\alpha, c, s) = r_\alpha \wedge \\ rewardAct(\beta, c, s) = r_\beta \wedge \alpha \neq \beta \wedge r = r_\alpha + 0.5 \cdot r_\beta) . & \end{aligned}$$

$rewardAct(\alpha, c, s)$  is here defined as follows:

$$\begin{aligned} rewardAct(\alpha, c, s) = r &\stackrel{def}{=} \exists o, x, y (at(\alpha, x, y, s) \wedge holding(\alpha, o, s) \wedge \\ \neg collision(c, s) \wedge base(\alpha, x, y) \wedge drop(\alpha, o) \in c \wedge (gold(o) \wedge r = 2 \vee & \\ wood(o) \wedge r = 1)) \vee \neg \exists o, x, y (at(\alpha, x, y, s) \wedge holding(\alpha, o, s) \wedge & \\ \neg collision(c, s) \wedge base(\alpha, x, y) \wedge drop(\alpha, o) \in c) \wedge r = 0 . & \end{aligned}$$

Note that the reward to agent  $\alpha$  is higher if  $\alpha$  itself can drop an object to the base. Thus, even though there is a joint interest in bringing objects to the base, each agent prefers to be the one who achieves the goal.

Consider the scenario shown in Fig. 3, where the two agents  $a_1$  and  $a_2$  are looking for a unit of gold  $g$ , and they are trying to bring it to their base. Suppose that the initial belief state of agent  $a_1$  (resp.,  $a_2$ ) is as in Fig. 3, left (resp., right) side. In particular, agent  $a_1$  (resp.,  $a_2$ ) initially believes that the unit of gold  $g$  is at either (6, 4) or (7, 4) (resp., (7, 3) or (7, 4)). Formally, let the belief state of agent  $a_1$  (resp.,  $a_2$ ) be given by  $b_{a_1} = \{(s_{6,4}, 0.2), (s_{7,4}, 0.8)\}$  (resp.,  $b_{a_2} = \{(s_{7,3}, 0.4), (s_{7,4}, 0.6)\}$ ), where every included  $s_{i,j}$  denotes one of the three possible situations shown in Fig. 3, and  $at(g, i, j, s_{i,j})$  is true. How should the two agents  $a_1$  and  $a_2$  now act in such an initial situation? Both agents believe that the unit of gold  $g$  could be in their position or in the position (7, 4). Thus, the probability of finding the unit of gold  $g$  in the position (7, 4) is higher, but if both agents decide to go there, then they could block each other's pick-up actions, since joint pick-up actions in the same position are not allowed. However, once one of the two agents  $a_1$  and  $a_2$  has caught the unit of gold  $g$ , it should go to its base and drop it. The following POGTGolog program, denoted *schema*, encodes

a possible way of acting of  $a_1$  and  $a_2$  in this scenario:

```

proc schema( $n$ )
  pickOrGo( $n$ );
   $\pi a, o (holding(a, o)? ; carryToBase(a))$ 
end.

```

In this program, each agent tries to pick up an object, solving a “pick up or go” dilemma when they aim at picking up the same object (through *pickOrGo*). Once an object is gathered, the agent has to find a way to bring it to the base (through *carryToBase*). The procedure *pickOrGo*( $n$ ) is defined as follows:

```

proc pickOrGo( $n$ )
if  $n > 0 \wedge \neg \exists \alpha, o (holding(\alpha, o))$  then [
   $\pi d_1, p_2, o_1, o_2 ((direction(p_1) \wedge direction(p_2) \wedge object(o_1) \wedge object(o_2))?) ;$ 
  choice( $a_1 : moveS(a_1, d_1) \mid moveS(a_1, stand) \mid pickUpS(a_1, o_1)$ )  $\parallel$ 
  choice( $a_2 : moveS(a_2, d_2) \mid moveS(a_2, stand) \mid pickUpS(a_2, o_2)$ );
  pickOrGo( $n-1$ )]
end.

```

The following procedure *carryToBase*( $a$ ) describes a partially specified behavior where the agent  $a$  is trying to move to its base to drop down an object:

```

proc carryToBase( $a$ )
choice( $a : moveS(a, N) \mid moveS(a, S) \mid moveS(a, E) \mid moveS(a, W)$ );
if atBase then  $\pi x (dropS(a, x))$ 
  else carryToBase( $a$ )
end.

```

Informally, in the program *schema*, the agents  $a_1$  and  $a_2$  first have to decide whether to move towards the most probable gold location, or to remain in their position, or to try to pick up the gold (through the *pickOrGo* procedure). We assume that they can try this maximally  $n$  times. Once a joint action of *pickOrGo* is executed, if one of the two agents  $a_1$  and  $a_2$  holds the gold, then it can start to move towards the base (through the *carryToBase* procedure).

The agent behavior is partially specified by the procedure *schema*( $n$ ), which is fully instantiated by the program interpreter given the current belief states of the two agents. Given the context in Fig. 3, we now focus on the possible instances of *schema*(2), that is, the agents can spend two attempts to pick up the object. From  $a_1$ 's point of view (Fig. 3, left side), a good policy  $\pi_{a_1}$  of *schema*(2) may



be to start with the joint action  $c = \{pickUpS(\mathbf{a}_1, g), moveS(\mathbf{a}_2, stand)\}$ , and then to act depending on  $\mathbf{a}_1$ 's observation:

$$\begin{aligned} \pi_{\mathbf{a}_1} = c; & \text{ if } obs_{\mathbf{a}_1}(success) \wedge obs_{\mathbf{a}_2}(success) \text{ then } \pi_{\mathbf{a}_1}^1 \\ & \text{ else if } obs_{\mathbf{a}_1}(success) \wedge obs_{\mathbf{a}_2}(failure) \text{ then } \pi_{\mathbf{a}_1}^2 \\ & \text{ else if } obs_{\mathbf{a}_1}(failure) \wedge obs_{\mathbf{a}_2}(success) \text{ then } \pi_{\mathbf{a}_1}^3 \\ & \text{ else if } obs_{\mathbf{a}_1}(failure) \wedge obs_{\mathbf{a}_2}(failure) \text{ then } \pi_{\mathbf{a}_1}^4. \end{aligned}$$

Indeed, acting in this way,  $\mathbf{a}_1$  could try to get its gold and then bring it to the base, avoiding a collision with  $\mathbf{a}_2$ . On the other hand, from  $\mathbf{a}_2$ 's point of view (Fig. 3, right side), a good policy  $\pi_{\mathbf{a}_2}$  of *schema* may be to start with the joint action  $c = \{moveS(\mathbf{a}_1, E), pickUpS(\mathbf{a}_2, g)\}$ . In fact, in this way, both  $\mathbf{a}_1$  and  $\mathbf{a}_2$  could try to pick up some gold and then move towards the base. However, the situation here is very complex. Indeed, each action could fail, and each failure determines a different context. The complete policy is generated by the POGTGolog interpreter. Assuming the horizon  $H = 5$ , an optimal 5-step joint policy  $\pi$  is given by  $DT \models OptG([schema(2); nil], (b_{\mathbf{a}_1}, b_{\mathbf{a}_2}), 5, \pi, (v_1, v_2), (pr_1, pr_2))$ . Since both  $\mathbf{a}_1$  and  $\mathbf{a}_2$  know each others' initial belief states, and they are both endowed with the same POGTGolog program *schema*, the interpreter of each of them independently produces the same joint 5-step policy. The obtained policy is a pure strategy that starts with the joint action  $c = \{moveS(\mathbf{a}_1, stand), moveS(\mathbf{a}_2, N)\}$ . Indeed, since agent  $\mathbf{a}_2$  is closer to the base, and the *moveS* action is not reliable, from the perspectives of  $\mathbf{a}_1$  and  $\mathbf{a}_2$ , an optimal joint strategy is to keep  $\mathbf{a}_1$  idle, letting  $\mathbf{a}_2$  try to pick up some gold from (7, 4). That is, the produced policy is the following:

$$\begin{aligned} \pi = c; & \text{ if } obs_{\mathbf{a}_2}(success) \text{ then } [ \\ & \{moveS(\mathbf{a}_1, stand), pickUpS(\mathbf{a}_2, g)\}; \\ & \text{ if } obs_{\mathbf{a}_2}(success) \text{ then } [ \\ & \{moveS(\mathbf{a}_1, stand), moveS(\mathbf{a}_2, E)\}; \\ & \text{ if } obs_{\mathbf{a}_2}(success) \text{ then } [ \\ & \{moveS(\mathbf{a}_1, stand), moveS(\mathbf{a}_2, E)\}; \\ & \text{ if } obs_{\mathbf{a}_2}(success) \text{ then } \\ & \{moveS(\mathbf{a}_1, stand), dropS(\mathbf{a}_2, g)\} \\ & \text{ else if } obs_{\mathbf{a}_2}(failure) \text{ then } \pi_1 ] \\ & \text{ else if } obs_{\mathbf{a}_2}(failure) \text{ then } \pi_2 ] \\ & \text{ else if } obs_{\mathbf{a}_2}(failure) \text{ then } \pi_3 ] \\ & \text{ else if } obs_{\mathbf{a}_2}(failure) \text{ then } \pi_4. \end{aligned}$$

In this policy, every  $\pi_n$  is a sequence of  $n$  joint idle actions, that is,  $\pi_n = \{moveS(\mathbf{a}_1, stand), moveS(\mathbf{a}_2, stand)\}; \pi_{n-1}$ . The rewards of the ob-

tained strategy  $\pi$  are 0.146 and 0.186 for  $\mathbf{a}_1$  and  $\mathbf{a}_2$ , respectively. Note that there are two other Nash equilibria of *schema*(2): A pure strategy  $\pi_1$ , which starts with  $c = \{moveS(\mathbf{a}_1, E), moveS(\mathbf{a}_2, stand)\}$ , and a mixed strategy  $\pi_2$ , which starts with  $\pi_{a_1} \cdot \pi_{a_2}$ , where  $\pi_{a_1} = \{(moveS(\mathbf{a}_1, stand), 0.4), (moveS(\mathbf{a}_1, E), 0.6)\}$  and  $\pi_{a_2} = \{(moveS(\mathbf{a}_2, N), 0.52), (moveS(\mathbf{a}_2, stand), 0.48)\}$  (both of them associated with a lower reward for both agents). The selected Nash equilibrium depends on the Nash selection function embedded in the POGTGolog interpreter, which can be implemented by an external Nash equilibrium solver. In this example, as two players are involved, a two-player solver for finding Nash equilibria in normal form games [40] can be deployed. In the general case, the Nash selection function is implemented by an n-person game solver (e.g., [30, 31]).

## 6. Related Work

In the situation calculus literature, we can find other frameworks for reasoning about actions in multi-agent contexts. In [62], the authors provide a situation-calculus-based framework suitable for reasoning about multi-agent beliefs, abilities, and multi-agent communicative actions. In this context, the agents' mental states are explicitly represented deploying a possible worlds semantics and introducing accessibility relations between situations (following [61]). In this setting, a Golog-based agent programming language is proposed for the specification and verification of complex multi-agent systems. In contrast, our main concern here is to use a relational representation inspired by partially observable stochastic games and functional abstractions for policy synthesis. Given this aim, our treatment of beliefs is different from the one in [62, 61, 1, 2, 3, 4]. In this paper, we are not interested in explicit representation of beliefs as modalities for reasoning about the mental states of the agents. Instead, belief states are outside the situation calculus object language, they are introduced to denote probability distributions over the situations and to assess situation calculus formulas with respect to these distributions. These belief states are exploited to assess specific properties (such as the probability vector  $\phi(b)$ , executable actions  $Poss(c, b)$ , and rewards *reward*), which are needed to specify the POGTGolog interpreter.

Explicit representations of the belief states are also proposed by the multi-agent epistemic planning literature [47, 37, 13, 14]. Differently from our approach, in these works, the belief states are explicitly represented in the language in order to enable reasoning about them within the proposed frameworks. In contrast, as already explained above, we introduce belief states outside the language as probability distributions over sets of situations. Specifically, the domain theory

only mentions situations, while we exploit meta-level expressions to define properties of the belief states. This means that we can only assess some properties of the belief state, which are tailored for our representation of partially observable stochastic games. For instance, we cannot represent and reason about nested and combined beliefs, which is beyond the scope of the paper. On the other hand, we can handle representations of rewards and probabilities, which are needed to represent and solve POSGs. Notice also that, in our framework, analogously to POSGs, stochastic actions provide observations; hence, sensing is here directly associated to their execution, whose effect is the belief state update.

More closely related to our framework are extensions of DTGolog [11, 17, 16] to the multi-agent setting. In [11], Lakemeyer and his group present ICPGolog, a multi-agent Golog framework for team playing. ICPGolog integrates different features like concurrency, exogenous actions, continuous change, and the possibility to project into the future. The framework is used in the robotic soccer domain. Multi-agent coordination is here achieved without communication by assuming that the world models of the agents do not differ too much. Differently from POGTGolog, the setting is fully observable and no game-theoretic mechanism is used. POMDPs are considered in [57], but the focus is on reasoning about decision-theoretic projections, while game-theoretic strategies are not considered. A high-level programming language that can manage belief states is also provided by [4]. The language, ALLEGRO, implements a deterministic fragment of Golog over online executions. The approach relies on the domain theory by [3]; hence, it enables to represent both discrete distributions and continuous densities; query mechanisms for nested belief properties are available, but also in this case, game-theoretic domains are not considered.

Another closely related work is Poole’s independent choice logic (ICL) [54], which is a representation and reasoning formalism for single- and multi-agent systems that is based on acyclic logic programs under different “choices”. Poole’s ICL can be used as a formalism for logically encoding games in extensive and normal form. Differently from POGTGolog, it aims more at representing games and generalized strategies, while the problem of policy synthesis is not addressed. Furthermore, partially observable stochastic games are also not treated.

In [44], high-level agent programming in the FLUX framework is considered in a multi-agent setting, modeling communicative actions among the agents. The agents can reason about the other agents’ knowledge and communication skills. Also in this case, the specification is based on modalities, while decision- and game-theoretic problems are not treated.

Logic-based multi-agent programming is also investigated in the BDI frame-

work, where the main focus is the specification and formal verification [6, 5] of BDI systems using BDI logical languages [67, 33].

Our approach is also different from that in [9, 10], which provides a situation-calculus-based framework for specifying games, but probabilities and utilities are not represented, while the focus is verification and synthesis. In [68, 69], the multi-agent epistemic situation calculus is exploited for strategy representation and reasoning for incomplete information concurrent games; also in this case, the focus is on strategic reasoning, while Markov models are not considered.

Several authors have investigated graphical representations of games [36, 39, 65], where each player’s reward function depends on a subset of players described in a graph structure, which exploit the locality of the interactions to obtain compact models and efficient algorithms. In these works, an  $n$ -player normal form game is explicitly described by an undirected graph on  $n$  vertices, representing the  $n$  players, and a set of  $n$  matrices, each representing a local subgame (involving only some of the players). In our system, the interaction structure is explicitly encoded, both in the action theory and in POGTGolog procedures. Hence, local dependencies among the players are also available. The multi-agent influence diagrams in [65] permit a structured and compact representation of extensive form games, involving time and information, using graphical models. Like in [54], the framework extends influence diagrams and Bayesian networks to the multi-agent setting. However, the focus of these works is very different from ours. In fact, their main concern is computational, and the structured representation is used to reduce the computational cost of finding equilibria. Instead, we propose an agent programming language suitable for multi-agent settings, integrating declarative and procedural features.

Other related works deal with multi-agent decision-theoretic planning in extensions of MDPs [51, 56, 29, 12]. In particular, decentralized POMDPs (Dec-POMDPs) [51, 29] have been explored, which are multi-agent POMDPs where the dynamic system is controlled by multiple distributed agents with common payoff, each with possibly different information about the current state of the world. Another approach to multi-agent POMDPs are communicative multi-agent team decision problems [56], which allow to subsume and analyze many existing models of multi-agent cooperative systems. Interestingly, both logic-based and decision-theoretic approaches can be embedded and assessed in this framework. The free communication model is also defined and analyzed in [56]. Closely related are Dec-POMDPs with communication [29], which allow for studying the tradeoff between the cost and the value of the information acquired in the communication process and its influence on the joint utility of the agents. Dec-POMDPs with

free communication are investigated in [60], where a free communication model is used at planning time to simplify the policy generation, while the problem of communication cost and limited resources is handled at the execution time. Differently from our free communication model, the agents in [56, 60] have a unique reward function. Identical payoff stochastic games [51] represent cooperative games by restricting each agent of the game to a single payoff representing the team reward. An algorithm that approximates POSGs as a series of smaller Bayesian games is proposed in [12]. Other approximate algorithms for POSGs can be found in [38, 34]. Interactive POMDPs [28] are a control paradigm that complements and generalizes the traditional (Nash) equilibrium approach. Like in our work, the main focus is on policy synthesis for agent control, but [28] addresses directly the synthesis problem in the space of states, while we focus on the representational problem aiming at providing a tool for specifying abstract domains and suitable for balancing the tradeoff between procedural programming and planning.

From the representational perspective, further related works focus on relational and first-order extensions of MDPs [7, 70, 43, 27], POMDPs [63], multi-agent MDPs [32], and Markov games [19]. In all these works, partial observable games are not addressed.

## 7. Conclusion

We have presented the agent programming language POGTGolog, which combines explicit agent programming in Golog with game-theoretic multi-agent planning in partially observable stochastic games. It allows for modeling one team of cooperative agents under partial observability, where the agents may also have different initial belief states and not necessarily the same rewards. POGTGolog allows to encode partial control programs in a high-level logical language, which are then completed by an interpreter in an optimal way. We have defined a formal semantics of POGTGolog programs in terms of Nash equilibria, and specified a POGTGolog interpreter that computes one of these Nash equilibria. We have also shown that POGTGolog programs faithfully extend partially observable stochastic games. We have illustrated the usefulness of this approach along several examples.

An interesting topic of future research is to generalize POGTGolog in the direction of weakening the free communication assumption. In particular, we are interested in mechanisms for limited and implicit communication among the agents [15]. Alternatively, one may also allow for explicit communication between the agents (for example, along the lines of [56, 29]) assuming a cost as-

sociated with communication actions. Another direction of future research is to generalize POGTGolog to two competing teams of cooperative agents under partially observability.

*Acknowledgments.* This paper is a significantly extended and revised version of a paper that has appeared in: *Proceedings KI-2006*, pp. 113–127. *LNCS/LNAI* 4314, Springer, 2007 [22]. This work was partially supported by the Austrian Science Fund (FWF) under the project P18146-N04, by the projects REFILLS (H2020-ICT-731590) and ICOSAF (PON R&I 2014–2020), by a Heisenberg Professorship of the German Research Foundation (DFG), by the Alan Turing Institute under the UK EPSRC grant EP/N510129/1, the AXA Research Fund, and by the EPSRC grants EP/R013667/1, EP/L012138/1, and EP/M025268/1.

## Appendix A: Proofs for Section 4

**Proof of Theorem 4.1.** Let  $DT = (AT, ST, OT)$  be a domain theory, let  $p$  be a POGTGolog program relative to  $DT$ , let  $b$  be a belief state, and let  $H \geq 0$  be a horizon. Observe first that  $DT \models \text{OptG}(\hat{p}, b, H, \pi, v, pr)$  implies  $DT \models G(\hat{p}, b, H, \pi, v, pr)$ . Hence, if  $DT \models \text{OptG}(\hat{p}, b, H, \pi, v, pr)$ , then  $\pi$  is a  $H$ -step policy of  $p$  in  $b$  relative to  $DT$ , and  $utility(v_i, pr_i)$  is its expected  $H$ -step utility to agent  $i \in I$ . Thus, it only remains to prove the following statement:  $(\star)$  if  $DT \models \text{OptG}(\hat{p}, b, H, \pi, v, pr)$ , then  $\pi$  is an  $H$ -step Nash equilibrium of  $p$  in  $b$  relative to  $DT$ . We give a proof by induction on the structure of  $\text{OptG}$ .

*Basis:* The statement  $(\star)$  trivially holds for the null program and zero horizon cases. Indeed, in these cases,  $\text{OptG}$  generates only the policy  $\pi = \text{nil}$ .

*Induction:* For every program construct that involves no action choice of one of the two agents, the statement  $(\star)$  holds by the induction hypothesis. We now prove the statement  $(\star)$  for the remaining constructs:

(1) *Nondeterministic action choice of agent  $i$ :* Let  $\hat{p} = [\text{choice}(i: a_1 | \dots | a_n); p']$ , and let  $\pi$  be the  $H$ -step policy associated with  $p$  via  $\text{OptG}$ . By the induction hypothesis, for every  $k \in \{1, \dots, n\}$ , it holds that  $DT \models \text{OptG}([a_k; p'], b, H, a_k; \pi_k, v_k, pr_k)$  implies that the policy  $a_k; \pi_k$  is an  $H$ -step Nash equilibrium of the program  $[a_k; p']$  in  $b$ . By construction,  $\pi$  is the policy with the maximal expected  $H$ -step utility among the  $a_k; \pi_k$ 's. Hence, any different action selection  $a_k$  would not be better for  $i$ , that is,  $U_i(H, b, a_q; \pi_q) \leq U_i(H, b, \pi)$  for all  $q \in \{1, \dots, n\}$ . That is, any first action deviation from  $\pi$  would not be better for the agent  $i$ . Moreover, since each  $a_k; \pi_k$  is an  $H$ -step Nash equilibrium of  $[a_k; p']$  in  $b$ , also any following deviation from  $\pi$  would not be better for  $i$ . In summary, this shows that  $U_i(H, b, \pi')$

$\leq U_i(H, b, \pi)$  for every  $H$ -step policy  $\pi'$  of  $p$  in  $b$  relative to  $DT$  obtained from  $\pi$  by modifying only actions of agent  $i$ . Also for any agent  $j \neq i$ , any unilateral deviation  $\pi'$  from  $\pi$  cannot be better. In fact, since  $j$  is not involved in the first action choice,  $j$  can deviate from  $\pi$  only after  $i$ 's selection of  $a_k; \pi_k$ , but this would not be better for  $j$  by the induction hypothesis. Hence,  $U_j(H, b, \pi') \leq U_j(H, b, \pi)$  for every  $H$ -step policy  $\pi'$  of  $p$  in  $b$  relative to  $DT$  obtained from  $\pi$  by modifying only actions of agent  $j$ .

(2) *Nondeterministic joint action choice*: Let  $\hat{p} = [\|_{j \in J} \mathbf{choice}(j: a_{j,1} | \dots | a_{j,n_j}); p']$ , and let  $\pi$  be the  $H$ -step policy that is associated with  $p$  via  $OptG$ . By the induction hypothesis, we have that  $DT \models OptG([\bigcup_{j \in J} a_j; p'], b, H, \bigcup_{j \in J} a_j; \pi_a, v_a, pr_a)$  implies that each  $\bigcup_{j \in J} a_j; \pi_a$  is an  $H$ -step Nash equilibrium of  $[\bigcup_{j \in J} a_j; p']$  in  $b$ . We now prove that  $\pi$  is an  $H$ -step Nash equilibrium of  $p$  in  $b$ . Observe first that, by construction,  $\pi$  is of the form  $\Pi_{j \in J} \pi_j; \pi'$ , where  $(\pi_j)_{j \in J}$  is a Nash equilibrium (computed via the Nash selection function *selectNash*) of the matrix game consisting of all  $r_a = utility(v_a, pr_a)|_J$  such that  $a \in A$ . Thus, if agent  $j$  deviates from  $\pi_j$  with  $\pi'_j$ , it would not do better, that is,  $U_j(H, b, \pi') \leq U_j(H, b, \pi)$ , where  $\pi'$  is obtained from  $\pi$  by replacing  $\pi_j$  by  $\pi'_j$ . That is, any first action deviation from  $\pi$  would not be better for agent  $j$ . Moreover, by the induction hypothesis, also any following deviation from  $\pi'$  would not be better for  $j$ . In summary, this shows that  $U_j(H, b, \pi') \leq U_j(H, b, \pi)$  for every  $H$ -step policy  $\pi'$  of  $p$  in  $b$  relative to  $DT$  that is obtained from  $\pi$  by modifying only actions of agent  $j$ .

(3) *Nondeterministic choice of two programs*: The line of argumentation is similar to the one in the case of nondeterministic action choice of agent  $i$  above.  $\square$

**Proof of Theorem 4.2.** Suppose that  $G = (I, Z, (A_i)_{i \in I}, (O_i)_{i \in I}, P, (R_i)_{i \in I})$  is a partially observable stochastic game. Without loss of generality, let the  $A_i$ 's be pairwise disjoint. We now construct a domain theory  $DT = (AT, ST, OT)$ , a set of situation constants  $\{S_z \mid z \in Z\}$ , a function  $g$  mapping any belief state  $b$  of  $G$  into the belief state  $B_b$  of  $DT$ , and a set of POGTGolog programs  $\{p^h \mid h \in \{0, \dots, H\}\}$  relative to  $DT$  such that  $\sigma = (\sigma_i)_{i \in I}$  is an  $H$ -step Nash equilibrium of  $G$ , where every  $\sigma_i(b, h) = \pi_i$ ,  $i \in I$ , is given by  $DT \models OptG(\hat{p}^h, B_b, h+1, \Pi_{i \in I} \pi_i; \pi', v, pr)$ , for all  $b$  and  $h \in \{0, \dots, H\}$ , and the expected  $H$ -step reward of  $\sigma$  in  $b$  to agent  $i$  is given by  $utility(v_i, pr_i)$ , where  $DT \models OptG(\hat{p}^H, B_b, H+1, \pi, v, pr)$ .

The basic action theory  $AT$  comprises a situation constant  $S_z$  for every state  $z \in Z$  and a fluent  $state(z, s)$  that associates with every situation  $s$  a state  $z \in Z$  such that  $state(z, S_z)$  for all  $z \in Z$ . Every state  $z \in Z$  is represented by a constant, and different states are interpreted in a different way. Informally, the set

of all situations is given by the set of all situations that are reachable from the situations  $S_z$  with  $z \in Z$  (and thus we neglect the situation  $S_0$ ), and  $Z$  partitions the set of all situations into equivalence classes (one for each  $z \in Z$ ) via the fluent  $state(z, s)$ . Given the situation constants  $S_z$  and the fluent  $state(z, S_z)$ , we map any belief state  $b = (b_i)_{i \in I}$  of  $G$  to a belief state  $B_b = (g(b_i))_{i \in I}$  through the function  $g$  such that every  $(z, \mu) \in b_i$  is mapped to  $(S_z, \mu) \in g(b_i)$ . The basic action theory  $AT$  also comprises a deterministic action  $n_{a,z}$  for every joint action  $a \in A = \times_{i \in I} A_i$  and  $z \in Z$ , which performs a transition into the situation  $S_z$ , that is,  $state(z, do(n_{a,z}, s))$  for all states  $z \in Z$  and situations  $s$ . The actions  $n_{a,z}$  are executable in every situation  $s$ , that is,  $Poss(n_{a,z}, s) \equiv \top$  for all states  $z \in Z$  and situations  $s$ .

The stochastic theory  $ST$  comprises the stochastic action  $\{a_i \mid i \in I\}$  for every joint action  $a \in A$  along with all axioms  $stochastic(\{a_i \mid i \in I\}, s, \{n_{a,z'}\}, o, P(z', o \mid z, a))$  such that (i)  $z, z' \in Z$ , (ii)  $s$  is a situation that satisfies  $state(z, s)$  and that contains at most  $H + 1$  actions, and (iii)  $o \in O = \times_{i \in I} O_i$ . Informally, the stochastic theory  $ST$  represents the transition probabilities encoded in  $P$ .

The optimization theory  $OT$  comprises all axioms  $reward(i, \{n_{a,z'}\}, s) = R_i(a, z)$  such that (i)  $i \in I$ , (ii)  $a \in A$ , (iii)  $z, z' \in Z$ , and (iv)  $s$  is a situation that satisfies  $state(z, s)$  with at most  $H + 1$  actions. Let  $f = selectNash$  be a Nash selection function for normal form games  $M = (I, (A_i)_{i \in I}, (S_i)_{i \in I})$ , and let the expected reward to agent  $i \in I$  under the Nash equilibrium  $f(M)$  be denoted by  $v_f^i(M)$ .

Finally, every POGTGolog program  $p^h$  is a sequence of  $h+1$  nondeterministic joint action choices of the form  $\|_{i \in I} \mathbf{choice}(i: a_{i,1} \mid \dots \mid a_{i,n_i})$ , where  $a_{i,1}, \dots, a_{i,n_i}$  are all the singleton subsets of  $A_i$  (representing all the actions in  $A_i$ ) for all  $i \in I$ .

Observe first that  $pr = 1$  for every success probability  $pr$  computed in  $OptG$  for such programs  $p^h$  (since the preconditions of all actions are always satisfied). Since  $utility(v, pr) = v \cdot pr$ , it then follows that  $utility(v, pr) = v$  for every expected reward  $v$  and success probability  $pr$  computed in  $OptG$  for the programs  $p^h$ .

We now prove the statement of the theorem by induction on the horizon  $H \geq 0$ . For every belief state  $b = (b_i)_{i \in I}$  of  $G$  and every  $h \in \{0, \dots, H\}$ , let the normal form game  $G[b, h] = (I, (A_i)_{i \in I}, (Q_i[b, h])_{i \in I})$  be defined by  $(Q_i[b, h])_{i \in I}(a) = v_a$ , where  $DT \models OptG([\{a_i \mid i \in I\}; \hat{p}^{h-1}], B_b, h+1, \pi_a, v_a, pr_a)$ . By induction on the horizon  $H \geq 0$ , we now prove that for all  $i \in I$ :

- ( $\star$ ) (i)  $Q_i[b, 0](a) = R_i(b_i, a)$  for every belief state  $b$  of  $G$ , and
- (ii)  $Q_i[b, h](a) = R_i(b_i, a) + \sum_{o \in O} v_f^i(G[(b_i^{a,o})_{i \in I}, h-1]) \cdot P_b(b_i^{a,o} \mid b_i, a)$  for



every belief state  $b$  of  $G$  and  $h \in \{1, \dots, H\}$ .

This then implies that  $(v_f^i(G[b, h]))_{i \in I} = v$  and  $f(G[b, h]) = (\pi_i)_{i \in I}$  are given by  $DT \models \text{Opt}G(\hat{p}^h, B_b, h+1, \Pi_{i \in I} \pi_i; \pi', v, pr)$ , for every  $b$  and every  $h \in \{0, \dots, H\}$ . Furthermore, by finite-horizon value iteration [35], the mixed policy  $\sigma = (\sigma_i)_{i \in I}$  that is defined by  $(\sigma_i(b, h))_{i \in I} = f(G[b, h])$ , for every  $b$  and every  $h \in \{0, \dots, H\}$ , is a  $H$ -step Nash equilibrium of  $G$ , and it holds that  $G_i(H, b, \sigma) = v_f^i(G[b, H])$  for every  $i \in I$ . This then proves the theorem. Hence, it only remains to show by induction on the horizon  $H \geq 0$  that the statement  $(\star)$  holds, which is done as follows:

*Basis:* Let  $H = 0$ , and thus we only have to consider the case  $h = 0$ . Let  $DT \models \text{Opt}G([\tilde{a}; \text{nil}], B_b, 1, \pi_a, v_a, pr_a)$ , where  $\tilde{a} = \{a_i \mid i \in I\}$ . Using the definition of  $\text{Opt}G$  for stochastic first program actions  $\tilde{a}$ , we then obtain  $v_a = \sum_{o \in O} v_{a,o} \cdot \text{prob}(\tilde{a}, B_b, o)$ , with  $DT \models \text{Opt}G([\tilde{a}_o; \text{nil}], B_b, 1, \tilde{a}_o; \pi_{a,o}, v_{a,o}, pr_{a,o})$ . Using the definition of  $\text{Opt}G$  for the case of stochastic first program action with observation, we then obtain  $v_{a,o} = v'_{a,o} + \text{reward}(\tilde{a}_o, B_b) = 0 + (R_i(b_i, a))_{i \in I}$ . It thus follows that  $v_a = (R_i(b_i, a))_{i \in I} \cdot \sum_{o \in O} \text{prob}(\tilde{a}, B_b, o) = (R_i(b_i, a))_{i \in I}$ .

*Induction:* Let  $H > 0$ . By the induction hypothesis, it holds that (i)  $Q_i[b, 0](a) = R_i(b_i, a)$  for every belief state  $b$  reachable from  $b_0$ , and (ii)  $Q_i[b, h](a) = R_i(b_i, a) + \sum_{o \in O} v_f^i(G[(b_i^{a,o})_{i \in I}, h-1]) \cdot P_b(b_i^{a,o} \mid b_i, a)$  for every belief state  $b$  reachable from  $b_0$  and every  $h \in \{1, \dots, H-1\}$ . Furthermore, as argued above,  $(v_f^i(G[b, h]))_{i \in I} = v$  and  $f(G[b, h]) = (\pi_i)_{i \in I}$  are given by  $DT \models \text{Opt}G(\hat{p}^h, B_b, h+1, \Pi_{i \in I} \pi_i; \pi', v, pr)$  for every  $b$  that is reachable from  $b_0$ , and every  $h \in \{0, \dots, H-1\}$ . Suppose now that  $DT \models \text{Opt}G([\tilde{a}; \hat{p}^{h-1}], B_b, h+1, \pi_a, v_a, pr_a)$ , with  $\tilde{a} = \{a_i \mid i \in I\}$ . Using the definition of  $\text{Opt}G$  for stochastic first program actions  $\tilde{a}$ , we then obtain that  $v_a = \sum_{o \in O} v_{a,o} \cdot \text{prob}(\tilde{a}, B_b, o)$ , where  $\text{prob}(\tilde{a}, B_b, o) = (P_b(b_i^{a,o} \mid b_i, a))_{i \in I}$ , and every  $v_{a,o}$  is given by  $DT \models \text{Opt}G([\tilde{a}_o; \hat{p}^{h-1}], B_b, h+1, \tilde{a}_o; \pi_{a,o}, v_{a,o}, pr_{a,o})$ . Using the definition of  $\text{Opt}G$  for the case of stochastic first program action with observation, we then obtain  $v_{a,o} = v'_{a,o} + \text{reward}(\tilde{a}_o, B_b) = v'_{a,o} + (R_i(b_i, a))_{i \in I}$ . By the induction hypothesis, we have that  $v'_{a,o} = (v_f^i(G[(b_i^{a,o})_{i \in I}, h-1]))_{i \in I}$ .

Therefore, we can conclude that  $v_{a,i} = R_i(b_i, a) + \sum_{o \in O} v_f^i(G[(b_i^{a,o})_{i \in I}, h-1]) \cdot P_b(b_i^{a,o} \mid b_i, a)$ .  $\square$

- [1] F. Bacchus, J. Y. Halpern, and H. J. Levesque. Reasoning about noisy sensors and effectors in the situation calculus. *Artif. Intell.*, 111(1/2):171–208, 1999.
- [2] V. Belle and G. Lakemeyer. Reasoning about imperfect information games in the epistemic situation calculus. In *Proceedings AAAI-2010*, pp. 255–260. AAAI Press, 2010.

- [3] V. Belle and H. J. Levesque. Reasoning about continuous uncertainty in the situation calculus. In *Proceedings IJCAI-2013*, pp. 732–738. IJCAI/AAAI, 2013.
- [4] V. Belle and H. J. Levesque. ALLEGRO: Belief-based programming in stochastic dynamical domains. In *Proceedings IJCAI-2015*, pp. 2762–2769. AAAI Press, 2015.
- [5] R. H. Bordini, M. Fisher, W. Visser, and M. Wooldridge. Model checking rational agents. *IEEE Intell. Syst.*, 19(5):46–52, 2004.
- [6] R. H. Bordini, M. Fisher, W. Visser, and M. Wooldridge. Verifying multi-agent programs by model checking. *Auton. Agents Multi-Agent Syst.*, 12(2):239–256, 2006.
- [7] C. Boutilier, R. Reiter, and B. Price. Symbolic dynamic programming for first-order MDPs. In *Proceedings IJCAI-2001*, pp. 690–700. Morgan Kaufmann, 2001.
- [8] C. Boutilier, R. Reiter, M. Soutchanski, and S. Thrun. Decision-theoretic, high-level agent programming in the situation calculus. In *Proceedings AAAI-2000*, pp. 355–362. AAAI Press/MIT Press, 2000.
- [9] G. De Giacomo, Y. Lesperance, and A. R. Pearce. Situation-calculus-based programs for representing and reasoning about game structures. In *Proceedings KR-2010*, pp. 445–455, AAAI Press, 2010.
- [10] G. De Giacomo, Y. Lesperance, and A. R. Pearce. Situation calculus game structures and GDL. In *Proceedings ECAI-2016*, pp. 408–416. IOS Press, 2016.
- [11] F. Dylla, A. Ferrein, and G. Lakemeyer. Specifying multirobot coordination in ICP-Golog – From simulation towards real robots. In *Proceedings AOS-2003*, 2003.
- [12] R. Emery-Montemerlo, G. Gordon, J. Schneider, and S. Thrun. Approximate solutions for partially observable stochastic games with common payoffs. In *Proceedings AAMAS-2004*, pp. 136–143. IEEE Computer Society, 2004.
- [13] T. Engesser, T. Bolander, R. Mattmüller, and B. Nebel. Cooperative epistemic multi-agent planning for implicit coordination. In *Proceedings of the M4M Workshop at ICLA-2017*, pp. 75–90, 2017.
- [14] T. Engesser, R. Mattmüller, B. Nebel, and M. Thielscher. Game description language and dynamic epistemic logic compared. In *Proceedings IJCAI-2018*, pp. 1795–1802. IJCAI, 2018.
- [15] A. Farinelli, A. Finzi, and T. Lukasiewicz. Team programming in Golog under partial observability. In *Proceedings IJCAI-2007*, pp. 2097–2102, Morgan Kaufmann, 2007.

- [16] A. Ferrein, C. Fritz, and G. Lakemeyer. On-line decision-theoretic Golog for unpredictable domains. In *Proceedings KI-2004*, volume 3238 of *LNCS/LNAI*, pp. 322–336. Springer, 2004.
- [17] A. Ferrein, C. Fritz, and G. Lakemeyer. Using Golog for deliberation and team coordination in robotic soccer. *Künstliche Intelligenz*, 1:24–43, 2005.
- [18] A. Finzi and T. Lukasiewicz. Game-theoretic agent programming in Golog. In *Proceedings ECAI-2004*, pp. 23–27. IOS Press, 2004.
- [19] A. Finzi and T. Lukasiewicz. Relational Markov games. In *Proceedings JELIA-2004*, volume 3229 of *LNCS/LNAI*, pp. 320–333. Springer, 2004.
- [20] A. Finzi and T. Lukasiewicz. Game-theoretic agent programming in Golog under partial observability In *Proceedings GTDT-2005*, 2005.
- [21] A. Finzi and T. Lukasiewicz. Game-theoretic Golog under partial observability. In *Proceedings AAMAS-2005*, pp. 1301–1302. ACM Press, 2005.
- [22] A. Finzi and T. Lukasiewicz. Game-theoretic agent programming in Golog under partial observability. In *Proceedings KI-2006*, volume 4314 of *LNCS/LNAI*, pp. 113–127. Springer, 2007.
- [23] A. Finzi and T. Lukasiewicz. Adaptive multi-agent programming in GTGolog. In *Proceedings ECAI-2006*, pp. 753–754. IOS Press, 2006.
- [24] A. Finzi and T. Lukasiewicz. Adaptive multi-agent programming in GTGolog. In *Proceedings KI-2006*, volume 4314 of *LNCS/LNAI*, pp. 389–403. Springer, 2007.
- [25] A. Finzi, F. Pirri, and R. Reiter. Open-world planning in the situation calculus. In *Proceedings AAAI-2000*, pp. 754–760, AAAI Press, 2000.
- [26] A. Finzi and F. Pirri. Combining probabilities, failures and safety in robot control. In *Proceedings IJCAI-2001*, pp. 1331–1336. Morgan Kaufmann, 2001.
- [27] N. H. Gardiol and L. Pack Kaelbling. Envelope-based planning in relational MDPs. In *Proceedings NIPS-2003*. MIT Press, 2003.
- [28] P. J. Gmytrasiewicz and P. Doshi. Interactive POMDPs: Properties and preliminary results. In *Proc. AAMAS-2004*, pp. 1374–1375. IEEE Computer Society, 2004.
- [29] C. V. Goldman and S. Zilberstein. Optimizing information exchange in cooperative multi-agent systems. In *Proceedings AAMAS-2003*, pp. 137–144. ACM Press, 2003.
- [30] S. Govindan and R. Wilson. A global Newton method to compute Nash equilibria. *J. Econ. Theory*, 110(1):65–86, 2003.

- [31] S. Govindan and R. Wilson. Computing Nash equilibria by iterated polymatrix approximation. *J. Econ. Dyn. Control*, 28:1229–1241, 2004.
- [32] C. Guestrin, D. Koller, C. Gearhart, and N. Kanodia. Generalizing plans to new environments in relational MDPs. In *Proceedings IJCAI-2003*, pp. 1003–1010. Morgan Kaufmann, 2003.
- [33] A. Herzig and N. Troquard. Knowing how to play: Uniform choices in logics of agency. In *Proceedings AAMAS-2006*, pp. 209–216. ACM Press, 2006.
- [34] K. Horák, B. Bosanský, and M. Pechoucek. Heuristic search value iteration for one-sided partially observable stochastic games. In *Proceedings AAAI-2017*, pp. 558–564. AAAI Press, 2017.
- [35] M. Kearns, Y. Mansour, and S. Singh. Fast planning in stochastic games. In *Proceedings UAI-2000*, pp. 309–316. Morgan Kaufmann, 2000.
- [36] M. Kearns, M. L. Littman, and S. Singh. Graphical models for game theory. In *Proceedings UAI-2001*, pp. 253–260. Morgan Kaufmann, 2001.
- [37] F. Kominis and H. Geffner. Beliefs in multiagent planning: From one agent to many. In *Proceedings AAAI-2015*, pp. 147–155. AAAI Press, 2015.
- [38] A. Kumar and S. Zilberstein. Dynamic programming approximations for partially observable stochastic games. In *Proceedings FLAIRS-2009*, pp. 547–552, 2009.
- [39] P. La Mura. Game networks. In *Proceedings UAI-2000*, pp. 335–342. Morgan Kaufmann, 2000.
- [40] C. E. Lemke and J. T. Howson, Jr. Equilibrium points of bimatrix games. *J. Soc. Ind. Appl. Math.*, 12(2):413–423, 1964.
- [41] M. L. Littman. Markov games as a framework for multi-agent reinforcement learning. In *Proceedings ICML-1994*, pp. 157–163. Morgan Kaufmann, 1994.
- [42] B. Marthi, S. J. Russell, D. Latham, and C. Guestrin. Concurrent hierarchical reinforcement learning. In *Proceedings IJCAI-2005*, pp. 779–785. Professional Book Center, 2005.
- [43] M. Martin and H. Geffner. Learning generalized policies from planning examples using concept languages. *Appl. Intell.*, 20(1):9–19, 2004.
- [44] Y. Martin, I. Narasamdya, and M. Thielscher. Knowledge of other agents and communicative actions in the fluent calculus. In *Proceedings KR-2004*, pp. 623–633. AAAI Press, 2004.

- [45] J. McCarthy and P. J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. In *Machine Intelligence*, volume 4, pp. 463–502. 1969.
- [46] R. D. McKelvey, A. M. McLennan, and T. L. Turocy. Gambit: Software tools for game theory, 2016. Available at: <http://www.gambit-project.org>.
- [47] C. Muise, V. Belle, P. Felli, S. A. McIlraith, T. Miller, A. R. Pearce, and L. Sonenberg. Planning over multi-agent epistemic states: A classical planning approach. In *Proceedings AAAI-2015*, pp. 3327–3334. AAAI Press, 2015.
- [48] R. Nair, M. Tambe, M. Yokoo, D. V. Pynadath, and S. Marsella. Taming decentralized POMDPs: Towards efficient policy computation for multiagent settings. In *Proceedings IJCAI-2003*, pp. 705–711. Morgan Kaufmann, 2003.
- [49] G. Owen. *Game Theory: Second Edition*. Academic Press, 1982.
- [50] L. Pack Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artif. Intell.*, 101(1/2):99–134, 1998.
- [51] L. Peshkin, K.-E. Kim, N. Meuleau, and L. Pack Kaelbling. Learning to cooperate via policy search. In *Proceedings UAI-2000*, pp. 489–496. Morgan Kaufmann, 2000.
- [52] J. Pinto. Integrating discrete and continuous change in a logical framework. *Comput. Intell.*, 14(1):39–88, 1998.
- [53] F. Pirri and R. Reiter. Some contributions to the metatheory of the situation calculus, *J. ACM*, 46(3):325–362, 1999.
- [54] D. Poole. The independent choice logic for modelling multiple agents under uncertainty. *Artif. Intell.*, 94(1/2):7–56, 1997.
- [55] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, 1994.
- [56] D. V. Pynadath and M. Tambe. The communicative multiagent team decision problem: Analyzing teamwork theories and models. *J. Artif. Intell. Res.*, 16:389–423, 2002.
- [57] G. Rens, T. A. Meyer, and G. Lakemeyer. A logic for reasoning about decision-theoretic projections. In *Proceedings ICAART-2015 (Revised Selected Papers)*, pp. 79–99, 2015.
- [58] R. Reiter. The frame problem in the situation calculus: A simple solution (sometimes) and a completeness result for goal regression. In *Artificial Intelligence and Mathematical Theory of Computation: Papers in Honor of John McCarthy*, pp. 359–380. Academic Press, 1991.

- [59] R. Reiter. *Knowledge in Action: Logical Foundations for Specifying and Implementing Dynamical Systems*. MIT Press, 2001.
- [60] M. Roth, R. Simmons, and M. Veloso. Reasoning about joint beliefs for execution-time communication decisions. In *Proceedings AAMAS-2005*, pp. 786–793. ACM Press, 2005.
- [61] R. B. Scherl and H. J. Levesque. The frame problem and knowledge-producing actions. In *Proceedings AAAI-1993*, pp. 689–697. AAAI Press/MIT Press, 1993.
- [62] S. Shapiro, Y. Lesperance, and H. J. Levesque. The cognitive agents specification language and verification environment for multiagent systems. In *Proceedings AAMAS-2002*, pp. 19–26. ACM Press, 2002.
- [63] S. Sanner and K. Kersting. Symbolic dynamic programming for first-order POMDPs. In *Proceedings AAAI-2010*, pp. 1140–1146, AAAI Press, 2010.
- [64] J. van der Wal. *Stochastic Dynamic Programming*, volume 139 of *Mathematical Centre Tracts*. Morgan Kaufmann, 1981.
- [65] D. Vickrey and D. Koller. Multi-agent algorithms for solving graphical games. In *Proceedings AAAI-2002*, pp. 345–351. AAAI Press, 2002.
- [66] J. von Neumann and O. Morgenstern. *The Theory of Games and Economic Behavior*. Princeton University Press, 1947.
- [67] M. Wooldridge. *Reasoning about rational agents*. MIT Press, Cambridge, MA, 2001.
- [68] L. Xiong and Y. Liu. Strategy representation and reasoning for incomplete information concurrent games in the situation calculus. In *Proceedings IJCAI-2016*, pp. 1322–1329. AAAI Press, 2016.
- [69] L. Xiong and Y. Liu. Strategy representation and reasoning in the situation calculus. In *Proceedings ECAI-2016*, pp. 982–990. AAAI Press, 2016.
- [70] S. W. Yoon, A. Fern, and R. Givan. Inductive policy selection for first-order MDPs. In *Proceedings UAI-2002*, pp. 568–576. Morgan Kaufmann, 2002.