

# Average Reward Timed Games<sup>\*</sup>

B. Thomas Adler<sup>1</sup>, Luca de Alfaro<sup>1</sup>, and Marco Faella<sup>1,2</sup>

<sup>1</sup> School of Engineering, University of California, Santa Cruz, USA

<sup>2</sup> Dipartimento di Scienze Fisiche, Università di Napoli “Federico II”, Italy

**Abstract.** We consider real-time games where the goal consists, for each player, in maximizing the average reward he or she receives per time unit. We consider zero-sum rewards, so that a reward of  $+r$  to one player corresponds to a reward of  $-r$  to the other player. The games are played on discrete-time game structures which can be specified using a two-player version of timed automata whose locations are labeled by reward rates. Even though the rewards themselves are zero-sum, the games are not, due to the requirement that time must progress along a play of the game.

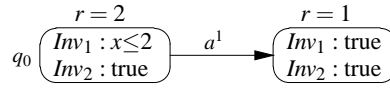
Since we focus on control applications, we define the value of the game to a player to be the maximal average reward per time unit that the player can ensure. We show that, in general, the values to players 1 and 2 do not sum to zero. We provide algorithms for computing the value of the game for either player; the algorithms are based on the relationship between the original, infinite-round game, and a derived game that is played for only finitely many rounds. As memoryless optimal strategies exist for both players in both games, we show that the problem of computing the value of the game is in  $\text{NP} \cap \text{coNP}$ .

## 1 Introduction

Games provide a setting for the study of control problems. It is natural to view a system and its controller as two players in a game; the problem of synthesizing a controller given a control goal can be phrased as the problem of finding a controller strategy that enforces the goal, regardless of how the system behaves [Chu63,RW89,PR89]. In the control of real-time systems, the games must not only model the interaction steps between the system and the controller, but also the amount of time that elapses between these steps. This leads to *timed games*, a model that was first applied to the synthesis of controllers for safety, reachability, and other  $\omega$ -regular goals [MPS95,AH97,AMAS98,HHM99,dAFH<sup>+</sup>03]. More recently, the problem of designing controllers for *efficiency* goals has been addressed, via the consideration of *priced* versions of timed games [BCFL04,ABM04]. In priced timed games, price rates (or, symmetrically, reward rates) are associated with the states of the game, and prices (or rewards) with its transitions. The problem that has so far been addressed is the synthesis of minimum-cost controllers for reachability goals [BCFL04,ABM04]. In this paper, we focus instead on the problem of synthesizing controllers that maximize the average

---

<sup>\*</sup> This is an extended version of a paper that appeared in the proceedings of *FORMATS 05: International Conference on Formal Modelling and Analysis of Timed Systems*, LNCS, Springer-Verlag, 2005. This research was supported in part by the NSF CAREER award CCR-0132780, by the ONR grant N00014-02-1-0671, and by the ARP award TO.030.MM.D.



**Fig. 1.** A game automaton where player 1 can freeze time to achieve a higher average reward.

reward<sup>1</sup> per time unit accrued along an infinite play of the game. This is an expressive and widely applicable efficiency goal, since many real-time systems are modeled as non-terminating systems which exhibit infinite behaviors.

We consider timed games played between two players over discrete-time game structures with finite state space. At each round, both players independently choose a move. We distinguish between *immediate moves*, which correspond to control actions or system transitions and take 0 time, and *timed moves*. There are two timed moves: the move  $\Delta_0$ , which signifies the intention to wait for 0 time, and the move  $\Delta_1$ , which signifies the intention of waiting for 1 time unit. The two moves chosen by the players jointly determine the successor state: roughly, immediate moves take the precedence over timed ones, and unit-length time steps occur only when both players play  $\Delta_1$ . Each state is associated with a reward rate, which specifies the reward obtained when staying at the state for one time unit. We consider zero-sum rewards, so that a reward of  $+r$  to one player corresponds to a reward of  $-r$  to the other player. These game structures can be specified using a notation similar to that of timed automata. Each location is labeled by a reward rate, and by two invariants (rather than one), which specify how long the two players can stay at the location; the actions labeling the edges correspond to the immediate moves of the players.

The goal of each player is to maximize the long-run average reward it receives per time unit; however, this goal is subordinate to the requirement that players should not block the progress of time by playing forever zero-delay moves (immediate moves, or  $\Delta_0$ ). As an example, consider the game of Figure 1. The strategy that maximizes the reward per time unit calls for player 1 staying forever at  $q_0$ : this yields an average reward per time unit of 4. However, such a strategy would block time, since the clock  $x$  would not be able to increase beyond the value 2, due to the player-1 invariant  $x \leq 2$  at  $q_0$ . If player 1 plays move  $a^1$ , time can progress, but the average reward per time unit is 1. To prevent players from blocking time in their pursuit of higher average reward, we define the value of a play of the game in a way that enforces time progress. If time diverges along the play, the value of the play is the average reward per time unit obtained along it. If time does not diverge along the play, there are two cases. If a player contributes to blocking the progress of time, then the value of the play to the player is  $-\infty$ ; if the progress of time is blocked entirely by the other player, then the value of the play to the player is  $+\infty$ . These definitions are based on the treatment of time divergence in timed games of [dAFH<sup>+</sup>03,dAHS02]. According to these definitions, even though the reward rate is zero-sum, and time-divergent plays have zero-sum values, the games are not zero-sum, due to the treatment of time divergence. Since we are interested in the problem of controller design, we define the value of a game to a player to be the maximal play value that the player is able to secure, regardless of how the adversary

<sup>1</sup> With a sign change, this is obviously equivalent to minimizing the average cost.

plays. The resulting games are not determined: that the values that the two players can secure do not sum to zero. We show that there is no symmetrical formulation that can at the same time enforce time progress, and lead to a determined setting.

We provide algorithms for computing the value of the game for either player. The algorithms are based on the relationship between the original, infinite-round, game, and a derived game that is played on the same discrete-time game structure, but for only finitely many rounds. As in [EM79], the derived game terminates whenever one of the two players closes a loop; our construction, however, differs from [EM79] in how it assigns a value to the loops, due to our different notion of value of a play. We show that a player can achieve the same value in the finite game, as in the original infinite-round game. Our proof is inspired by the argument in [EM79], and it closes some small gaps in the proof of [EM79].

The equivalence between finite and infinite games provides a PSPACE algorithm for computing the value of average reward discrete-time games. We improve this result by showing that both finite and infinite games admit memoryless optimal strategies for each player. Once we fix a memoryless strategy for a player, the game is reduced to a graph. We provide a polynomial-time algorithm that enables the computation of the value of the graph for the other player. The algorithm is based on polynomial-time graph transformations, followed by the application of Karp’s algorithm for computing the minimum/maximal average cost of a cycle [Kar78]. The existence of memoryless strategies, together with this algorithm, provide us with a polynomial witness and with a polynomial-time algorithm for checking the witness. Since this analysis can be done both for the winning strategies of a player, and for the “spoiling” strategies of the opponent, we conclude that the problem of computing the value of an average-reward timed game, for either player, is in  $NP \cap coNP$ . This matches the best known bounds for several other classes of games, among which are turn-based deterministic parity games [EJ91] and turn-based stochastic reachability games [Con92]. Since the maximum average reward accumulated in the first  $n$  time units cannot be computed by iterating  $n$  times a dynamic-programming operator, the weakly-polynomial algorithm of [ZP96] cannot be adapted to our games; the existence of polynomial algorithms is an open problem.

The goal of minimizing the long-run average cost incurred during the life of a real-time system has been considered previously in [BBL04]. There, the underlying model is a timed automaton, and the paper solves the verification problem (“what is the minimum long-run average cost achievable?”), or equivalently, the control problem for a fully deterministic system. In contrast, the underlying computational model in this paper is a timed game, and the problem solved is the control of a nondeterministic real-time system.

Compared to other work on priced timed games [BCFL04,ABM04], our models for timed games are simplified in two ways. First, rewards can only be accrued by staying at a state, and not by taking transitions. Second, we study the problem in discrete time. On the other hand, our models are more general in that, unlike [BCFL04,ABM04], we do not impose structural constraints on the game structures that ensure the progress of time. There is a tradeoff between imposing structural constraints and allowing rewards for transitions: had we introduced constraints that ensure time progress, we could have easily accommodated for rewards on the transitions. The restriction to discrete-time lim-

its somewhat the expressiveness of the models. Nevertheless, control problems where the control actions can be issued only at discrete points in time are very common: most real controllers are driven by a periodic clock; hence, the discrete-time restriction is not unduly limiting as far as the controller actions are concerned. We note that there are also many cases where the system actions can be considered to occur in discrete-time: this is the case, for instance, whenever the state of the system is sampled regularly in time.

## 2 Discrete-Time Game Structures

We define *discrete-time game structures* as a discrete-time version of the timed game structures of [dAFH<sup>+</sup>03]. A discrete-time game structure represents a game between two players, which we denote by 1, 2; we indicate by  $\sim i$  the opponent of  $i \in \{1, 2\}$  (that is, player  $3 - i$ ). A *discrete-time game structure* is a tuple  $\mathcal{G} = (S, Acts_1, Acts_2, \Gamma_1, \Gamma_2, \delta, r)$ , where:

- $S$  is a finite set of states.
- $Acts_1$  and  $Acts_2$  are two disjoint sets of actions for player 1 and player 2, respectively. We assume that  $\Delta_0, \Delta_1 \notin Acts_i$  and write  $M_i = Acts_i \cup \{\Delta_0, \Delta_1\}$  for the sets of moves of player  $i \in \{1, 2\}$ .
- For  $i \in \{1, 2\}$ , the function  $\Gamma_i : S \mapsto 2^{M_i} \setminus \emptyset$  is an enabling condition, which assigns to each state  $s$  a set  $\Gamma_i(s)$  of moves available to player  $i$  in that state.
- $\delta : S \times (M_1 \cup M_2) \mapsto S$  is a destination function that, given a state and a move of either player, determines the next state in the game.
- $r : S \mapsto \mathbb{Z}$  is a function that associates with each state  $s \in S$  the *reward rate* of  $s$ : this is the reward that player 1 earns for staying for one time unit at  $s$ .

The move  $\Delta_0$  represents an always-enabled stuttering move that takes 0 time: we require that for  $s \in S$  and  $i \in \{1, 2\}$ , we have  $\Delta_0 \in \Gamma_i(s)$  and  $\delta(s, \Delta_0) = s$ . The moves in  $\{\Delta_0\} \cup Acts_1 \cup Acts_2$  are known as the *zero-time* moves. The move  $\Delta_1$  represents the decision of waiting for 1 time unit. We do not require that  $\Delta_1$  be always enabled: if we have  $\Delta_1 \notin \Gamma_i(s)$  for player  $i \in \{1, 2\}$  at a state  $s \in S$ , then player  $i$  cannot wait, but must immediately play a zero-time move. We define the *size* of a discrete-time game structure by  $|\mathcal{G}| = \sum_{s \in S} (|\Gamma_1(s)| + |\Gamma_2(s)|)$ .

### 2.1 Move Outcomes, Runs, and Strategies

A timed game proceeds as follows. At each state  $s \in S$ , player 1 chooses a move  $a^1 \in \Gamma_1(s)$ , and simultaneously and independently, player 2 chooses a move  $a^2 \in \Gamma_2(s)$ . The set of successor states  $\tilde{\delta}(s, a^1, a^2) \subseteq S$  is then determined according to the following rules.

- *Actions take precedence over stutter steps and time steps.* If  $a^1 \in Acts_1$  or  $a^2 \in Acts_2$ , then the game takes an action  $a$  selected nondeterministically from  $A = \{a^1, a^2\} \cap (Acts_1 \cup Acts_2)$ , and  $\tilde{\delta}(s, a^1, a^2) = \{\delta(s, a) \mid a \in A\}$ .
- *Stutter steps take precedence over time steps.* If  $a^1, a^2 \in \{\Delta_0, \Delta_1\}$ , there are two cases.
  - If  $a^1 = \Delta_0$  or  $a^2 = \Delta_0$ , the game performs a stutter step, and  $\tilde{\delta}(s, a^1, a^2) = \{s\}$ .

- If  $a^1 = a^2 = \Delta_1$ , then the game performs a time step of duration 1, and the game proceeds to  $\tilde{\delta}(s, a^1, a^2) = \{\delta(s, \Delta_1)\}$ .

An *infinite run* (or simply *run*) of the discrete-time game structure  $\mathcal{G}$  is a sequence  $s_0, \langle a_1^1, a_1^2 \rangle, s_1, \langle a_2^1, a_2^2 \rangle, s_2, \dots$  such that  $s_k \in S$ ,  $a_{k+1}^1 \in \Gamma_1(s_k)$ ,  $a_{k+1}^2 \in \Gamma_2(s_k)$ , and  $s_{k+1} \in \tilde{\delta}(s_k, a_{k+1}^1, a_{k+1}^2)$  for all  $k \geq 0$ . A *finite run*  $\sigma$  is a finite prefix of a run that terminates at a state  $s$ , we then set  $last(\sigma) = s$ . We denote by *FRuns* the set of all finite runs of the game structure, and by *Runs* the set of its infinite runs. For a finite or infinite run  $\sigma$ , and a number  $k < |\sigma|$ , we denote by  $\sigma_{\leq k}$  the prefix of  $\sigma$  up to and including state  $\sigma_k$ . A state  $s'$  is *reachable* from another state  $s$  if there exists a finite run  $\mathfrak{s}, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n$  such that  $s_0 = s$  and  $s_n = s'$ .

A *strategy*  $\pi_i$  for player  $i \in \{1, 2\}$  is a mapping  $\pi_i : FRuns \mapsto M_i$  that associates with each finite run  $\mathfrak{s}, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n$  the move  $\pi_i(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n)$  to be played at  $s_n$ . We require that the strategy only selects enabled moves, that is,  $\pi_i(\sigma) \in \Gamma_i(last(\sigma))$  for all  $\sigma \in FRuns$ . For  $i \in \{1, 2\}$ , let  $\Pi_i$  denote the set of all player  $i$  strategies. A strategy  $\pi_i$  for player  $i \in \{1, 2\}$  is *memoryless* if for all  $\sigma, \sigma' \in FRuns$  we have that  $last(\sigma) = last(\sigma')$  implies  $\pi_i(\sigma) = \pi_i(\sigma')$ . For strategies  $\pi_1 \in \Pi_1$  and  $\pi_2 \in \Pi_2$ , we say that a run  $s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots$  is *consistent* with  $\pi_1$  and  $\pi_2$  if, for all  $n \geq 0$  and  $i = 1, 2$ , we have  $\pi_i(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n) = a_{n+1}^i$ . We denote by  $Outcomes(s, \pi_1, \pi_2)$  the set of all runs that start in  $s$  and are consistent with  $\pi_1, \pi_2$ . Note that in our timed games, two strategies and a start state yield a *set* of outcomes, because if the players both propose actions, a nondeterministic choice between the two moves is made. According to this definition, strategies can base their choices on the entire history of the game, consisting of both past states and moves.

## 2.2 Discrete-Time Game Automata

We specify discrete-time game structures via *discrete-time game automata*, which are a discrete-time version of the *timed automaton games* of [dAFH<sup>+</sup>03]; both models are two-player versions of timed automata [AD94]. A *clock condition* over a set  $C$  of clocks is a boolean combination of formulas of the form  $x \preceq c$  or  $x - y \preceq c$ , where  $c$  is an integer,  $x, y \in C$ , and  $\preceq$  is either  $<$  or  $\leq$ . We denote the set of all clock conditions over  $C$  by  $ClkConds(C)$ . A *clock valuation* is a function  $\kappa : C \mapsto \mathbb{R}_{\geq 0}$ , and we denote by  $K(C)$  the set of all clock valuations for  $C$ .

A *discrete-time game automaton* is a tuple  $\mathcal{A} = (Q, C, Acts_1, Acts_2, E, \theta, \rho, Inv_1, Inv_2, Rew)$ , where:

- $Q$  is a finite set of locations.
- $C$  is a finite set of clocks.
- $Acts_1$  and  $Acts_2$  are two disjoint, finite sets of actions for player 1 and player 2, respectively.
- $E \subseteq Q \times (Acts_1 \cup Acts_2) \times Q$  is an edge relation.
- $\theta : E \mapsto ClkConds(C)$  is a mapping that associates with each edge a clock condition that specifies when the edge can be traversed. We require that for all  $(q, a, q_1), (q, a, q_2) \in E$  with  $q_1 \neq q_2$ , the conjunction  $\theta(q, a, q_1) \wedge \theta(q, a, q_2)$  is unsatisfiable. In other words, the game move and clock values determine uniquely the successor location.

- $\rho : E \mapsto 2^C$  is a mapping that associates with each edge the set of clocks to be reset when the edge is traversed.
- $Inv_1, Inv_2 : Q \mapsto ClkConds(C)$  are two functions that associate with each location an invariant for player 1 and 2, respectively.
- $Rew : Q \mapsto \mathbb{Z}$  is a function that assigns a reward  $Rew(q) \in \mathbb{Z}$  with each  $q \in Q$ .

Given a clock valuation  $\kappa : C \mapsto \mathbb{R}_{\geq 0}$ , we denote by  $\kappa + 1$  the valuation defined by  $(\kappa + 1)(x) = \kappa(x) + 1$  for all clocks  $x \in C$ . The clock valuation  $\kappa : C \mapsto \mathbb{R}_{\geq 0}$  satisfies the clock constraint  $\alpha \in ClkConds(C)$ , written  $\kappa \models \alpha$ , if  $\alpha$  holds when the clocks have the values specified by  $\kappa$ . For a subset  $C' \subseteq C$  of clocks,  $\kappa[C' := 0]$  denotes the valuation defined by  $\kappa[C' := 0](x) = 0$  if  $x \in C'$ , and by  $\kappa[C' := 0](x) = \kappa(x)$  otherwise.

The discrete-time game automaton  $\mathcal{A}$  induces a discrete-time game structure  $\llbracket \mathcal{A} \rrbracket$ , whose states consist of a location of  $\mathcal{A}$  and a clock valuation over  $C$ . The idea is the following. The move  $\Delta_0$  is always enabled at all states  $\langle q, \kappa \rangle$ , and leads again to  $\langle q, \kappa \rangle$ . The move  $\Delta_1$  is enabled for player  $i \in \{1, 2\}$  at state  $\langle q, \kappa \rangle$  if  $\kappa + 1 \models Inv_i(q)$ ; the move leads to state  $\langle q, \kappa + 1 \rangle$ . For player  $i \in \{1, 2\}$  and  $a \in Acts_i$ , the move  $a$  is enabled at a state  $\langle q, \kappa \rangle$  if there is a transition  $(q, a, q')$  in  $E$  which is enabled at  $\langle q, \kappa \rangle$ , and if the invariant  $Inv_i(q')$  holds for the destination state  $\langle q', \kappa[\rho(q, a, q') := 0] \rangle$ . If the values of the clocks can grow unboundedly, this translation would yield an infinite-state discrete-time game structure. However, we can define *clock regions* similarly to timed automata [AD94], and we can include in the discrete-time game structure only one state per clock region; as usual, this leads to a finite state space.

### 3 The Average Reward Condition

In this section, we consider a discrete-time game structure  $\mathcal{G} = (S, Acts_1, Acts_2, \Gamma_1, \Gamma_2, \delta, r)$ , unless otherwise noted.

#### 3.1 The Value of a Game

We consider games where the goal for player 1 consists in maximizing the average reward per time unit obtained along a game outcome. The goal for player 2 is symmetrical, and it consists in minimizing the average reward per time unit obtained along a game outcome. To make these goals precise, consider a finite run  $\sigma = \sigma_0, \langle \sigma_1^1, \sigma_1^2 \rangle, \sigma_1, \dots, \sigma_n$ . For  $k \geq 1$ , the time  $D_k$  elapsed at step  $k$  of the run is defined by  $D_k(\sigma) = 1$  if  $\sigma_k^1 = \sigma_k^2 = \Delta_1$ , and  $D_k(\sigma) = 0$  otherwise; the reward  $R_k$  accrued at step  $k$  of the run is given by  $R_k(\sigma) = r(\sigma_{k-1}) \cdot D_k(\sigma)$ . The time elapsed during  $\sigma$  and the reward achieved during  $\sigma$  are defined in the obvious way, by  $D(\sigma) = \sum_{k=1}^n D_k(\sigma)$  and  $R(\sigma) = \sum_{k=1}^n R_k(\sigma)$ . Finally, we define the long-run average reward of an infinite run  $\sigma'$  by:

$$\bar{r}(\sigma') = \liminf_{n \rightarrow \infty} \frac{R(\sigma'_{\leq n})}{D(\sigma'_{\leq n})}.$$

A first attempt to define the goal of the game consists in asking for the maximum value of this long-run average reward that player 1 can secure. According to this approach, the value for player 1 of the game at a state  $s$  would be defined by

$$\tilde{v}(\mathcal{G}, s) = \sup_{\pi_1 \in \Pi_1} \inf_{\pi_2 \in \Pi_2} \inf \{ \bar{r}(\sigma) \mid \sigma \in Outcomes(s, \pi_1, \pi_2) \}.$$

However, this approach fails to take into account the fact that, in timed games, players must not only play in order to achieve the goal, but must also play realistic strategies that guarantee the advancement of time. As an example, consider the game of Figure 1. We have  $\tilde{v}(\langle q_0, [x := 0] \rangle) = 4$ , and the optimal strategy of player 1 consists in staying at  $q_0$  forever, never playing the move  $a^1$ . Due to the invariant  $x \leq 2$ , such a strategy blocks the progress of time: once  $x = 2$ , the only move player 1 can play is  $\Delta_0$ . It is easy to see that the only strategies of player 1 that do not block time eventually play move  $a^1$ , and have value 1. Note that the game does not contain any blocked states, i.e., from every reachable state there is a run that is time-divergent: the lack of time progress of the above-mentioned strategy is due to the fact that player 1 values more obtaining high average reward, than letting time progress.

To ensure that winning strategies do not block the progress of time, we modify the definition of value of a run, so that ensuring time divergence has higher priority than maximizing the average reward. Following [dAFH<sup>+</sup>03], we introduce the following predicates:

- For  $i \in \{1, 2\}$ , we denote by  $blameless^i(\sigma)$  (“*blameless i*”) the predicate defined by  $\exists n \geq 0. \forall k > n. \sigma_k^i = \Delta_1$ . Intuitively,  $blameless^i(\sigma)$  holds if, along  $\sigma$ , player  $i$  beyond a certain point cannot be blamed for blocking time.
- We denote by  $td(\sigma)$  (“time-divergence”) the predicate defined by  $\forall n \geq 0. \exists k > n. [(\sigma_k^1 = \Delta_1) \wedge (\sigma_k^2 = \Delta_1)]$ .

We define the value of a run  $\sigma \in Runs$  for player  $i \in \{1, 2\}$  by:

$$w_i(\sigma) = \begin{cases} +\infty & \text{if } blameless^i(\sigma) \wedge \neg td(\sigma); \\ (-1)^{(i+1)} \bar{r}(\sigma) & \text{if } td(\sigma); \\ -\infty & \text{if } \neg blameless^i(\sigma) \wedge \neg td(\sigma). \end{cases} \quad (1)$$

It is easy to check that, for each run, exactly one of the three cases of the above definition applies. Notice that if  $td(\sigma)$  holds, then  $w_1(\sigma) = -w_2(\sigma)$ , so that the value of time-divergent runs is defined in a zero-sum fashion. We define the value of the game for player  $i$  at  $s \in S$  as follows:

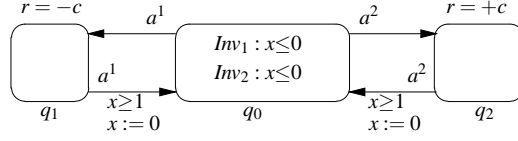
$$v_i(\mathcal{G}, s) = \sup_{\pi_i \in \Pi_i} \inf_{\pi_{\sim i} \in \Pi_{\sim i}} \inf\{w_i(\sigma) \mid \sigma \in Outcomes(s, \pi_1, \pi_2)\}. \quad (2)$$

We omit the argument  $\mathcal{G}$  from  $v_i(\mathcal{G}, s)$  when clear from the context.

We say that a state  $s \in S$  is *well-formed* if, for all  $i \in \{1, 2\}$ , we have  $v_i(s) > -\infty$ . From (1) and (2), a state is well-formed if both players can ensure that time progresses from that state, unless blocked by the other player: this is the same notion of well-formedness introduced in [dAHS02, dAFH<sup>+</sup>03]. Since we desire games where time progresses, we consider only games consisting of well-formed states.

### 3.2 Determinacy

A game is *determined* if, for all  $s \in S$ , we have  $v_1(s) + v_2(s) = 0$ : this means that if player  $i \in \{1, 2\}$  cannot enforce a reward  $c \in \mathbb{R}$ , then player  $\sim i$  can enforce at least reward  $-c$ . The following theorem provides a strong non-determinacy result for average-reward discrete-time games.



**Fig. 2.** A game automaton. Unspecified guards and invariants are “true”.

**Theorem 1.** (*non-determinacy*) For all  $c > 0$ , there exists a game structure  $\mathcal{G} = (S, Acts_1, Acts_2, \Gamma_1, \Gamma_2, \delta, r)$  with a state  $s \in S$ , and two “spoiling” strategies  $\pi_1^* \in \Pi_1$ ,  $\pi_2^* \in \Pi_2$ , such that the following holds:

$$\begin{aligned} \sup_{\pi_1 \in \Pi_1} \sup\{w_1(\sigma) \mid \sigma \in Outcomes(s, \pi_1, \pi_2^*)\} &\leq -c \\ \sup_{\pi_2 \in \Pi_2} \sup\{w_2(\sigma) \mid \sigma \in Outcomes(s, \pi_1^*, \pi_2)\} &\leq -c. \end{aligned}$$

As a consequence,  $v_1(s) \leq -c$  and  $v_2(s) \leq -c$ .

Note that in the theorem we take  $\sup$ , rather than  $\inf$  as in (2), over the set of outcomes arising from the strategies. Hence, the theorem states that even if the choice among actions is resolved in favor of the player trying to achieve the value, there is a game with a state  $s$  where  $v_1(s) + v_2(s) \leq -2c < 0$ . Moreover, in the theorem, the adversary strategies are fixed, again providing an advantage to the player trying to achieve the value.

*Proof.* Consider the game of Figure 2. We take for  $\pi_1^* \in \Pi_1$  and  $\pi_2^* \in \Pi_2$  the strategies that play always  $\Delta_0$  in  $q_0$ , and  $\Delta_1$  elsewhere. Let  $s_0 = \langle q_0, [x := 0] \rangle$ , and consider the value

$$\hat{v}_1(s_0) = \sup_{\pi_1 \in \Pi_1} \sup\{w_1(\sigma) \mid \sigma \in Outcomes(s_0, \pi_1, \pi_2^*)\}.$$

There are two cases. If eventually player 1 plays forever  $\Delta_0$  in  $s_0$ , player 1 obtains the value  $-\infty$ , as time does not progress, and player 1 is not blameless. If player 1, whenever at  $s_0$ , eventually plays  $a^1$ , then the value of the game to player 1 is  $-c$ . Hence, we have  $\hat{v}_1(s_0) = -c$ . The analysis for player 2 is symmetrical. ■

The example of Figure 2, together with the above analysis, indicates that we cannot define the value of an average reward discrete-time game in a way that is symmetrical, leads to determinacy, and enforces time progress. In fact, consider again the case in which player 2 plays always  $\Delta_0$  at  $s_0$ . If, beyond some point, player 1 plays forever  $\Delta_0$  in  $s_0$ , time does not progress, and the situation is symmetrical wrt. players 1 and 2: they both play forever  $\Delta_0$ . Hence, we must rule out this combination of strategies (either by assigning value  $-\infty$  to the outcome, as we do, or by some other device). Once this is ruled out, the other possibility is that player 1, whenever in  $s_0$ , eventually plays  $a^1$ . In this case, time diverges, and the average value to player 1 is  $-c$ . As the analysis is symmetrical, the value to both players is  $-c$ , contradicting determinacy.



## 4 Solution of Average Reward Timed Games

In this section, we solve the problem of computing the value of an average reward timed game with respect to both players. First, we define a turn-based version of the timed game. Such version is equivalent to the first game when one is concerned with the value achieved by a specific player. Then, following [EM79], we define a finite game and we prove that it has the same value as the turn-based infinite game. This will lead to a PSPACE algorithm for computing the value of the game. We then show that the finite and, consequently, the infinite game admit memoryless optimal strategies for both players; as mentioned in the introduction, this will enable us to show that the problem of computing the value of the game is in  $\text{NP} \cap \text{coNP}$ .

In the remainder of this section, we consider a fixed discrete-time game structure  $\mathcal{G} = (S, \text{Acts}_1, \text{Acts}_2, \Gamma_1, \Gamma_2, \delta, r)$ , and we assume that all states are well-formed. We focus on the problem of computing  $v_1(s)$ , as the problem of computing  $v_2(s)$  is symmetrical. For a finite run  $\sigma$  and a finite or infinite run  $\sigma'$  such that  $\text{last}(\sigma) = \text{first}(\sigma')$ , we denote by  $\sigma \cdot \sigma'$  their concatenation, where the common state is included only once.

### 4.1 Turn-based Timed Game

We describe a turn-based version of the timed game, where at each round player 1 chooses his move before player 2. Player 2 can thus use her knowledge of player 1's move to choose her own. Moreover, when both players choose an action, the action chosen by player 2 is carried out. This accounts for the fact that in the definition of  $v_1(s)$ , nondeterminism is resolved in favor of player 2 (see (2)). Notice that if player 2 prefers to carry out the action chosen by player 1, she can reply with the stuttering move  $\Delta_0$ . Definitions pertaining this game have a “ $\infty$ ” superscript that stands for “turn-based infinite”. We define the *turn-based joint destination function*  $\tilde{\delta}^\infty : S \times M_1 \times M_2 \mapsto S$  by

$$\tilde{\delta}^\infty(s, a^1, a^2) = \begin{cases} \delta(s, \Delta_1) & \text{if } a^1 = a^2 = \Delta_1 \\ \delta(s, \Delta_0) & \text{if } \{a^1, a^2\} \subseteq \{\Delta_0, \Delta_1\} \text{ and } a^1 = \Delta_0 \text{ or } a^2 = \Delta_0 \\ \delta(s, a^1) & \text{if } a^1 \in \text{Acts}_1 \text{ and } a^2 \in \{\Delta_0, \Delta_1\} \\ \delta(s, a^2) & \text{if } a^2 \in \text{Acts}_2 \end{cases}$$

As before, a run is an infinite sequence  $s_0, \langle a_1^1, a_1^2 \rangle, s_1, \langle a_2^1, a_2^2 \rangle, s_2, \dots$  such that  $s_k \in S$ ,  $a_{k+1}^1 \in \Gamma_1(s_k)$ ,  $a_{k+1}^2 \in \Gamma_2(s_k)$ , and  $s_{k+1} \in \tilde{\delta}^\infty(s_k, a_{k+1}^1, a_{k+1}^2)$  for all  $k \geq 0$ . A *1-run* is a finite prefix of a run ending in a state  $s$ , while a *2-run* is a finite prefix of a run ending in a move  $a \in M_1$ . For a 2-run  $\sigma = s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n, \langle a_{n+1}^1 \rangle$ , we set  $\text{last}(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n, \langle a_{n+1}^1 \rangle) = s_n$  and  $\text{last}_a(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n, \langle a_{n+1}^1 \rangle) = a_{n+1}^1$ . For  $i \in \{1, 2\}$ , we denote by  $\text{FRuns}_i$  the set of all  $i$ -runs. Intuitively,  $i$ -runs are runs where it is player  $i$ 's turn to move. In the turn-based game, a *strategy*  $\pi_i$  for player  $i \in \{1, 2\}$  is a mapping  $\pi_i : \text{FRuns}_i \mapsto M_i$  such that  $\pi_i(\sigma) \in \Gamma_i(\text{last}(\sigma))$  for all  $\sigma \in \text{FRuns}_i$ . For  $i \in \{1, 2\}$ , let  $\Pi_i^\infty$  denote the set of all player  $i$  strategies; notice that  $\Pi_1^\infty = \Pi_1$ . Player-1 memoryless strategies are defined as usual. We say that a player-2 strategy  $\pi \in \Pi_2^\infty$  is *memoryless* iff, for all  $\sigma, \sigma' \in \text{FRuns}_2$ ,  $\text{last}(\sigma) = \text{last}(\sigma')$  and  $\text{last}_a(\sigma) = \text{last}_a(\sigma')$  imply  $\pi(\sigma) = \pi(\sigma')$ .

For strategies  $\pi_1 \in \Pi_1^\infty$  and  $\pi_2 \in \Pi_2^\infty$ , we say that a run  $s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots$  is *consistent* with  $\pi_1$  and  $\pi_2$  if, for all  $n \geq 0$  and  $i = 1, 2$ , we have  $\pi_1(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n) = a_{n+1}^1$

and  $\pi_2(s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n, \langle a_{n+1}^1 \rangle) = a_{n+1}^2$ . Since  $\tilde{\delta}^t$  is deterministic, for all  $s \in S$ , there is a unique run that starts in  $s$  and is consistent with  $\pi_1$  and  $\pi_2$ . We denote this run by  $outcomes^{t\infty}(s, \pi_1, \pi_2)$ . The value assigned to a run, to a strategy and to the whole game are defined as follows. We set  $w_1^{t\infty}(\sigma) = w_1(\sigma)$ , and

$$v_1^{t\infty}(s, \pi_1) = \inf_{\pi_2 \in \Pi_2^t} w_1^{t\infty}(outcomes^{t\infty}(s, \pi_1, \pi_2)); \quad v_1^{t\infty}(s) = \sup_{\pi_1 \in \Pi_1^t} v_1^{t\infty}(s, \pi_1).$$

The following theorem follows from the definition of turn-based game and from (2).

**Theorem 2.** For all  $s \in S$ , it holds  $v_1(s) = v_1^{t\infty}(s)$ .

## 4.2 Turn-based Finite Game

We now define a finite turn-based game that can be played on a discrete-time game structure. Definitions pertaining this game have a “tf” superscript that stands for “turn-based finite”. The finite game ends as soon as a loop is closed. A *maximal run* in the finite game is a 1-run  $\sigma = s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n$  such that  $s_n$  is the first state that is repeated in  $\sigma$ . Formally,  $n$  is the least number such that  $s_n = s_j$ , for some  $j < n$ . We set  $loop(\sigma)$  to be the suffix of  $\sigma$ :  $s_j, \langle a_{j+1}^1, a_{j+1}^2 \rangle, \dots, s_n$ . For  $\pi_1 \in \Pi_1^t$ ,  $\pi_2 \in \Pi_2^t$ , and  $s \in S$ , we denote by  $outcomes^{tf}(s, \pi_1, \pi_2)$  the unique maximal run that starts in  $s$  and is consistent with  $\pi_1$  and  $\pi_2$ .

In the finite game, a maximal run  $\sigma$  ending with the loop  $\lambda$  is assigned the value of the infinite run obtained by repeating  $\lambda$  forever. Formally,  $w_1^{tf}(\sigma) = w_1(\sigma \cdot \lambda^\omega)$ , where  $\lambda^\omega$  denotes the concatenation of numerably many copies of  $\lambda$ . The value assigned to a strategy  $\pi_1 \in \Pi_1^t$  and the value assigned to the whole game are defined as follows.

$$v_1^{tf}(s, \pi_1) = \inf_{\pi_2 \in \Pi_2^t} w_1^{tf}(outcomes^{tf}(s, \pi_1, \pi_2)); \quad v_1^{tf}(s) = \sup_{\pi_1 \in \Pi_1^t} v_1^{tf}(s, \pi_1).$$

Notice that since this game is finite and turn-based, for all  $s \in S$ , it holds:

$$\sup_{\pi_1 \in \Pi_1^t} \inf_{\pi_2 \in \Pi_2^t} w_1^{tf}(outcomes^{tf}(s, \pi_1, \pi_2)) = \inf_{\pi_2 \in \Pi_2^t} \sup_{\pi_1 \in \Pi_1^t} w_1^{tf}(outcomes^{tf}(s, \pi_1, \pi_2)). \quad (3)$$

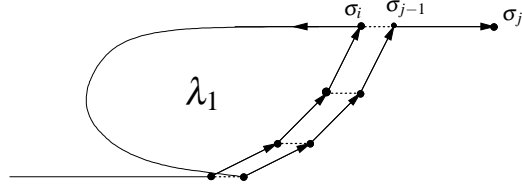
## 4.3 Mapping Strategies

We introduce definitions that allow us to relate the finite game to the infinite one. For a 1-run  $\sigma = s_0, \langle a_1^1, a_1^2 \rangle, s_1, \dots, s_n$ , let  $firstloop(\sigma)$  be the operator that returns the first simple loop (if any) occurring in  $\sigma$ . Similarly, let  $loopcut(\sigma)$  be the operator that removes the first simple loop (if any) from  $\sigma$ . Formally, if  $\sigma$  is a simple run (i.e. it contains no loops) we set  $firstloop(\sigma) = \varepsilon$  (the empty sequence), and  $loopcut(\sigma) = \sigma$ . Otherwise, let  $k \geq 0$  be the smallest number such that  $\sigma_j = \sigma_k$ , for some  $j < k$ ; we set

$$\begin{aligned} firstloop(\sigma) &= \sigma_j, \langle a_{j+1}^1, a_{j+1}^2 \rangle, \dots, \langle a_k^1, a_k^2 \rangle, \sigma_k; \\ loopcut(\sigma) &= \sigma_0, \langle a_1^1, a_1^2 \rangle, \dots, \sigma_j, \langle a_{k+1}^1, a_{k+1}^2 \rangle, \dots, \sigma_n. \end{aligned}$$

We now define the *quasi-segmentation*  $QSeg(\sigma)$  to be the sequence of simple loops obtained by applying  $firstloop$  repeatedly to  $\sigma$ .

$$QSeg(\sigma) = \begin{cases} \varepsilon & \text{if } firstloop(\sigma) = \varepsilon \\ firstloop(\sigma) \cdot QSeg(loopcut(\sigma)) & \text{otherwise} \end{cases}$$



**Fig. 3.** Nodes linked by dashed lines represent the same state of the game.

For an infinite run  $\sigma$ , we set  $QSeg(\sigma) = \lim_{n \rightarrow \infty} QSeg(\sigma_{\leq n})$ . Given a finite run  $\sigma$ , *loopcut* can only be applied a finite number of times before it converges to a fixpoint. We call this fixpoint  $resid(\sigma)$ . Notice that for all runs  $\sigma$ ,  $resid(\sigma)$  is a simple path and therefore its length is bounded by  $|S|$ .

For simplicity, we developed the above definitions for 1-runs. The corresponding definitions of  $resid(\sigma)$  and  $QSeg(\sigma)$  for 2-runs  $\sigma$  are similar.

For all  $i \in \{1, 2\}$  and all strategies  $\pi \in \Pi_i^t$ , we define the strategy  $\tilde{\pi}$  as  $\tilde{\pi}(\sigma) = \pi(resid(\sigma))$  for all  $\sigma \in FRuns_i$ . Intuitively,  $\tilde{\pi}$  behaves like  $\pi$  until a loop is formed. At that point,  $\tilde{\pi}$  *forgets* the loop, behaving as if the whole loop had not occurred. We now give some technical lemmas.

**Lemma 1.** *Let  $\pi_1 \in \Pi_1^t$ ,  $\pi_2 \in \Pi_2^t$ , and  $\sigma = outcomes^{t\infty}(s, \tilde{\pi}_1, \pi_2)$ . For all  $k > 0$ ,  $resid(\sigma_{\leq k})$  is a prefix of a finite run consistent with  $\pi_1$ . Formally, there is  $\pi'_2 \in \Pi_2^t$  and  $\sigma' = outcomes^{t\infty}(s, \pi_1, \pi'_2)$  such that  $\sigma' = resid(\sigma_{\leq k}) \cdot \rho$ .*

*Similarly, let  $\sigma = outcomes^{t\infty}(s, \pi_1, \tilde{\pi}_2)$ . For all  $k > 0$ , there is  $\pi'_1 \in \Pi_1^t$  and  $\sigma' = outcomes^{t\infty}(s, \pi'_1, \pi_2)$  such that  $\sigma' = resid(\sigma_{\leq k}) \cdot \rho$ .*

*Proof.* We prove the first statement, as the second one is analogous. We proceed by induction on the length of  $QSeg(\sigma_{\leq k})$ . If  $QSeg(\sigma_{\leq k})$  is the empty sequence (i.e.  $\sigma_{\leq k}$  contains no loops), the result is easily obtained, as  $\tilde{\pi}_1$  coincides with  $\pi_1$  until a loop is formed. So, we can take  $\pi'_2 = \pi_2$  and obtain the conclusion.

On the other hand, suppose  $QSeg(\sigma_{\leq k}) = \lambda_1, \dots, \lambda_n$ . For simplicity, suppose  $\lambda_1 \neq \lambda_2$ . As illustrated in Figure 3, let  $\sigma_j$  be the first state after  $\lambda_1$  that does not belong to  $\lambda_1$ . Then,  $\sigma_{j-1}$  belongs to  $\lambda_1$  and there is another index  $i < j - 1$  such that  $\sigma_i = \sigma_{j-1}$ . So, the game went twice through  $\sigma_{j-1}$  and two different successors were taken. However, player 1 must have chosen the same move in  $\sigma_i$  and  $\sigma_{j-1}$ , as by construction  $\tilde{\pi}_1(\sigma_{\leq i}) = \tilde{\pi}_1(\sigma_{\leq j-1})$ . Therefore, the change must be due to a different choice of  $\pi_2$ . It is easy to devise  $\pi'_2$  that coincides with  $\pi_2$ , except that  $\lambda_1$  may be skipped, and at  $\sigma_i$ , the successor  $\sigma_j$  is chosen. We can then obtain a run  $\rho = outcomes^{t\infty}(s, \tilde{\pi}_1, \pi'_2)$  and an integer  $k' \geq 0$  such that  $QSeg(\rho_{\leq k'}) = \lambda_2, \dots, \lambda_n$  and  $resid(\rho_{\leq k'}) = resid(\rho)$ . The thesis is obtained by applying the inductive hypothesis to  $\rho$  and  $k'$ . ■

Using this lemma, we can show that for all  $\pi_1 \in \Pi_1$ , each loop occurring in the infinite game under  $\tilde{\pi}_1$  can also occur in the finite game under  $\pi_1$ .

**Lemma 2.** *Let  $\pi_1 \in \Pi_1^t$ ,  $\pi_2 \in \Pi_2^t$ , and  $\sigma = outcomes^{t\infty}(s, \tilde{\pi}_1, \pi_2)$ . For all  $\lambda \in QSeg(\sigma)$ ,  $\lambda$  can occur as the final loop in a maximal run of the finite game. Formally, there is  $\pi'_2 \in \Pi_2^t$  and  $\sigma' = outcomes^{t\infty}(s, \pi_1, \pi'_2)$  such that  $\lambda = loop(\sigma')$ .*

Similarly, let  $\sigma = \text{outcomes}^{t\infty}(s, \pi_1, \tilde{\pi}_2)$ . For all  $\lambda \in QSeg(\sigma)$ , there is  $\pi'_1 \in \Pi_1^t$  and  $\sigma' = \text{outcomes}^{t\infty}(s, \pi'_1, \pi_2)$  such that  $\lambda = \text{loop}(\sigma')$ .

The next lemma states that if the strategy  $\pi_1$  of player 1 achieves value  $v$  in the fi nite turn-based game, the strategy  $\tilde{\pi}_1$  achieves at least as much in the infi nite turn-based game.

**Lemma 3.** For all  $s \in S$  and  $\pi_1 \in \Pi_1^t$ , it holds  $v_1^{t\infty}(s, \tilde{\pi}_1) \geq v_1^{t\infty}(s, \pi_1)$ .

*Proof.* Let  $v = v_1^{t\infty}(s, \pi_1)$ . We show that  $\tilde{\pi}_1$  can ensure reward  $v$  in the infi nite game. The result is trivially true if  $v = -\infty$ . So, in the following we assume that  $v > -\infty$ .

Fix a player 2 strategy  $\pi_2 \in \Pi_2^t$ , and let  $\sigma = \text{outcomes}^{t\infty}(s, \tilde{\pi}_1, \pi_2)$ . Let  $QSeg(\sigma) = \lambda_1, \lambda_2 \dots$ . We distinguish two cases, according to whether time diverges or not in  $\sigma$ . If time diverges, all loops  $\lambda_j$  that contain no tick give no contribution to the value of  $\sigma$  and can therefore be ignored.

For all  $\lambda_j$  containing (at least) a time step, by Lemma 2,  $\lambda_j$  is a possible terminating loop for the fi nite game under  $\pi_1$ . Thus,  $R(\lambda_j) \geq v \cdot D(\lambda_j)$ . Now, the value of  $\sigma$  can be split as the value due to loops containing time steps, plus the value due to the residual. For all  $n \geq 0$ , let  $m_n$  be the number of loops in  $QSeg(\sigma_{\leq n})$ . We obtain:

$$w_1^{t\infty}(\sigma) = \liminf_{n \rightarrow \infty} \frac{R(\sigma_{\leq n})}{D(\sigma_{\leq n})} = \liminf_{n \rightarrow \infty} \frac{R(\text{resid}(\sigma_{\leq n})) + \sum_{j=1}^{m_n} R(\lambda_j)}{D(\text{resid}(\sigma_{\leq n})) + \sum_{j=1}^{m_n} D(\lambda_j)} = \liminf_{n \rightarrow \infty} \frac{\sum_{j=1}^{m_n} R(\lambda_j)}{\sum_{j=1}^{m_n} D(\lambda_j)} \geq v.$$

Consider now the case when  $\sigma$  contains only fi nitely many time steps. Let  $k \geq 0$  be such that no time steps occur in  $\sigma$  after  $\sigma_k$ . Consider a loop  $\lambda_j$  entirely occurring after  $\sigma_k$ . Obviously  $\lambda_j$  contains no time steps. Moreover, by Lemma 2,  $\lambda_j$  is a terminating loop for a maximal run  $\rho$  in the fi nite game under  $\pi_1$ . Since  $v_1^{t\infty}(s, \pi_1) > -\infty$ , it must be  $w_1^{t\infty}(\rho) = +\infty$ . Consequently, it holds  $\text{blameless}^1(\rho)$  and in particular player 1 is blameless in all edges in  $\lambda_j$ .

Now, let  $k' \geq 0$  be such that each state (and edge) after  $\sigma_{k'}$  will eventually be part of a loop of  $QSeg(\sigma)$ . Let  $k'' = \max\{k, k'\}$ . Then, all edges that occur after  $k''$  will eventually be part of a loop where player 1 is blameless. Consequently,  $k''$  is a witness to the fact that  $\text{blameless}^1(\sigma)$ , and therefore  $w_1^{t\infty}(\sigma) = +\infty \geq v$ . ■

**Lemma 4.** For all  $s \in S$  and  $\pi_2 \in \Pi_2^t$ , it holds  $v_1^{t\infty}(s, \tilde{\pi}_2) \leq v_1^{t\infty}(s, \pi_2)$ .

*Proof.* Let  $v = v_1^{t\infty}(s, \pi_2)$ . Similarly to Lemma 3, we can rule out the case  $v = +\infty$  as trivial. Fix a player 1 strategy  $\pi_1$ , and let  $\sigma = \text{outcomes}^{t\infty}(s, \pi_1, \tilde{\pi}_2)$ . We show that  $w_1^{t\infty}(\sigma) \leq v$ . If time diverges on  $\sigma$ , the proof is similar to the analogous case in Lemma 3. Otherwise, let  $k \geq 0$  be such that no time steps occur in  $\sigma$  after  $\sigma_k$ . Consider a loop  $\lambda \in QSeg(\sigma)$ , entirely occurring after  $\sigma_k$ . Obviously  $\lambda$  contains no time steps. Moreover, by Lemma 2,  $\lambda$  is a terminating loop for a maximal run  $\rho$  in the fi nite game under  $\pi_1$ . Since  $v_1^{t\infty}(s, \pi_1) < +\infty$ , it must be  $w_1^{t\infty}(\rho) = -\infty$ . Consequently, it holds  $\neg \text{blameless}^1(\rho)$  and in particular player 1 is blamed in some edge of  $\lambda$ . This shows that  $\neg \text{blameless}^1(\sigma)$ , and consequently  $w_1^{t\infty}(\sigma) = -\infty \leq v$ . ■

Lemmas 3 and 4 show that the infi nite game is no harder than the fi nite one, for both players. Considering also (3), we obtain the following result.

**Theorem 3.** For all  $s \in S$ ,  $v_1^{\text{inf}}(s) = v_1^{\text{ff}}(s)$ .

Theorems 2 and 3 allow us to use the finite game to compute the value of the original timed game. The length of the finite game is bounded by  $|S|$ . It is well-known that a recursive, backtracking algorithm can compute the value of such game in PSPACE.

**Theorem 4.** For all  $s \in S$ ,  $v_1(s)$  can be computed in PSPACE.

#### 4.4 Memory

By following the “forgetful game” construction and proofs used by [EM79], we can derive a similar result on the existence of memoryless strategies for both players. The proof depends on the fact that the value of forgetful game is the same as the turn-based finite game (and hence, the same as the infinite game, from Theorem 3), and follows the same inductive steps as provided in [EM79].

**Theorem 5.** For all  $i \in \{1, 2\}$ , and  $t \in S$ , there exists a memoryless optimal strategy for player  $i$ . Formally, there exists  $\pi_i \in \Pi_i$  such that  $v_1(t, \pi_i) = v_1(t)$ .

#### 4.5 Improved Algorithms

We show that, given  $s \in S$ ,  $v \in \mathbb{Q}$  and  $i \in \{1, 2\}$ , the problem of checking whether  $v_i^{\text{ff}}(s) \geq v$  is in  $\text{NP} \cap \text{coNP}$ . The decision problem  $v_1^{\text{ff}}(s) \geq v$  is in NP because a memoryless strategy for player 1 acts as a polynomial-time witness: once such a strategy  $\pi_1$  is fixed, we can compute in polynomial time the value  $v_1^{\text{ff}}(s, \pi_1)$ . The problem is also in coNP because, once a memoryless strategy of player 2 is fixed, we can compute in polynomial time the value  $v_1^{\text{ff}}(s, \pi_2)$ .

Once we fix a memoryless strategy for player  $i \in \{1, 2\}$ , the finite game is reduced to a multigraph where all the choices belong to player  $\sim i$ . It is convenient to define the set of vertices of the multigraph as  $U = \{\{s\} \mid s \in S\}$ , rather than simply as  $S$ . Let  $E$  be the set of edges of the multigraph. Each edge  $e \in E$  is labeled with the pair of moves  $\langle a^1, a^2 \rangle \in M_1 \times M_2$  played by the players along  $e$ . We label  $e$  with *tick* whenever  $a^1 = a^2 = \Delta_1$ , and with  $bl_i$  whenever  $a^i \in \text{Acts}_i \cup \{\Delta_0\}$ ; every edge  $e$  from  $\{s\}$  to  $\{t\}$  is also associated with reward  $r(s)$  if it has label *tick*, and reward 0 otherwise. We indicate paths in this graph by  $u_0, e_1, u_1, e_2, \dots, u_n$ , where  $e_i$  is an edge from  $u_{i-1}$  to  $u_i$ , for  $1 \leq i \leq n$ . Given a strongly connected component (SCC)  $(V, F)$ , where  $V \subseteq U$  and  $F \subseteq E$ , we *collapse*  $(V, F)$  as follows: (i) we replace in  $U$  the vertices in  $V$  by the single vertex  $\bigcup V$ ; (ii) we remove all edges in  $F$ ; (iii) we replace every edge from  $v \in V$  to  $u \in U \setminus V$  (resp. from  $u \in U \setminus V$  to  $v \in V$ ) with an edge of the same label from  $\bigcup V$  to  $u$  (resp. from  $u$  to  $\bigcup V$ ); (iv) we replace every edge  $e \notin F$  from  $v \in V$  to  $v' \in V$  with a self-loop of the same label from  $\bigcup V$  to  $\bigcup V$ .

To determine the value of this multigraph to player 1, we first transform the multigraph so that all edges are labeled with *tick*, and we then apply Karp’s algorithm for computing the loop with minimum or maximum average reward [Kar78]. We proceed depending on whether player 1, or player 2, fixes a memoryless strategy. When player 1 fixes a memoryless strategy:

1. Find a maximal SCC  $(V, F)$ , where  $V \subseteq U$  and  $F \subseteq E$ , such that all edges in  $F$  are labeled with  $\neg tick$  and  $\neg bl_1$ . Player 2 will want to avoid following this SCC forever; thus, we collapse it. Repeat until no more SCCs can be collapsed.

2. If a vertex  $u \in U$  has no outgoing edges, it means that player 2 could not avoid entering and following one of the SCCs collapsed above. Hence, for each  $u \in U$  without outgoing edges, remove  $u$  from the graph along with all incoming edges, and assign value  $+\infty$  to all  $s \in u$ . Repeat until no more vertices can be removed.
3. Find all the loops whose edges are all labeled with  $\neg tick$ . Due to the collapsing in the above steps, each of these loops contains at least one edge labeled  $bl_1$ , so its value when followed forever is  $-\infty$ . Remove all such vertices from the graph, and assign value  $-\infty$  to the corresponding states.
4. From the resulting multigraph  $G$ , construct a multigraph  $G'$  with the same vertices as  $G$ . For each simple path in  $G$  of the form  $u_0, e_1, u_1, \dots, u_n, e_{n+1}, u_{n+1}$  where the edges  $e_1, \dots, e_n$  are labeled by  $\neg tick$ , and the edge  $e_{n+1}$  is labeled by  $tick$ , we insert in  $G'$  an edge from  $u_0$  to  $u_{n+1}$  labeled by the same reward as  $e_{n+1}$ .
5. Use the algorithm of [Kar78] to find the loop with minimal average reward in  $G'$  (the algorithm of [Kar78] is phrased for graphs, but it can be trivially adapted to multigraphs). If  $r$  is the average reward of the loop thus found, all the vertices of the loop, and all the vertices that can reach the loop, have value  $r$ . Remove them from  $G'$ , and assign value  $r$  to the corresponding states. Repeat this step until all vertices have been removed.

Similarly (but not symmetrically), if player 2 fixes a memoryless strategy, we can compute the value for player 1 as follows:

1. Find all the loops where all the edges are labeled with  $\neg tick$  and  $\neg bl_1$ . These loops, and all the vertices that can reach them, have value  $+\infty$ . Remove them from the graph, and assign value  $+\infty$  to the corresponding states.
2. Find a maximal SCC  $(V, F)$ , where  $V \subseteq U$  and  $F \subseteq E$ , such that all edges in  $F$  are labeled with  $\neg tick$ . Due to the previous step, every loop in  $(V, F)$  contains at least one edge labeled  $bl_1$ , and player 1 will want to avoid following forever such an SCC: thus, we collapse  $(V, F)$ .
3. For each  $u \in U$  without outgoing edges, remove  $u$  from the graph along with all incoming edges, and assign value  $-\infty$  to all  $s \in u$ . Repeat until no more vertices can be removed.
4. From the resulting multigraph  $G$ , construct a multigraph  $G'$  as in step 4 of the previous case.
5. This step is the same as step 5 of the previous case, except that in each iteration we find the loop with *maximal* average reward.

Since the algorithm of [Kar78], as well as the above graph manipulations, can all be done in polynomial time, we have the following result.

**Theorem 6.** *The problem of computing the value to player  $i \in \{1, 2\}$  of a discrete-time average reward game is in  $NP \cap coNP$ .*

We note that the maximal reward that a player can accrue in the first  $n$  time units cannot be computed by iterating  $n$  times a dynamic-programming operator, as is the case for untimed games. In fact, each player can play an unbounded number of zero-time moves in the first  $n$  time units, so that even the finite time-horizon version of our games requires the consideration of time divergence. Hence, it does not seem possible to adapt the approach of [ZP96] to obtain a weakly-polynomial algorithm. Whether polynomial-time algorithms can be achieved by other means is an open problem.

## References

- [ABM04] R. Alur, M. Bernadsky, and P. Madhusudan. Optimal reachability for weighted timed games. In *Proc. 31st Int. Colloq. Aut. Lang. Prog.*, volume 3142 of *Lect. Notes in Comp. Sci.* Springer-Verlag, 2004.
- [AD94] R. Alur and D.L. Dill. A theory of timed automata. *Theor. Comp. Sci.*, 126:183–235, 1994.
- [AH97] R. Alur and T.A. Henzinger. Modularity for timed and hybrid systems. In *CONCUR 97: Concurrency Theory. 8th Int. Conf.*, volume 1243 of *Lect. Notes in Comp. Sci.*, pages 74–88. Springer-Verlag, 1997.
- [AMAS98] E. Asarin, O. Maler, A.Pnueli, and J. Sifakis. Controller synthesis for timed automata. In *Proc. IFAC Symposium on System Structure and Control*, pages 469–474. Elsevier, 1998.
- [BBL04] P. Bouyer, E. Brinksma, and K.G. Larsen. Staying alive as cheaply as possible. In *Proc. of 7th Intl. Workshop on Hybrid Systems: Computation and Control (HSCC)*, volume 2993 of *Lect. Notes in Comp. Sci.*, pages 203–218. Springer-Verlag, 2004.
- [BCFL04] P. Bouyer, F. Cassez, E. Fleury, and K.G. Larsen. Optimal strategies in priced timed game automata. In *Found. of Software Technology and Theoretical Comp. Sci.*, volume 3328 of *Lect. Notes in Comp. Sci.*, pages 148–160. Springer-Verlag, 2004.
- [Chu63] A. Church. Logic, arithmetics, and automata. In *Proc. International Congress of Mathematicians, 1962*, pages 23–35. Institut Mittag-Leffler, 1963.
- [Con92] A. Condon. The complexity of stochastic games. *Information and Computation*, 96:203–224, 1992.
- [dAFH<sup>+</sup>03] L. de Alfaro, M. Faella, T.A. Henzinger, R. Majumdar, and M. Stoelinga. The element of surprise in timed games. In *CONCUR 03: Concurrency Theory. 14th Int. Conf.*, volume 2761 of *Lect. Notes in Comp. Sci.*, pages 144–158. Springer-Verlag, 2003.
- [dAHS02] L. de Alfaro, T.A. Henzinger, and M. Stoelinga. Timed interfaces. In *Proceedings of the Second International Workshop on Embedded Software (EMSOFT 2002)*, volume 2491 of *Lect. Notes in Comp. Sci.*, pages 108–122. Springer-Verlag, 2002.
- [EJ91] E.A. Emerson and C.S. Jutla. Tree automata, mu-calculus and determinacy (extended abstract). In *Proc. 32nd IEEE Symp. Found. of Comp. Sci.*, pages 368–377. IEEE Computer Society Press, 1991.
- [EM79] A. Ehrenfeucht and J. Mycielski. Positional strategies for mean payoff games. *Int. Journal of Game Theory*, 8(2):109–113, 1979.
- [HHM99] T.A. Henzinger, B. Horowitz, and R. Majumdar. Rectangular hybrid games. In *CONCUR'99: Concurrency Theory. 10th Int. Conf.*, volume 1664 of *Lect. Notes in Comp. Sci.*, pages 320–335. Springer-Verlag, 1999.
- [Kar78] R.M. Karp. A characterization of the minimum cycle mean in a digraph. *Discrete Mathematics*, 23:309–311, 1978.
- [MPS95] O. Maler, A. Pnueli, and J. Sifakis. On the synthesis of discrete controllers for timed systems. In *Proc. of 12th Annual Symp. on Theor. Asp. of Comp. Sci.*, volume 900 of *Lect. Notes in Comp. Sci.*, pages 229–242. Springer-Verlag, 1995.
- [PR89] A. Pnueli and R. Rosner. On the synthesis of a reactive module. In *Proceedings of the 16th Annual Symposium on Principles of Programming Languages*, pages 179–190. ACM Press, 1989.
- [RW89] P.J.G. Ramadge and W.M. Wonham. The control of discrete event systems. *IEEE Transactions on Control Theory*, 77:81–98, 1989.
- [ZP96] U. Zwick and M. Paterson. The complexity of mean payoff games on graphs. *Theor. Comp. Sci.*, 158:343–359, 1996.