

## Social laws in alternating time: effectiveness, feasibility, and synthesis

Wiebe van der Hoek · Mark Roberts ·  
Michael Wooldridge

Received: 18 May 2006 / Accepted: 12 June 2006  
Published online: 16 September 2006  
© Springer Science+Business Media B.V. 2006

**Abstract** Since it was first proposed by Moses, Shoham, and Tennenholtz, the *social laws* paradigm has proved to be one of the most compelling approaches to the offline coordination of multiagent systems. In this paper, we make four key contributions to the theory and practice of social laws in multiagent systems. First, we show that the *Alternating-time Temporal Logic* (ATL) of Alur, Henzinger, and Kupferman provides an elegant and powerful framework within which to express and understand social laws for multiagent systems. Second, we show that the *effectiveness*, *feasibility*, and *synthesis* problems for social laws may naturally be framed as ATL model checking problems, and that as a consequence, existing ATL model checkers may be applied to these problems. Third, we show that the complexity of the feasibility problem in our framework is no more complex in the general case than that of the corresponding problem in the Shoham–Tennenholtz framework (it is NP-complete). Finally, we show how our basic framework can easily be extended to permit social laws in which constraints on the legality or otherwise of some action may be explicitly required. We illustrate the concepts and techniques developed by means of a running example.

**Keywords** Multi-Agent Systems · Social Laws · Coordination · Alternating-time Temporal Logic · Model checking

---

This paper was presented at the *Social Software* conference in May 2004, Copenhagen, organised by PHILOG. We thank the organisers for providing the opportunity, and the participants for their useful feedback.

---

W. van der Hoek (✉) · M. Roberts · M. Wooldridge  
Department of Computer Science,  
University of Liverpool,  
Liverpool, L69 7ZF, UK  
e-mail: wiebe@csc.liv.ac.uk

## 1 Introduction

One of the defining problems in multiagent systems research is that of *coordination*—managing the interdependencies between the actions of multiple interacting agents (Durfee, 1999; Wooldridge, 2002). Techniques to coordinate activity in multiagent systems may be broadly categorised by whether they are “online” or “offline”. Online techniques aim to equip agents with the ability to dynamically coordinate their activities, for example by explicitly reasoning about coordination at run-time. In contrast, offline techniques aim at developing a coordination regime at design-time, and hard-wiring this regime into a system for use at run-time. There are arguments in favour of both approaches: the former is potentially more flexible, and may be more robust against unanticipated events, while the latter approach benefits from offline reasoning about coordination, thereby reducing the run-time decision-making burden on agents (Wooldridge, 2002).

One of the most successful approaches to “offline” coordination is the *social laws* paradigm, introduced largely through the work of Shoham, Tennenholtz, and Moses (Moses & Tennenholtz, 1995; Shoham & Tennenholtz, 1992; Shoham & Tennenholtz, 1996; Shoham & Tennenholtz, 1997). A social law can be understood as a set of rules imposed upon a multiagent system with the goal of ensuring that some desirable global (coordination) behaviour will result. Social laws work by *constraining* the behaviour of the agents in the system—by *forbidding* agents from performing certain actions in certain circumstances. In this sense, they can be understood as *norms* (Dignum, 1999). Shoham et al. investigated a number of issues surrounding the development and use of social laws, including the computational complexity of their synthesis, and the possibility of the development of social laws or conventions by the agents within the system themselves.

In this paper, we make four key contributions to the literature on social laws. First, we demonstrate that *Alternating-time Temporal Logic* (ATL)—the temporal logic of cooperation developed by Alur, Henzinger, and Kupferman (2002)—provides a rich and natural technical framework within which to investigate social laws and their properties. Second, we show that the *effectiveness*, *feasibility*, and *synthesis* problems for social laws in the ATL framework can be posed as ATL *model checking* problems (Clarke, Grumberg, & Peled, 2002) and that existing model checkers for ATL can hence be directly applied to these problems. Third, we show that, despite its apparent expressive power, the complexity of the feasibility problem in our framework is no more complex than the corresponding problem in the Shoham–Tennenholtz framework: it is NP-complete. We also identify cases where the problem becomes tractable. Finally, we show how our basic framework can easily be extended to permit social laws in which constraints on the legality or otherwise of some action may be explicitly required: for example, we show how the feasibility and synthesis problems for *dictatorship* social laws can be formulated.

We begin by introducing Action-based Alternating Transition Systems, (AATSS), the semantic structures that we employ throughout the remainder of the paper; we then introduce ATL itself, and give some properties of the logic that will be used. In Sect. 4, we introduce and formally define our social laws framework with respect to AATSS and ATL, present the feasibility and synthesis problems, and show how these can be understood as ATL model checking problems. A discussion of related work is presented in Sect. 5, and some conclusions are presented in 6.

## 2 Action-based alternating transition systems

The semantic structures underpinning ATL are known as *Alternating Transition Systems* (ATSS), and several essentially equivalent variations of these structures have been described in the ATL literature. As the notions of *action* and *action pre-condition* play such a key role in our framework, we find it convenient to work with another version of ATSS, in which actions and pre-conditions are firstclass citizens. We refer to these structures as *Action-based Alternating Transition Systems* (AATSS), and emphasise that they are equivalent to “conventional” ATSS. We first assume that the systems of interest to us may be in any of a finite set  $Q$  of possible *states*, with some  $q_0 \in Q$  designated as the *initial state*. Systems are populated by a set  $Ag$  of *agents*; a *coalition* of agents is simply a set  $G \subseteq Ag$ , and the set of all agents is known as the *grand coalition*. Each agent  $i \in Ag$  is associated with a set  $Ac_i$  of possible actions, and we assume that these sets of actions are pairwise disjoint (i.e., actions are unique to agents). We denote the set of actions associated with a coalition  $G \subseteq Ag$  by  $Ac_G$ , so  $Ac_G = \bigcup_{i \in G} Ac_i$ . A *joint action* for a coalition  $G$  is a tuple  $\langle \alpha_1, \dots, \alpha_k \rangle$ , where  $\alpha_i \in Ac_i$ , for each  $i \in G$ . We denote the set of all joint actions for coalition  $G$  by  $J_G$ , so  $J_G = \prod_{i \in G} Ac_i$ . Given an element  $j$  of  $J_G$  and agent  $i \in G$ , we denote  $i$ 's component of  $j$  by  $j_i$ .

An *Action-based Alternating Transition System*—hereafter referred to simply as an AATS—is an  $(n + 7)$ -tuple

$$S = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho, \tau, \Phi, \pi \rangle$$

where:

- $Q$  is a finite, non-empty set of *states*;
- $q_0 \in Q$  is the *initial state*;
- $Ag = \{1, \dots, n\}$  is a finite, non-empty set of *agents*;
- $Ac_i$  is a finite, non-empty set of *actions*, for each  $i \in Ag$ , where  $Ac_i \cap Ac_j = \emptyset$  for all  $i \neq j \in Ag$ ;
- $\rho: Ac_{Ag} \rightarrow 2^Q$  is an *action precondition function*, which for each action  $\alpha \in Ac_{Ag}$  defines the set of states  $\rho(\alpha)$  from which  $\alpha$  may be executed;
- $\tau: Q \times J_{Ag} \rightarrow Q$  is a partial *system transition function*, which defines the state  $\tau(q, j)$  that would result by the performance of  $j$  from state  $q$ —note that, as this function is partial, not all joint actions are possible in all states (cf. the pre-condition function above);
- $\Phi$  is a finite, non-empty set of *atomic propositions*; and
- $\pi: Q \rightarrow 2^\Phi$  is an interpretation function, which gives the set of primitive propositions satisfied in each state: if  $p \in \pi(q)$ , then this means that the propositional variable  $p$  is satisfied (equivalently, true) in state  $q$ .

We require that AATSS satisfy the following two coherence constraints:

1. *Non-triviality* [13]. Agents always have at least one legal action:

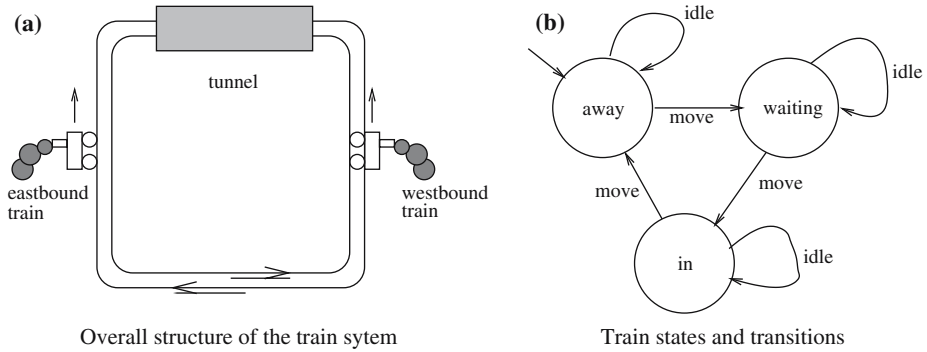
$$\forall q \in Q, \forall i \in Ag, \exists \alpha \in Ac_i \text{ s.t. } q \in \rho(\alpha)$$

2. *Consistency*. The  $\rho$  and  $\tau$  functions agree on actions that may be performed:

$$\forall q, \forall j \in J_{Ag}, (q, j) \in \text{dom } \tau \text{ iff } \forall i \in Ag, q \in \rho(j_i)$$

We denote the set of infinite sequences over  $Q$  by  $Q^\omega$ .

The next scenario explains our running example, to which we will refer many times in this paper:



**Fig. 1** The train system

**Example 1** There are two trains, one of which ( $E$ ) is Eastbound, the other of which ( $W$ ) is Westbound. Each train occupies its own circular track. At one point, both tracks pass through a narrow tunnel—a crash will occur if both trains are in the tunnel at the same time. Unlike the original version of this scenario (van der Hoek & Wooldridge, 2002), we do not assume that there is a “controller” agent, whose purpose is to ensure that collisions do not occur. Instead, we will be concerned with social laws that achieve this end.

We model each train  $i \in Ag = \{E, W\}$  as an automaton that can be in one of three states (see Fig. 1(b)): “ $away_i$ ” (the initial state of the train); “ $waiting_i$ ” (waiting to enter the tunnel); and “ $in_i$ ” (the train is in the tunnel). Each train  $i \in \{E, W\}$  has two actions available:  $Ac_i = \{move_i, idle_i\}$ . The  $idle_i$  action is the identity, which causes no change in the train’s state (i.e., it stays where it is). If a train  $i$  executes a  $move_i$  action while it is  $away_i$ , then it goes to a  $waiting_i$  state; executing a  $move_i$  while  $waiting_i$  causes a transition to an  $in_i$  state; and finally, executing a  $move_i$  while  $in_i$  causes a transition to  $away_i$  as long as the other train was not in the tunnel, while if both trains are in the tunnel, then they have crashed, and are forced to idle indefinitely. Initially, both trains are away.

The overall state of the system at any given time can be characterised in terms of the propositional variables  $\{away_E, away_W, waiting_E, waiting_W, in_E, in_W\}$ , where these variables have the obvious interpretation. The overall structure of the train system, and the model of trains is illustrated in Fig. 1; a formal definition of the train system AATS is given in Fig. 2 (the function  $\rho$  is left implicit, but can be read off from  $\tau$ : e.g.,  $\rho(move_W) = Q \setminus \{q8\}$ , etc.).

Of course, not all combinations of propositional variables correspond to reachable system states (i.e., states that the system could possibly enter). For example, an agent  $i$  cannot be both  $waiting_i$  and  $in_i$  simultaneously. There are in fact just nine reachable states of the system; see Fig. 2.

Given an agent  $i \in Ag$  and a state  $q \in Q$ , we denote the *options* available to  $i$  in  $q$ —the actions that  $i$  may perform in  $q$ —by  $options(i, q)$ :

$$options(i, q) = \{\alpha \mid \alpha \in Ac_i \text{ and } q \in \rho(\alpha)\}.$$

We then say that a *strategy* for an agent  $i \in Ag$  is a function:

$$\sigma_i: Q \rightarrow Ac_i$$

which must satisfy the *legality* constraint that  $\sigma_i(q) \in options(i, q)$  for all  $q \in Q$ .

States and Initial States:

$$Q = \{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8\} \quad \text{Initial state } q_0$$

Agents, Actions, and Joint Actions:

$$Ag = \{E, W\} \quad Ac_E = \{idle_E, move_E\} \quad Ac_W = \{idle_W, move_W\}$$

$$J_{Ag} = \{ \underbrace{\langle idle_E, idle_W \rangle}_{j_0}, \underbrace{\langle idle_E, move_W \rangle}_{j_1}, \underbrace{\langle move_E, idle_W \rangle}_{j_2}, \underbrace{\langle move_E, move_W \rangle}_{j_3} \}$$

Transitions/Pre-conditions:

$$\tau(q_0, j) = \begin{cases} q_0 & \text{if } j = j_0 \\ q_1 & \text{if } j = j_1 \\ q_3 & \text{if } j = j_2 \\ q_5 & \text{if } j = j_3 \end{cases} \quad \tau(q_4, j) = \begin{cases} q_4 & \text{if } j = j_0 \\ q_7 & \text{if } j = j_1 \\ q_0 & \text{if } j = j_2 \\ q_1 & \text{if } j = j_3 \end{cases}$$

$$\tau(q_1, j) = \begin{cases} q_1 & \text{if } j = j_0 \\ q_2 & \text{if } j = j_1 \\ q_5 & \text{if } j = j_2 \\ q_6 & \text{if } j = j_3 \end{cases} \quad \tau(q_5, j) = \begin{cases} q_5 & \text{if } j = j_0 \\ q_6 & \text{if } j = j_1 \\ q_7 & \text{if } j = j_2 \\ q_8 & \text{if } j = j_3 \end{cases}$$

$$\tau(q_2, j) = \begin{cases} q_2 & \text{if } j = j_0 \\ q_0 & \text{if } j = j_1 \\ q_6 & \text{if } j = j_2 \\ q_3 & \text{if } j = j_3 \end{cases} \quad \tau(q_6, j) = \begin{cases} q_6 & \text{if } j = j_0 \\ q_3 & \text{if } j = j_1 \\ q_8 & \text{if } j = j_2 \\ q_4 & \text{if } j = j_3 \end{cases}$$

$$\tau(q_3, j) = \begin{cases} q_3 & \text{if } j = j_0 \\ q_5 & \text{if } j = j_1 \\ q_4 & \text{if } j = j_2 \\ q_7 & \text{if } j = j_3 \end{cases} \quad \tau(q_7, j) = \begin{cases} q_7 & \text{if } j = j_0 \\ q_8 & \text{if } j = j_1 \\ q_1 & \text{if } j = j_2 \\ q_2 & \text{if } j = j_3 \end{cases}$$

$$\tau(q_8, j_0) = q_8$$

Propositional Variables:

$$\Phi = \{away_E, away_W, waiting_E, waiting_W, in_E, in_W\}$$

Interpretation Function:

$$\pi(q_0) = \{away_E, away_W\} \quad \pi(q_4) = \{in_E, away_W\}$$

$$\pi(q_1) = \{away_E, waiting_W\} \quad \pi(q_5) = \{waiting_E, waiting_W\}$$

$$\pi(q_2) = \{away_E, in_W\} \quad \pi(q_6) = \{waiting_E, in_W\}$$

$$\pi(q_3) = \{waiting_E, away_W\} \quad \pi(q_7) = \{in_E, waiting_W\}$$

$$\pi(q_8) = \{in_E, in_W\}$$

**Fig. 2** The AATS for the trains scenario

A strategy profile for a coalition  $G = \{a_1, \dots, a_k\} \subseteq Ag$  is a tuple of strategies  $(\sigma_1, \dots, \sigma_k)$ , one for each agent  $a_i \in G$ . We denote by  $\Sigma_G$  the set of all strategy profiles for coalition  $G \subseteq Ag$ ; if  $\sigma_G \in \Sigma_G$  and  $i \in G$ , then we denote  $i$ 's component of  $\sigma_G$  by  $\sigma_G^i$ . Given a strategy profile  $\sigma_G \in \Sigma_G$  and state  $q \in Q$ , let  $out(\sigma_G, q)$  denote the set of possible states that may result by the members of the coalition  $G$  acting as defined by their components of  $\sigma_G$  for one step from  $q$ :

$$out(\sigma_G, q) = \{q' \mid \tau(q, j) = q' \text{ where } (q, j) \in \text{dom } \tau \text{ and } \sigma_G^i(q) = j_i \text{ for } i \in G\}$$

Notice that, for any grand coalition strategy profile  $\sigma_{Ag}$  and state  $q$ , the set  $out(\sigma_{Ag}, q)$  will be a singleton.

A *computation* is an infinite sequence of states  $\lambda = q_0, q_1, \dots$ . A computation  $\lambda \in Q^\omega$  starting in state  $q$  is referred to as a *q-computation*; if  $u \in \mathbb{N}$ , then we denote by  $\lambda[u]$  the component indexed by  $u$  in  $\lambda$  (thus  $\lambda[0]$  denotes the first element,  $\lambda[1]$  the second, and so on).

Given a strategy profile  $\sigma_G$  for some coalition  $G$ , and a state  $q \in Q$ , we define  $comp(\sigma_G, q)$  to be the set of possible runs that may occur if every agent  $a_i \in G$  follows the corresponding strategy  $\sigma_i$ , starting when the system is in state  $q \in Q$ . That is, the set  $comp(\sigma_G, q)$  will contain all possible  $q$ -computations that the coalition  $G$  can “enforce” by cooperating and following the strategies in  $\sigma_G$ .

$$comp(\sigma_G, q) = \{\lambda \mid \lambda[0] = q \text{ and } \forall u \in \mathbb{N} : \lambda[u + 1] \in out(\sigma_G, \lambda[u])\}.$$

Again, note that for any state  $q \in Q$  and any grand coalition strategy  $\sigma_{Ag}$ , the set  $comp(\sigma_{Ag}, q)$  will be a singleton, consisting of exactly one infinite computation.

### 3 ATL

Alternating-time Temporal Logic (ATL), can be understood as a generalisation of the well-known branching time temporal logic CTL (Emerson, 1990), in which path quantifiers are replaced by *cooperation modalities*. A cooperation modality  $\langle\langle G \rangle\rangle\varphi$ , where  $G$  is a coalition, expresses the fact that the coalition  $G$  can cooperate to ensure that  $\varphi$ ; more precisely, that there exists a strategy profile for  $G$  such that by following this strategy profile,  $G$  can ensure  $\varphi$ . Thus, for example, the system requirement “agents 1 and 2 can cooperate to ensure that the system never enters a fail state” may be captured by the ATL formula  $\langle\langle 1, 2 \rangle\rangle \Box \neg fail$ . The “ $\Box$ ” temporal operator means “now and forever more”: additional temporal connectives in ATL are “ $\Diamond$ ” (“either now or at some point in the future”), “ $\mathcal{U}$ ” (“until”), and “ $\bigcirc$ ” (“in the next state”).

Formally, the set of ATL formulae, formed with respect to a set of agents  $Ag$ , and a set of primitive propositions  $\Phi$ , is given by the following grammar:

$$\varphi ::= p \mid \neg\varphi \mid \varphi \vee \varphi \mid \langle\langle G \rangle\rangle \varphi \mid \bigcirc \varphi \mid \langle\langle G \rangle\rangle \Box \varphi \mid \langle\langle G \rangle\rangle \varphi \mathcal{U} \varphi$$

where  $p \in \Phi$  is a propositional variable and  $G \subseteq Ag$  is a set of agents. We identify two fragments of ATL: the  $\bigcirc$ -fragment is the class of ATL-formulae containing only the  $\bigcirc$  temporal operator; and *propositional logic* formulae are those containing no cooperation modalities.

We can now give the rules defining the satisfaction relation “ $\models$ ” for ATL, which holds between pairs of the form  $S, q$  (where  $S$  is an AATS and  $q$  is a state in  $S$ ), and formulae of ATL:

$$S, q \models p \text{ iff } p \in \pi(q) \quad (\text{where } p \in \Phi);$$

$$S, q \models \neg\varphi \text{ iff } S, q \not\models \varphi;$$

$$S, q \models \varphi \vee \psi \text{ iff } S, q \models \varphi \text{ or } S, q \models \psi;$$

$$S, q \models \langle\langle G \rangle\rangle \bigcirc \varphi \text{ iff } \exists \sigma_G \in \Sigma_G, \text{ such that } \forall \lambda \in comp(\sigma_G, q), \text{ we have } S, \lambda[1] \models \varphi;$$

$S, q \models \langle\langle G \rangle\rangle \Box \varphi$  iff  $\exists \sigma_G \in \Sigma_G$ , such that  $\forall \lambda \in \text{comp}(\sigma_G, q)$ , we have  $S, \lambda[u] \models \varphi$   
for all  $u \in \mathbb{N}$ ;

$S, q \models \langle\langle G \rangle\rangle \varphi \mathcal{L} \psi$  iff  $\exists \sigma_G \in \Sigma_G$ , such that  $\forall \lambda \in \text{comp}(\sigma_G, q)$ , there exists some  $u \in \mathbb{N}$   
such that  $S, \lambda[u] \models \psi$ , and for all  $0 \leq v < u$ , we have  $S, \lambda[v] \models \varphi$ .

The remaining classical logic connectives (“ $\wedge$ ”, “ $\rightarrow$ ”, “ $\leftrightarrow$ ”) are assumed to be defined as abbreviations in terms of  $\neg, \vee$ , in the conventional manner. For readability, we omit set brackets in cooperation modalities, for example writing  $\langle\langle 1 \rangle\rangle$  instead of  $\langle\langle \{1\} \rangle\rangle$ .

Two cooperation modalities play a special role in the remainder of the paper, and are worth singling out for special attention. The cooperation modality  $\langle\langle \rangle\rangle$  (“the emptyset of agents can cooperate to...”) asserts that its argument is true on all computations, and thus acts like CTL’s universal path quantifier “A”. Similarly, the cooperation modality  $\langle\langle Ag \rangle\rangle$  asserts that its argument is satisfied on at least one computation, and thus acts like the CTL path quantifier “E”.

The *model checking problem* for ATL is the problem of determining, for any given ATL formula  $\varphi$ , AATS  $S$ , and state  $q$  in  $S$ , whether or not  $S, q \models \varphi$ . If  $S$  is an AATS and  $\varphi$  is a formula then we say that  $\varphi$  is *initially satisfied* in  $S$  if,  $S, q_0 \models \varphi$ ; we indicate this by writing  $S \models \varphi$ . A formula  $\varphi$  is *satisfiable* if there is some AATS  $S$  and state  $q$  in  $S$  such that  $S, q \models \varphi$ . The *satisfiability problem* for ATL is the problem of determining, for any given ATL formula, whether this formula is satisfiable or not.

### 3.1 Some properties of ATL

We begin by proving some properties of ATL that we use in subsequent proofs. For any  $S = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho, \tau, \Phi, \pi \rangle$ , let the relation  $R_{Ag}$  between states in  $Q$  be defined as:  $q_1 R_{Ag} q_2$  iff for some  $j \in J_{Ag}, \tau(q_1, j) = q_2$ . In other words,  $q_1 R_{Ag} q_2$  is true if the grand coalition can enforce a transition from  $q_1$  to  $q_2$ . Now, let  $S = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho, \tau, \Phi, \pi \rangle$  and  $S' = \langle Q', q'_0, Ag, Ac'_1, \dots, Ac'_n, \rho', \tau', \Phi', \pi' \rangle$  be two AATSs. We say that  $S'$  is a subsystem of  $S$ , (notation  $S' \sqsubseteq S$ ), or that  $S$  is a super-system of  $S'$ , (notation  $S \supseteq S'$ ), if there exists a relation  $\mathfrak{R} \subseteq Q \times Q'$  such that:

1.  $q_0 \mathfrak{R} q'_0$
2.  $\forall q \in Q, q' \in Q' : q \mathfrak{R} q' \Rightarrow (\pi(q) = \pi(q'))$
3.  $\forall q_1 \in Q, q'_1, q'_2 \in Q' : ((q_1 \mathfrak{R} q'_1 \& q'_1 R'_{Ag} q'_2) \Rightarrow (\exists q_2 \in Q : q_1 R_{Ag} q_2 \& q_2 \mathfrak{R} q'_2))$

A special case is obtained when  $S' = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho', \tau', \Phi, \pi \rangle$ , in which case we write  $S' \ll S$ . The relation  $\mathfrak{R}$  is essentially half of a *bisimulation* between two Kripke models (Blackburn, Rijke, & Venema, 2001, p. 64): the final clause represents the “backward clause” for such relations, with the “forward clause” not being present in the conditions on  $\mathfrak{R}$ . We are interested in which formulas are preserved when going from  $S$  to  $S'$ , where  $S' \sqsubseteq S$ . To this end, we define a *universal* and an *existential* sublanguage of ATL, denoted  $\mathcal{L}^u$  and  $\mathcal{L}^e$ , respectively. These languages are defined by the following grammars:

$$\begin{aligned} \mathcal{L}^u \quad v &::= p | \neg p | v \wedge v | v \vee v | \langle\langle \rangle\rangle \circ v | \langle\langle \rangle\rangle \diamond v | \langle\langle \rangle\rangle \Box v | \langle\langle \rangle\rangle v \mathcal{L}^u \\ \mathcal{L}^e \quad \epsilon &::= p | \neg p | \epsilon \wedge \epsilon | \epsilon \vee \epsilon | \langle\langle Ag \rangle\rangle \circ \epsilon | \langle\langle Ag \rangle\rangle \diamond \epsilon | \langle\langle Ag \rangle\rangle \Box \epsilon | \langle\langle Ag \rangle\rangle \epsilon \mathcal{L}^e \end{aligned}$$

where  $p \in \Phi$ .

**Lemma 1** Suppose we have  $S \sqsupseteq S'$  and  $v \in \mathcal{L}^u$ ,  $\epsilon \in \mathcal{L}^e$ . Then:

1.  $\forall q \in \mathcal{Q}, q' \in \mathcal{Q}'$  with  $q \mathfrak{R} q' : S, q \models v \Rightarrow S', q' \models v$ .
2.  $\forall q \in \mathcal{Q}, q' \in \mathcal{Q}'$  with  $q \mathfrak{R} q' : S', q' \models \epsilon \Rightarrow S, q \models \epsilon$ .

*Proof* We only prove the first item; the second follows from it and the observation that every  $\epsilon \in \mathcal{L}^e$  is equivalent to the negation of some  $v \in \mathcal{L}^u$ . The proof is by structural induction. For the case  $v$  is  $p$  or  $\neg p$  where  $p \in \Phi$ , the claim follows immediately from the second condition in the definition of subsystems. The cases of conjunction and disjunction follow directly, so assume that  $v$  equals  $\langle\langle\rangle\rangle \diamond v'$ , and that the claim is proven for  $v'$ . Assume that  $S, q \models \langle\langle\rangle\rangle \diamond v'$ , and  $q \mathfrak{R} q'$ . The former implies that for all strategy profiles  $\sigma \in \Sigma_{Ag}$ , for the unique  $\lambda_\sigma \in \text{comp}(\sigma, q)$ , we have that for some  $i_\sigma, S, \lambda[i_\sigma] \models v'$ . In order to derive a contradiction, suppose that  $S', q' \not\models \langle\langle\rangle\rangle \diamond v'$ . It would mean  $S', q' \models \langle\langle Ag \rangle\rangle \Box \neg v'$ , i.e., for some strategy profile  $\sigma' \in \Sigma_{Ag}$ , and the computation  $\lambda' \in \text{comp}(\sigma', q')$  and all  $i \in \mathbb{N}$ , we have  $S', \lambda'[i] \models \neg v'$ . Consider the run  $\lambda'[0]R_{Ag}\lambda'[1] \dots \lambda'[i] \dots$  that is induced by  $\lambda'$ , with  $\lambda'[0] = q$ . By the backward condition on  $\mathfrak{R}$ , we hence also find  $q_0 = q, q_1, \dots, q_i, \dots$  with  $q_i R_{Ag} q_{i+1}$  and  $q_i \mathfrak{R} \lambda'[i]$ . By the inductive assumption, we have  $S, q_i \models \neg v'$  for all  $i$ . But then, the agents  $Ag$  have a strategy to always ensure  $\neg v'$  from  $q$ : they just choose the actions that induce this path. This then means that  $S, q \models \langle\langle Ag \rangle\rangle \Box \neg v'$ , which contradicts the fact that  $S, q \models \langle\langle\rangle\rangle \diamond v'$ . The remaining cases are similar.  $\square$

It is not possible within the scope of this paper to examine other properties of ATL in detail; instead, we simply note some results that are pertinent to the goals of this paper.

**Theorem 1** (Alur et al., 2002; van Drimmelen, 2003; Pauly, 2001) The satisfiability problem for ATL is EXPTIME-complete (van Drimmelen, 2003), while the satisfiability problem for the  $\bigcirc$ -fragment of ATL is PSPACE-complete (Pauly, 2001). The model checking problem for “full” ATL can be solved in time  $O(|Q| \cdot |\varphi|)$ , where  $|Q|$  is the number of states, and  $|\varphi|$  is the size of the formula to be checked (Alur et al., 2002).

## 4 Social laws

We now introduce the formal framework of social laws, which we build upon throughout the remainder of the paper. Intuitively, a social law consists of two parts:

- An *objective*, which represents what the society aims to achieve by the introduction, or adoption of the law. In human societies, for example, the objective of a law might be to eliminate alcohol-related road accidents.
- A *behavioural constraint*, which corresponds to the requirements that a law places on the members of a society. A behavioural constraint corresponding to the objective of eliminating alcohol-related road accidents might be to forbid the action of driving a car after drinking alcohol.

We represent an objective for a social law as a formula of ATL, the intuition being that a social law is *effective* if it ensures that the objective is satisfied. (We give the formal definition shortly.) ATL provides a natural and powerful language for expressing the objectives of social laws, not least because such laws frequently refer to the *powers* or *rights* (or, conversely, the limits to powers) that agents have. The original social laws framework of Shoham and Tennenholtz illustrates this (Shoham & Tennenholtz,



1992). In this framework, each agent  $i \in Ag$  was associated with a set  $F_i \subseteq Q$  of focal states, the idea being that a successful social law would be one which ensured that if ever an agent  $i \in Ag$  found the environment was in a state  $q \in F_i$ , then  $i$  should have the power to ensure that, the environment eventually entered state  $q' \in F_i$ . This objective can naturally be expressed within our framework: assume that, for each  $q \in Q$ , we have a proposition  $\mathbf{q}$  that is satisfied iff the system is currently in state  $q$ . Then, given focal state sets  $F_1, \dots, F_n$ , we can express the Shoham–Tennenholtz objective as follows:

$$\bigwedge_{i \in Ag} \left( \bigwedge_{\mathbf{q} \in F_i} \left[ \mathbf{q} \rightarrow \bigwedge_{\mathbf{q}' \in F_i} \langle\langle i \rangle\rangle \diamond \mathbf{q}' \right] \right) \tag{1}$$

The ATL framework allows us to express much richer objectives, however. For example, from CTL it inherits the ability to express liveness and safety properties, and moreover we can reason about what certain coalitions can bring about.

**Example 2** Recall the trains scenario introduced earlier. The most obvious requirement for a social law is that the trains do not crash. The objective for this social law,  $O_1$ , is thus:

$$O_1 = \neg(in_E \wedge in_W)$$

The basic system  $S_1$  does not ensure that ( $O_1$ ) is satisfied—there are initial computations of the basic system on which both trains enter the tunnel simultaneously:

$$S_1, q_0 \not\models \langle\langle \rangle\rangle \square \neg(in_E \wedge in_W).$$

We model a behavioural constraint,  $\beta$ , as a function

$$\beta: Ac_{Ag} \rightarrow 2^Q$$

with the intended interpretation that if  $q \in \beta(\alpha)$ , then action  $\alpha$  may *not* be performed when the system is in state  $q$ —that is,  $\alpha$  is “forbidden” in state  $q$ . (As an aside, notice the duality between the pre-condition function  $\rho$ , and behavioural constraints  $\beta$ : if  $q \in \rho(\alpha)$ , then  $\alpha$  is permitted in  $q$ , whereas if  $q \in \beta(\alpha)$ , then  $\alpha$  is forbidden in  $q$ .) We will require that any behavioural constraint is “reasonable”, in that it always permits an agent to have at least one action left that can be performed in any state:

$$\forall i \in Ag, \forall q \in Q, \exists \alpha \in options(i, q) \text{ s.t. } q \notin \beta(\alpha).$$

We can now see what it means for a behavioural constraint  $\beta$  to be *implemented* in an AATS  $S$ : it simply means eliminating from  $S$  all transitions that are forbidden by  $\beta$ . The operation of implementing a behavioural constraint is thus an *update* on AATSS, in the sense that it results in a new AATS, which potentially satisfies different formulae. We denote the AATS obtained from  $S$  by implementing  $\beta$  by  $S \dagger \beta$ . Formally, if  $S = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho, \tau, \Phi, \pi \rangle$  is an AATS, and  $\beta$  is a behavioural constraint on  $S$ , then

$$S \dagger \beta = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho', \tau', \Phi, \pi \rangle,$$

where:

1.  $\forall \alpha \in Ac,$

$$\rho'(\alpha) = \rho(\alpha) \setminus \beta(\alpha)$$

$$2. \quad \forall q \in Q, \forall j \in J_{Ag},$$

$$\tau'(q, j) = \begin{cases} \tau(q, j) & \text{if } (q, j) \in \text{dom } \tau \text{ and } \forall i \in Ag, q \notin \beta(j_i) \\ \text{undefined} & \text{otherwise} \end{cases}$$

3. All other components of  $S \dagger \beta$  are as in  $S$ .

It is natural to ask what properties the implementation operator “ $\dagger$ ” has. First, notice that  $\text{AATSS}$  are closed under the implementation of behavioural constraints. That is, if  $S$  is an  $\text{AATS}$  and  $\beta$  is a behavioural constraint over  $S$ , then  $S \dagger \beta$  is an  $\text{AATS}$ —it satisfies the non-triviality and consistency coherence constraints given earlier. Second, notice that  $(S \dagger \beta) \dagger \beta = S \dagger \beta$ . Third, consider what properties of  $\text{AATSS}$  might be preserved by the implementation of social laws. To answer this question, first recall the notion of a *subsystem* as defined in Sect. 3.1; we have the following result.

**Lemma 2** *Let  $S$  be an  $\text{AATS}$ , and let  $\beta$  be a behavioural constraint over  $S$ . Then  $S \dagger \beta$  is a subsystem of  $S$ , i.e.,  $S \dagger \beta \sqsubseteq S$ . (In fact,  $S \dagger \beta \ll S$ .)*

From Lemma 1, we immediately obtain the following.

**Lemma 3** *Suppose we have an  $\text{AATS}$   $S$ , a behavioural constraint  $\beta$  over  $S$ , a state  $q$  in  $S$ , and formulae  $v \in \mathcal{L}^u, \epsilon \in \mathcal{L}^e$ . Then:*

1. *If  $S, q \models v$  then  $S \dagger \beta, q \models v$ .*
2. *If  $S \dagger \beta, q \models \epsilon$  then  $S, q \models \epsilon$ .*

The first of these two results tells us that implementing a behavioural constraint *guarantees to preserve the universal properties of a system*. However, it is easy to see that implementing a behavioural constraint does not guarantee to preserve existential properties. The second result, in contrast, tells us that if an existential property holds of a system in which some behavioural constraint has been implemented, then it must have held of the original system, before the constraint was imposed. Thus, implementing a behavioural constraint cannot *create* existential properties in a system.

Now, a *social law* over an  $\text{AATS}$   $S$  is a pair:

$$(\varphi, \beta)$$

where:

- $\varphi$  is an ATL formula called the *objective* of the law; and
- $\beta: Ac_{Ag} \rightarrow 2^Q$  is a behavioural constraint on  $S$ .

A social law  $(\varphi, \beta)$  is *effective* in  $\text{AATS}$   $S$  if, after implementing  $\beta$  in  $S$ , we know that  $\langle\langle \rangle\rangle \Box \varphi$  will be initially satisfied in  $S$ , i.e., if  $S \dagger \beta \models \langle\langle \rangle\rangle \Box \varphi$ .

There are three key computational questions with respect to social laws, which we shall investigate throughout the remainder of this paper:

1. *Effectiveness.* Given an  $\text{AATS}$   $S$  and a social law  $(\varphi, \beta)$  over  $S$ , determine whether  $(\varphi, \beta)$  is effective in  $S$ .
2. *Feasibility.* Given an  $\text{AATS}$   $S$  and a formula  $\varphi$  of ATL representing an objective, does there exist a behavioural constraint  $\beta$  such that  $(\varphi, \beta)$  is an effective social law in  $S$ .
3. *Synthesis.* Given an  $\text{AATS}$   $S$  and a formula  $\varphi$  of ATL representing an objective, exhibit a behavioural constraint  $\beta$  such that  $(\varphi, \beta)$  is an effective social law in  $S$  if such a constraint exists, otherwise answer “no”.

Our first result, with respect to the effectiveness problem, is now immediate.

**Lemma 4** *The effectiveness problem for social laws may be solved in time polynomial in the size of  $S$  and  $\varphi$ .*

*Proof* Generate  $S' = S \dagger \beta$ , and check that  $S', q_0 \models \langle\langle \rangle\rangle \Box \varphi$ . The first step can obviously be done in polynomial time: it simply requires eliminating every forbidden transition from  $\tau$ , and modify  $\rho$  similarly. From Theorem 1, so can the second step.  $\square$

With respect to the trains example, is there a behavioural constraint  $\beta_1$ , such that  $(O_1, \beta_1)$  is an effective social law? Clearly there is. The constraint  $\beta_1$  must ensure that the system never enters state  $q_8$ : from examination of the state transition function  $\tau$  (see Fig. 2), we can see that  $\tau(q_5, j_3) = \tau(q_6, j_2) = \tau(q_7, j_1) = q_8$ , and there are no other transitions leading to  $q_8$  (apart from when the trains have already crashed, which we need not consider!) Consider the behavioural constraint  $\beta_1$  as follows.

$$\beta_1(\alpha) = \begin{cases} \emptyset & \text{if } \alpha = \textit{idle}_E \\ \emptyset & \text{if } \alpha = \textit{idle}_W \\ \{q_5, q_6\} & \text{if } \alpha = \textit{move}_E \\ \{q_7\} & \text{if } \alpha = \textit{move}_W \end{cases}$$

The constraint ensures that:

- when both agents are waiting to enter the tunnel, the eastbound train is prevented from moving;
- when the westbound train is already in the tunnel and the eastbound train is waiting to enter the tunnel, then the eastbound train is prevented from moving; and
- when the eastbound train is already in the tunnel and the westbound train is waiting to enter the tunnel, then the westbound train is prevented from moving.

Notice that this constraint is, in a sense, asymmetric, as it constrains the eastbound train rather than the westbound train: we could equally well replace the first constraint with the requirement that if both trains are waiting to enter the tunnel, then the *westbound* train is prevented from moving, thus enabling the eastbound train to enter. Now, let  $S_2 = S_1 \dagger \beta_1$ . We claim that  $S_2, q_0 \models \langle\langle \rangle\rangle \Box \neg (in_E \wedge in_W)$ ; this can be established by inspection. Thus:

**Proposition 1**  $(O_1, \beta_1)$  is an effective social law in  $S_1$ .

Of course, there are other, less “sensible” behavioural constraints that are effective in  $S_1$  for  $O_1$ . Consider  $\beta_2$ :

$$\beta_2(\alpha) = \begin{cases} Q & \text{if } \alpha = \textit{move}_E \text{ or } \alpha = \textit{move}_W \\ \emptyset & \text{otherwise} \end{cases}$$

This behavioural constraint prevents both trains from moving, and yet:

**Proposition 2**  $(O_1, \beta_2)$  is an effective social law in  $S_1$ .

Clearly, our original objective needs some refinement. Consider objective  $O_2$ :

$$O_2 = O_1 \wedge \bigwedge_{i \in \{E, W\}} (\textit{waiting}_i \rightarrow \langle\langle i \rangle\rangle \Diamond (in_i \wedge O_1))$$

This objective ensures that, not only do the trains never crash, but that both trains can eventually safely enter the tunnel if they are waiting. Consider  $\beta_3$ , which works by forbidding trains from lingering in the tunnel, but is otherwise the same as  $\beta_1$ :

$$\beta_3(\alpha) = \begin{cases} \{q_4, q_7\} & \text{if } \alpha = \textit{idle}_E \\ \{q_2, q_6\} & \text{if } \alpha = \textit{idle}_W \\ \{q_5, q_6\} & \text{if } \alpha = \textit{move}_E \\ \{q_7\} & \text{if } \alpha = \textit{move}_W \end{cases}$$

**Proposition 3** ( $O_2, \beta_3$ ) is an effective social law in  $S_1$ .

However, objective  $O_2$  may still need further refinement, as less sensible behavioural constraints can still be effective in  $S_1$ :

**Proposition 4** ( $O_2, \beta_2$ ) is an effective social law in  $S_1$ .

So, consider objective  $O_3$ :

$$O_3 = O_1 \wedge \bigwedge_{i \in \{E, W\}} \left( \begin{array}{l} (\textit{away}_i \rightarrow \langle\langle i \rangle\rangle \diamond \textit{waiting}_i) \quad \wedge \\ (\textit{waiting}_i \rightarrow \langle\langle i \rangle\rangle \diamond (\textit{in}_i \wedge O_1)) \quad \wedge \\ (\textit{in}_i \rightarrow \langle\langle i \rangle\rangle \bigcirc \textit{away}_i) \end{array} \right)$$

This objective ensures that, not only do the trains never crash, but both trains will at some point enter the  $\textit{waiting}_i$  state. The objective also ensures that, once in the  $\textit{waiting}_i$  state, the train does eventually enter the state where it has safely entered the tunnel. Once the train is in the tunnel, it will not linger there, but its next state will be where it has safely left the tunnel and gone into the  $\textit{away}_i$  state.

**Proposition 5** ( $O_3, \beta_2$ ) is not an effective social law in  $S_1$ , but ( $O_3, \beta_3$ ) is.

The first issue to which we address ourselves is the computational complexity of the feasibility problem. As it turns out, the feasibility problem in our framework is no harder, for objectives expressed as arbitrary ATL formulae, than the corresponding problem studied by Shoham and Tennenholtz, where objectives were expressed as sets of focal states (see above); and for objectives propositional logic objectives, it is easier.

**Theorem 2** The feasibility problem for objectives expressed as arbitrary ATL formulae is NP-complete, and remains NP-complete for the CTL fragment of ATL, and also the universal ( $\mathcal{L}^u$ ) fragment. However, for objectives expressed as propositional formulae, the feasibility problem is decidable in polynomial time.

*Proof* With respect to the NP-completeness results, membership in NP may be seen by the following non-deterministic algorithm:

1. guess a behavioural constraint  $\beta$ ;
2. verify that  $\beta$  is effective.

Since  $\text{dom } \beta = Ac_{Ag}$ , step (1) can be done in (non-deterministic) polynomial time  $O(|Ac_{Ag} \times Q|)$ , while from Lemma 4, step (2) requires only polynomial time.

That the feasibility problem for ATL objectives is NP-hard follows immediately from the fact that the Useful Social Law problem of Shoham–Tennenholtz, (proven to be NP-complete in [Shoham & Tennenholtz, 1996, p. 612]), can be directly reduced to our feasibility problem, by encoding focal state sets as shown in Eq. 1: the reduction

is clearly possible in polynomial time. However, to see that the effectiveness problem for the  $\mathcal{L}^u$  fragment is NP-complete, we need to work a little harder. We reduce the directed Hamiltonian cycle problem (DHC) [(Papadimitriou, 1994), p. 209] to the effectiveness problem for  $\mathcal{L}^u$  objectives.

An instance of DHC is given by a directed graph  $H = \langle V, E \subseteq V \times V, v_0 \in V \rangle$ , where  $v_0$  is a distinguished “start” node. The aim is to determine whether or not  $H$  contains a directed Hamiltonian cycle starting from  $v_0$ , i.e., a cycle containing every vertex  $v \in V$ , in which no vertices are repeated.<sup>1</sup> The idea of the reduction is to encode the graph  $H$  directly in the state transformer function  $\tau$ : actions correspond to edges of the graph. Formally, given a directed graph  $H = \langle V, E \rangle$ , we define a single-agent AATS  $S^H = \langle Q^H, q_0^H, Ag^H, Ac_1^H, \rho^H, \tau^H, \Phi^H, \pi^H \rangle$  as follows:

- for each vertex  $v \in V$ , we create a state  $q_v \in Q^H$ , and in addition we create a “sink” state  $q_{sink}$ ;
- we set  $q_0 = v_0$ ;
- we create a single agent,  $Ag = \{a_1\}$ ;
- for each edge  $(v, v') \in E^H$  where  $v' \neq v_0$ , we define an action  $\alpha_{(v,v')}$ , define  $\tau^H(q_v, \alpha_{(v,v')}) = q_{v'}$ , and define  $\rho^H$  accordingly;
- for each edge  $(v, v') \in E^H$  where  $v' = v_0$ , we define an action  $\alpha_{(v,sink)}$ , define  $\tau^H(q_v, \alpha_{(v,sink)}) = q_{sink}$ , and define  $\rho^H$  accordingly;
- we create an action  $\alpha_{loop}$ , which may be performed only in state  $q_{sink}$ , such that  $\tau^H(q_{sink}, \alpha_{loop}) = q_{sink}$ , and define  $\rho^H$  accordingly;
- for each vertex  $v \in V$ , we create a proposition  $p_v \in \Phi^H$ ; and finally,
- we define  $\pi^H(q_v) = \{p_v\}$ .

We then define the objective  $O^H$  by

$$O^H \hat{=} \varphi_1^H \wedge \varphi_2^H$$

where:

$$\begin{aligned} \varphi_1^H &\hat{=} \bigwedge_{v \in V} \langle \langle \rangle \rangle \Diamond p_v \\ \varphi_2^H &\hat{=} \bigwedge_{v \in V} p_v \rightarrow \langle \langle \rangle \rangle \bigcirc \langle \langle \rangle \rangle \Box \neg p_v \end{aligned}$$

We now claim that there exists a DHC in graph  $H$  starting from  $v_0$  iff the objective  $O^H$  is feasible in AATS  $S^H$ . To see this, note that constraint  $\varphi_1^H$  ensures that every vertex is visited, while  $\varphi_2^H$  ensures that no vertex is visited more than once. Thus, the paths through any transition system resulting from the imposition of a behavioural constraint that is effective for  $O^H$  will be DHCS in  $H$ , with the agent ending up in  $v_{sink}$  and repeating the action  $\alpha_{loop}$  infinitely often. Since  $O^H$  is a formula of  $\mathcal{L}^u$ , we are done.

The final result—that the feasibility problem for objectives expressed as propositional logic formulae may be decided in polynomial time – is an immediate corollary of the following result, which shows that, for objectives expressed as formulae of propositional logic, the feasibility problem reduces directly to model checking in ATL.  $\square$

**Theorem 3** *Suppose  $\varphi$  is a propositional logic formula (representing an objective), and  $S$  is an AATS. Then  $S \models \langle \langle Ag \rangle \rangle \Box \varphi$  iff  $\varphi$  is feasible in  $S$ .*

<sup>1</sup> Of course, we strictly speaking do not need the start node, since if the graph contains a Hamiltonian cycle then we can start from any state—but it simplifies the exposition.

*Proof*

( $\rightarrow$ ) Assume  $S \models \langle\langle Ag \rangle\rangle \Box \varphi$ . By the semantics of ATL, we know that  $\exists \sigma_{Ag} \in \Sigma_{Ag}$  such that  $\forall \lambda \in \text{comp}(\sigma_{Ag}, q_0)$ , we have  $S, \lambda[u] \models \varphi$  for all  $u \in \mathbb{N}$ . In fact, since  $\sigma_{Ag}$  is a grand coalition strategy,  $\text{comp}(\sigma_{Ag}, q_0)$  will be a singleton; let  $\lambda^*$  denote its member; so  $\forall u \in \mathbb{N} : S, \lambda^*[u] \models \varphi$ . We must show that this implies there exists a behavioural constraint  $\beta$  such that  $S \dagger \beta \models \langle\langle \rangle\rangle \Box \varphi$ ; we construct  $\beta$  as follows:

$$\begin{aligned} &\text{for each } i \in Ag, \\ &\quad \text{for each } \alpha \in Ac_i, \\ &\quad \quad \text{set } \beta(\alpha) = \{q \mid \sigma_i(q) \neq \alpha\}. \end{aligned}$$

We claim that, for any  $q_0$ -computation  $\lambda$  of  $S \dagger \beta$ , we have  $S \dagger \beta, \lambda[u] \models \varphi$  for all  $u \in \mathbb{N}$ . To see this, observe that there will in fact be a single  $q_0$  computation of  $S \dagger \beta$ , namely  $\lambda^*$ , and we know that  $\forall u \in \mathbb{N}, S, \lambda^*[u] \models \varphi$ . We appeal to the fact that  $\varphi$  is a propositional formula, and that  $\pi$  is the same in  $S$  and  $S \dagger \beta$ , and conclude that for all  $u \in \mathbb{N} : S \dagger \beta, \lambda^*[u] \models \varphi$ .

( $\leftarrow$ ) Assume  $\varphi$  is feasible in  $S$ . Then there exists a behavioural constraint  $\beta$  such that  $S \dagger \beta, q_0 \models \langle\langle \rangle\rangle \Box \varphi$ . Let  $\Lambda$  denote the set of  $q_0$ -computations of  $S \dagger \beta$ . Then by the semantics of ATL, for all  $\lambda \in \Lambda$  and for all  $u \in \mathbb{N}$ , we have  $S \dagger \beta, \lambda[u] \models \varphi$ . We need to show that this implies  $S \models \langle\langle Ag \rangle\rangle \Box \varphi$ . By the semantics of ATL,  $S, q_0 \models \langle\langle Ag \rangle\rangle \Box \varphi$  iff  $\exists \sigma_{Ag} \in \Sigma_{Ag}$  such that  $\forall \lambda \in \text{comp}(q_0, \sigma_{Ag})$ , and  $\forall u \in \mathbb{N}$ , we have  $S, \lambda[u] \models \varphi$ . We show how to construct such a  $\sigma_{Ag}$ . We start by constructing from  $\beta$  a non-deterministic strategy  $\sigma_i^n : Q \rightarrow 2^{Ac_i}$  for each agent  $i \in Ag$ , as follows:

$$\begin{aligned} &\text{for each } i \in Ag, \\ &\quad \text{for each } q \in Q, \\ &\quad \quad \text{set } \sigma_i^n(q) = \{\alpha \mid q \notin \beta(\alpha) \wedge q \in \rho(\alpha)\}. \end{aligned}$$

Now, say a (conventional deterministic) strategy  $\sigma_i$  is consistent with  $\sigma_i^n$  if  $\sigma_i(q) \in \sigma_i^n(q)$ , for all  $q \in Q$ . Now, let  $\sigma_{Ag} = \langle \sigma_1, \dots, \sigma_n \rangle$  be any grand coalition strategy profile such that each  $\sigma_i$  is consistent with the corresponding non-deterministic strategy  $\sigma_i^n$ . We claim that, thus defined,  $\forall \lambda \in \text{comp}(\sigma_{Ag}, q_0)$  we have  $S, \lambda[u] \models \varphi$  for all  $u \in \mathbb{N}$ . To see this, observe that by construction we have  $\text{comp}(\sigma_{Ag}, q_0) \subseteq \Lambda$ , as defined above. We know that  $\forall \lambda \in \Lambda$ , and for all  $u \in \mathbb{N}$ , we have  $S \dagger \beta, \lambda[u] \models \varphi$ . Since  $\varphi$  is propositional, its valuation depends only upon the state in which it is interpreted. Since  $\pi$  is unchanged in  $S$  and  $S \dagger \beta$ , it must be that  $\forall \lambda \in \Lambda, \forall u \in \mathbb{N} : S, \lambda[u] \models \varphi$ .  $\square$

Notice that, as a direct corollary to Theorem 3, the synthesis problem for propositional logic objectives may also be solved in polynomial time: for if the answer to the model checking problem is “yes”, then the witness to this will be a strategy for the grand coalition from which, as the proof of the theorem illustrates, we can extract a behavioural constraint implementing the objective.

To appreciate that Theorem 3 does not hold for arbitrary formulae, consider the following simple one agent AATS (one may also assume this agent to be a grand coalition  $Ag$ ). Suppose we have to states,  $q_0$  and  $q_1$ , in the first, the atom  $p$  is false, in the second it is true. Moreover, we have four actions,  $a_0, a_1, b_0$  and  $b_1$ . Both  $a_i$  actions are possible in  $q_0$ , both  $b_i$  actions in  $q_1$ . Also, any action  $x_i$  ( $x \in \{a, b\}$ ) takes the system to  $q_i$ . More formally:  $\tau(q_0, a_0) = \tau(q_1, b_0) = q_0$ , and  $\tau(q_0, a_1) = \tau(q_1, b_1) = q_1$ . Thus, our agent  $i$  can always take the system to a  $p$ -state, but also he can ensure the system’s

next state will be one in which  $\neg p$  is true. Let  $\varphi = \neg p \wedge \langle\langle i \rangle\rangle \diamond p$ . In this system  $S$  we have  $S \models \langle\langle Ag \rangle\rangle \Box \varphi$  which expresses that the agent has a strategy to ensure that always  $\neg p$  is true, at the same time always having the opportunity to make  $p$  true. However, if we want any behavioural constraint  $\beta$  to be such that  $S \dagger \beta \models \langle\langle \rangle\rangle \Box \varphi$ , we see that this is impossible: to ensure that  $\varphi$  is always true in all runs,  $p$  has to be always false, so  $\beta$  has to forbid any transition to the state  $q_1$ , in which case the second conjunct of  $\varphi$ , i.e.,  $\langle\langle i \rangle\rangle \diamond p$  can never be made true anymore.

Theorem 3 illustrates a close relationship between the feasibility problem and model checking, which begs the question as to what extent the result can be extended for objectives beyond propositional logic. We have the following.

**Theorem 4** *Suppose  $v \in \mathcal{L}^u$  is a universal ATL formula (representing an objective), and  $S$  is an AATS. Then  $S \models \langle\langle Ag \rangle\rangle \Box v$  implies  $v$  is feasible in  $S$ .*

*Proof (Sketch)* The basic structure of the proof is as Theorem 3, but as  $v$  is no longer a propositional formula, but a universal ATL formula, we make use of Lemma 3.  $\square$

#### 4.1 Objectives with explicit action constraints

We now consider objectives for laws that explicitly refer to the legality or otherwise of actions. The following example illustrates the idea.

**Example 3** If  $\alpha$  is an action, then let the proposition  $\ell(\alpha)$  mean “action  $\alpha$  is legal”. Consider the following three objectives in the context of the trains system, as discussed earlier.

$$\begin{aligned} O_5 &= \neg(in_E \wedge in_W) \wedge \ell(move_E) \wedge \ell(move_W) \\ O_6 &= \neg(in_E \wedge in_W) \wedge \ell(move_E) \\ O_7 &= \neg(in_E \wedge in_W) \wedge \ell(move_W) \end{aligned}$$

Objective  $O_5$  requires that not only do the trains never crash, but that both trains are always able to move. Any behavioural constraint for this objective must not prevent the trains from moving if they choose to do so. Objective  $O_6$  is similar, but only requires that the Eastbound train is always able to move; and  $O_7$  only requires that the Westbound train is always able to move.

In our basic social laws framework, we have no way of expressing or reasoning about social laws with such objectives. We now show how the basic framework can be extended to include such constraints. In particular, we show how this can be done in the context of the MOCHA model checking system (Alur, Henzinger, Mang, Qadeer, Rajamani, & Taşiran, 1998). MOCHA takes as input a specification of an AATS, expressed in the REACTIVE MODULES language, and a formula of ATL, and is capable of either checking whether this formula is true in the AATS, or else giving a counter example. The actual syntax of the REACTIVE MODULES language used in MOCHA is rather complex, and so we adopt the following simplified syntax in the interests of easy comprehension. A REACTIVE MODULES agent (they are called “atoms” in the REACTIVE MODULES literature) has the following structure:

$$\begin{aligned} \text{agent } name \text{ reads } in \text{ writes } out \\ P_1 \mapsto \alpha_1; \\ P_2 \mapsto \alpha_2; \\ \dots \\ P_k \mapsto \alpha_k. \end{aligned} \tag{2}$$

where  $name$  is the name of the agent,  $in \subseteq \Phi$  is a list of boolean variables that the agent observes,  $out \subseteq \Phi$  is a list of boolean variables that it controls, and provides to the rest of the system, and the  $P_j \mapsto \alpha_j$  structures are *guarded commands*, where  $P_j$  is a predicate over the variables the agent observes and controls (i.e., a boolean expression over variables  $in \cup out$ ), and  $\alpha_j \in Ac_{name}$  is an action. Actions can be understood as functions that take as input the variables visible to the agent ( $in \cup out$ ), and produce as output an assignment for the variables controlled by the agent ( $out$ ). The idea is that at each time step, each agent generates the set of rules whose pre-conditions are satisfied by the current state of the system. One of these rules is non-deterministically chosen for execution – this involves simply executing the corresponding action, which produces an assignment for the variables under its control ( $out$ ), which is the value they take in the next state of the system. Of course, a guarded command may have  $\top$  as a pre-condition, in which case it is always enabled for execution.

It should be clear how a collection of such agents maps to an AATS, as defined in Sect. 2. The most important point is that the pre-conditions to guarded commands correspond to the pre-condition function  $\rho$ : given a state of the system  $q$  and guarded command  $P \mapsto \alpha$ , then  $q \in \rho(\alpha)$  iff  $q \models P$ .

We now show how an AATS, expressed in such a framework, can be extended to permit objectives referring to the legality or otherwise of actions. The idea is to define a transformation on AATSS:

$$\cdot^\circ : \text{AATS} \rightarrow \text{AATS}$$

such that the transformed system includes the  $\ell(\alpha)$  propositions as before, but otherwise has the same properties.

Given an AATS  $S = \langle Q, q_0, Ag, Ac_1, \dots, Ac_n, \rho, \tau, \Phi, \pi \rangle$ , defined using the REACTIVE MODULES language as above, the system  $S^\circ$  is created as follows.

First, for each action  $\alpha \in Ac$ , we create:

- A new propositional variable  $\ell(\alpha)$ , with the intended interpretation that  $\ell(\alpha)$  will be true in a state if the action  $\alpha$  is legal according to the behavioural constraint that we are looking for.
- A new agent, which controls the variable  $\ell(\alpha)$ , as follows.

$$\begin{aligned} \text{agent } \alpha\text{-controller reads } \Phi \text{ writes } \ell(\alpha) \\ \top \mapsto \ell(\alpha) := \top; \\ \top \mapsto \ell(\alpha) := \perp. \end{aligned}$$

Thus, in every possible state, the agent  $\alpha$ -controller has two actions available: make  $\ell(\alpha)$  true, or make  $\ell(\alpha)$  false. Note that because the condition on the guard of each of these actions is  $\top$ , the  $\alpha$ -controller agent can *always* perform these actions. We will denote by *controllers* the set of all controller agents introduced in this way.

We then take each of the original agents  $i$  in the system, and transform it as follows.

- We replace each guarded command  $P \mapsto \alpha$  for agent  $i$  with a rule as follows.

$$\ell(\alpha)' \wedge P \mapsto \alpha$$

Notice that the first conjunct of the guard is “primed”; this notation in REACTIVE MODULES means the value that  $\ell(\alpha)$  has been assigned *in the current round*, i.e., the value of this variable after it has been assigned by the  $\alpha$ -controller agent. Thus agent  $i$  may perform the action  $\alpha$  iff the original guard condition  $P$  is true in the



current state, and the new agent controlling  $\ell(\alpha)$  has assigned this variable the value  $\top$  in the current state. In this way, the  $\alpha$ -controller agent can determine whether or not  $\alpha$  is performed.

- Next, if  $in$  is the set of variables that  $i$  reads, then we replace  $in$  by  $in \cup \{\ell(\alpha) \mid \alpha \in Ac_i\}$ . In other words, the agent  $i$  reads those  $\ell$ -variables that determine whether its actions are legal.

Notice that agent  $i$  is not dependent on the  $\ell$ -variables of any other agent’s actions, and hence any strategy computed for  $i$  by MOCHA will not be dependent on the  $\ell$ -variables of other agents.

Now, we can prove:

**Theorem 5** *Suppose  $\varphi$  is a propositional logic formula (representing an objective), and  $S$  is an AATS. Then*

$$S^\circ \models \langle\langle \text{controllers} \rangle\rangle \Box \left( \varphi \wedge \left[ \bigwedge_{i \in Ag} \bigvee_{\alpha \in Ac_i} \ell(\alpha) \right] \right)$$

iff  $\varphi$  is feasible in  $S$ .

*Proof (Sketch)* The structure of the proof is as Theorem 3: if  $S^\circ, q_0 \models \langle\langle \text{controllers} \rangle\rangle \Box \varphi$ , then the controller agents *controllers* have a strategy for enabling/disabling actions such that if this strategy is followed, the law is implemented. From this strategy, we can extract a behavioural constraint that implements the law. The additional condition ensures that every agent has at least one action available in every state, and so the strategy controlling the  $\ell(\alpha)$  variables really does map to a behavioural constraint, as defined earlier. Conversely, if there is a behavioural constraint  $\beta$  such that  $(\varphi, \beta)$  is an effective social law in  $S^\circ$ , then we can extract a strategy for the *controllers* such that this strategy ensures  $\varphi$  is always true.  $\square$

Notice that the obvious analogue of Theorem 4 also holds, although space restrictions prevent more discussion.

**Proposition 6** *Objective  $O_5$  is not feasible in  $S_1$ , while both objectives  $O_6$  and  $O_7$  are feasible.*

Using this framework, we can also investigate more general properties of social laws. Consider behavioural constraints that act as *dictatorships* (cf. [Pauly, 2002, p. 17]). A behavioural constraint is a dictatorship if, once it is implemented, it never presents an agent with more than one possible action at any given time; that is, if an agent has no choice about what action it may perform. Formally, a behavioural constraint  $\beta$  with respect to an AATS  $S$  is a dictatorship if

$$\forall i \in Ag, \forall q \in Q, |options(i, q)| = 1$$

where  $options(i, q)$  is evaluated in  $S \dagger \beta$ .

**Example 4** The following is a dictatorship behavioural constraint for the trains scenario:

$$\beta_4(\alpha) = \begin{cases} \{q_0, q_4, q_6\} & \text{if } \alpha = idle_E \\ \{q_0, q_5, q_6\} & \text{if } \alpha = idle_W \\ \{q_5\} & \text{if } \alpha = move_E \\ \{q_4\} & \text{if } \alpha = move_W \end{cases}$$

**Theorem 6** Suppose  $\varphi$  is a propositional logic formula (representing an objective), and  $S$  is an  $\Delta$ ATS. Then  $\varphi$  is feasible in  $S$  by a dictatorship behavioural constraint iff

$$S^\circ \models \langle\langle \text{controllers} \rangle\rangle \Box \left( \varphi \wedge \bigwedge_{i \in \text{Ag}} \left[ \bigvee_{\alpha \in \text{Act}_i} \ell(\alpha) \right] \wedge \left[ \bigwedge_{\substack{\alpha_1 \in \text{Act}_i \\ \alpha_2 \in \text{Act}_i \\ \alpha_1 \neq \alpha_2}} \neg(\ell(\alpha_1) \wedge \ell(\alpha_2)) \right] \right).$$

*Proof* (Sketch) As before, except that we have an additional constraint, which ensures that no more than one action from an agent's action set is enabled at any given time.  $\square$

**Proposition 7**  $(O_2, \beta_4)$  and  $(O_3, \beta_4)$  are effective social laws in  $S_1$ .

## 5 Related work

The closest approach to ours in the literature is the original framework of Moses, Shoham, and Tennenholtz. Shoham and Tennenholtz were the first to precisely articulate the notion of social laws for multiagent systems, and set up a basic formal framework within which computational questions about social laws could be formulated (Shoham & Tennenholtz, 1992; Shoham & Tennenholtz, 1996; Shoham & Tennenholtz, 1997). The particular application domain was that of traffic laws for robotic agents. The basic framework was extended by Fitoussi and Tennenholtz, to consider *simple* social laws – essentially, social laws that could not be any simpler without failing (Fitoussi & Tennenholtz, 2000). In related work, Moses and Tennenholtz developed a deontic epistemic logic for representing properties of multiagent systems with normative structures in place (Moses & Tennenholtz, 1995). Although semantically similar to ATL (and ATEL (van der Hoek & Wooldridge, 2003b)), their logic was quite different to ATL in terms of the syntactic constructs it provided, and the emphasis was primarily on deriving axioms capturing static aspects of artificial social systems and social laws.

## 6 Conclusions

In this paper, we have demonstrated how the Alternating-time Temporal Logic of Alur, Henzinger, and Kupferman provides a natural and powerful framework within which to express and reason about social laws. Following a formulation of social laws within ATL, we demonstrated how the effectiveness, feasibility, and synthesis problems could be understood as model checking problems for ATL, and also demonstrated how our basic framework could naturally be extended to include social laws that require explicit constraints on actions.

There are many possible routes for future investigation. One obvious question is the extent to which other notions such as knowledge can be incorporated into the framework. For example, one could imagine a social law with the objective that certain facts about agents must remain private to those agents; epistemic extensions to ATL seem to provide a natural language for representing such objectives (van der Hoek & Wooldridge, 2003a, b).

## References

- Alur, R., Henzinger, T. A., & Kupferman, O. (2002). Alternating-time temporal logic. *Journal of the ACM*, 49(5), 672–713.
- Alur, R., Henzinger, T. A., Mang, F. Y. C., Oadeer, S., Rajamani, S. K., & Taşiran, S. (1998). Mocha: Modularity in model checking. In *CAV 1998: Tenth International Conference on Computer-aided Verification (LNCS Volume 1427)* (pp. 521–525). Berlin, Germany: Springer-Verlag.
- Blackburn, P., de Rijke, M., & Venema, Y. (2001). *Modal Logic*. Cambridge, England: Cambridge University Press.
- Clarke, E. M., Grumberg, O., & Peled, D. A. (2000). *Model Checking*. Cambridge, MA: The MIT Press.
- Dignum, F. (1999). Autonomous agents with norms. *Artificial Intelligence and Law*, 7, 69–79.
- van Drimmelen, G. (2003). Satisfiability in alternating-time temporal logic. In *Eighteenth Annual IEEE Symposium on Logic in Computer Science (LICS 2003)* (pp. 208–217). Ottawa, Canada.
- Durfee, E. H. (1999). Distributed problem solving and planning. In G. Weiß (Ed.), *Multiagent systems* (pp. 121–164). Cambridge, MA: The MIT Press.
- Emerson, E. A. (1990). Temporal and modal logic. in J. van Leeuwen (Ed.), *Handbook of theoretical computer science volume B: Formal models and semantics* (pp. 996–1072). Amsterdam, The Netherlands: Elsevier Science Publishers B.V.
- Fitoussi, D., & Tennenholtz, M. (2000). Choosing social laws for multi-agent systems: Minimality and simplicity. *Artificial Intelligence*, 119(1–2), 61–101.
- van der Hoek, W., & Wooldridge, M. (2002). Tractable multiagent planning for epistemic goals. In *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS-2002)* (pp. 1167–1174). Bologna, Italy.
- van der Hoek, W., & Wooldridge, M. (2003a). Model checking cooperation, knowledge, and time—a case study. *Research in Economics*, 57(3), 235–265.
- van der Hoek, W., & Wooldridge, M. (2003b). Time, knowledge, and cooperation: Alternating-time temporal epistemic logic and its applications. *Studia Logica*, 75(1), 125–157.
- Moses, Y., & Tennenholtz, M. (1995). Artificial social systems. *Computers and AI*, 14(6), 533–562.
- Papadimitriou, C. H. (1994). *Computational complexity*. Reading, MA: Addison-Wesley.
- Pauly, M. (2001). *Logic for Social Software*. PhD thesis, University of Amsterdam, ILLC Dissertation Series 2001–10.
- Pauly, M. (2002). A modal logic for coalitional power in games. *Journal of Logic and Computation*, 12(1), 149–166.
- Shoham, Y., & Tennenholtz, M. (1992). On the synthesis of useful social laws for artificial agent societies. In *Proceedings of the Tenth National Conference on Artificial Intelligence (AAAI-92)*, San Diego, CA.
- Shoham, Y., & Tennenholtz, M. (1996). On social laws for artificial agent societies: Off-line design. In P. E. Agre, & S. J. Rosenschein (Eds.), *Computational theories of interaction and agency* (pp. 597–618). Cambridge, MA: The MIT Press.
- Shoham, Y., & Tennenholtz, M. (1997). On the emergence of social conventions: Modelling, analysis, and simulations. *Artificial Intelligence*, 94(1–2), 139–166.
- Wooldridge, M. (2002). *An introduction to multiagent systems*. John Wiley & Sons.