

# A formalisation of violation, error recovery, and enforcement in the bit transmission problem

Alessio Lomuscio<sup>1</sup>      Marek Sergot<sup>2</sup>

<sup>1</sup>Department of Computer Science, King's College London  
London WC2R 2LS, UK. alessio@dcs.kcl.ac.uk

<sup>2</sup>Department of Computing, Imperial College London  
London SW7 2BZ, UK. m.sergot@imperial.ac.uk

## Abstract

The design of complex multi-agent systems is increasingly having to confront the possibility that agents may not behave as they are supposed to. In addition to analysing the properties that hold if protocols are followed correctly, it is also necessary to predict, test, and verify the properties that would hold if these protocols were to be violated. We illustrate how the formal machinery of deontic interpreted systems can be applied to the analysis of such problems by considering three variations of the bit transmission problem. The first, an example in which an agent may fail to do something it is supposed to do, shows how we deal with violations of protocols and specifications generally. The second, an example in which an agent may do something it is not supposed to do, shows how it is possible to specify and analyse remedial or error-recovery procedures. The third combines both kinds of faults and introduces a new component into the system, a controller whose role is to enforce compliance with the protocol. In each case the formal analysis is used to test whether critical properties of the system are compromised, in this example, the reliable communication of information from one agent to the other.

## 1 Introduction

The design of complex multi-agent systems is increasingly having to confront the possibility that agents may not behave as they are supposed to. In e-commerce, in security, in automatic negotiation, in any application where agents are programmed by different parties with competing interests, it is unrealistic to assume that all agents will behave according to some given protocol or standard of behaviour. In addition to analysing the properties that hold if protocols are followed correctly, it is also necessary to predict, test, and verify the properties that would hold if these protocols were to be violated, and to test the effectiveness of introducing proposed enforcement mechanisms.

Consider an online auction mechanism modelled and implemented as a multi-agent system, where agents play the parts of auctioneers and bidders. At each round the protocol of the system gives the opportunity to a selected number of agents to bid for some goods, according to the specific auction protocol being employed. Suppose that an agent bids an amount which it is unable to fulfil

(either by mistake, or perhaps with the sole intention of raising the bid so that an opponent will have to pay a higher price). If that bid is made in the prescribed manner then it would still count as a valid bid—even if it were to be regarded as illegal, anti-social, unethical to make a bid in those circumstances. Undesirable behaviour of this kind will often have adverse effects on the system as a whole. Perhaps the agent will be forced to de-commit from the commitment it entered upon by bidding, as a result of which the seller will (at least temporarily) lose the deal, which in turn may have consequences on other deals it has entered upon with other agents. At the very least, the resulting disputes will have to be resolved at a cost.

One way out of such problems is to attempt to devise tighter controls in the protocols, such as stipulating that a bid by an agent is valid only if made with previously cleared funds. The opportunities for inventing such controls are limited, however, and in any case have associated costs too (in the example, the costs of the clearing mechanism and the consequently diminished liquidity of the market). Moreover forcing agents of an open multi-agent system to behave according to a rigid set of rules is seen as increasingly unfeasible because of the inherent distribution of the system. One alternative, having specified what is correct, desirable, or acceptable behaviour of the agents within a given context, is to strive to handle and reason about violations of these norms within the system. The required formalisms not only have to be capable of handling *local violations*, but must also be capable of representing the usual attributes ascribed to agents in a multi-agent setting, such as knowledge, intentions, and goals (see e.g., [Woo00]).

It is important to emphasise that the focus of the approach presented in this paper is not on formalisms used *internally* by the agents as part of their reasoning mechanisms, but rather for specification of properties of the system as a whole, ideally accompanied by verification mechanisms. The question of how an agent may reason internally about the norms that constrain its behaviour is an interesting one, but it is a different question that will not be pursued in this paper. Note that it is legitimate and perfectly meaningful to make this separation. Consider the familiar problem of controlling access to sensitive or confidential data. There is a specification of what access is permitted or forbidden at the system level but the agents who actually make the access need not be aware, and usually are not aware, of what this specification is.

These are inherently complex issues, and we do not expect to have methods able to give comprehensive answers to all of them in the near future. One approach being explored by the authors as a step in this direction is an attempt to extend interpreted systems [FHMV95], a mainstream and well-developed semantics for reasoning about knowledge in distributed and multi-agent systems, with a deontic component. Deontic interpreted systems [LS03] are designed to represent correct and incorrect functioning behaviour of agents as well as their epistemic properties.

While a considerable amount of theoretical work considers constructs not dissimilar to those employed here (among others: [ML85, Mey88, KM89, FM91, GMP92, Coe93, JS93, CJ96]), we are not aware of any study that puts these theoretical constructs to the test in an actual application. It is surprisingly difficult to find concrete examples in the literature, let alone case studies showing how the formal approaches may be applied in detail. Our aim in this paper is to advance the state-of-the-art by constructing a fully worked out formal analysis

of a concrete example, and to illustrate and evaluate the formalism of deontic interpreted systems in the process.

Specifically, we show how the formalism of deontic interpreted systems [LS03] may be applied to model and reason about protocol violations in a simple example, the bit-transmission problem [HZ92]. This is a much discussed scenario in distributed computing involving two agents attempting to communicate the value of a bit over a faulty communication channel. Of course, the example is trivial compared to the kinds of applications alluded to in the introductory sections. Still, we argue that it is rich enough to raise many important issues such as error correction and control/enforcement mechanisms, and in this sense is representative. Furthermore, using a small, well understood, and widely discussed problem offers the possibility of testing rigorously a formal apparatus and investigating how it performs in detail.

We will examine three variations of the bit-transmission problem. The first (Section 5) is an example in which an agent may fail to do something it is supposed to do, and shows how we deal with violations of protocols and specifications generally. The second (Section 6) is an example in which an agent may do something it is not supposed to do, and shows how it is possible to specify and analyse remedial or error-recovery procedures. The third (Section 7) combines both kinds of faults and introduces a new component into the system, a controller whose role is to *enforce* compliance with the protocol. In each case the formal analysis is used to test whether critical properties of the system are compromised, in this example, the reliable communication of information from one agent to the other.

## 2 Preliminaries

In this paper we use the machinery of interpreted systems as presented in [FHMV95], and the extensions for modelling correct and incorrect functioning behaviours as developed in [LS03]. We present the main definitions here, but refer to the cited literature for more details.

Consider  $n$  non-empty sets  $L_1, \dots, L_n$  of local states, one for every agent of the system, and a set of states for the environment  $L_E$ . Elements of  $L_i$  will be denoted by  $l_1, l'_1, l_2, l'_2, \dots$ . Elements of  $L_E$  will be denoted by  $l_E, l'_E, \dots$ .

**Definition 1 (System of global states)** *A system of global states for  $n$  agents  $\mathcal{S}$  is a non-empty subset of a Cartesian product  $L_1 \times \dots \times L_n \times L_E$ .*

*An interpreted system of global states is a pair  $IS = (\mathcal{S}, h)$  where  $\mathcal{S}$  is a system of global states and  $h: \mathcal{S} \rightarrow 2^{\mathcal{P}}$  is an interpretation function for a set of propositional variables  $\mathcal{P}$ .*

*When  $g = (l_1, \dots, l_n, l_E)$  is a global state of a system  $\mathcal{S}$ ,  $l_i(g)$  denotes the local state of agent  $i$  in global state  $g$ .  $l_E(g)$  denotes the local state of the environment in global state  $g$ .*

Systems of global states can be used to interpret epistemic modalities  $K_i$ , one for each agent.

$$(IS, g) \models K_i \varphi \quad \text{if for all } g' \text{ we have that } l_i(g) = l_i(g') \text{ implies } (IS, g') \models \varphi.$$

Alternatively one can consider generated models  $(\mathcal{S}, \sim_1, \dots, \sim_n, h)$  of the standard relational (Kripke) kind, where the equivalence relations  $\sim_i$  are defined on

equivalence of local states, i.e.,  $g \sim_i g'$  iff  $l_i(g) = l_i(g')$ . Modalities are then interpreted on these models in the usual way. The resulting logic for modalities  $K_i$  is of type  $S5_n$ ; this models agents with complete introspection capabilities and veridical knowledge [Hin62, MH95], though it is important to see that the notion of knowledge modelled here is an information-theoretic one. An agent ‘knows’  $\varphi$  when from an *external* perspective it has enough information to determine that  $\varphi$  is true.

The notion of interpreted systems can be extended to incorporate the idea of correct functioning behaviour of some or all of the components [LS03], as follows.

**Definition 2 (Deontic system of global states)** *Given  $n$  agents and  $n + 1$  non-empty sets  $G_1, \dots, G_n, G_E$ , a deontic system of global states is any system of global states defined on  $L_1 \supseteq G_1, \dots, L_n \supseteq G_n, L_E \supseteq G_E$ . For any agent  $i$ ,  $G_i$  is called the set of green states for agent  $i$ .  $G_E$  is called the set of green states for the environment. The complement of  $G_i$  with respect to  $L_i$  (respectively  $G_E$  with respect to  $L_E$ ) is called the set of red states for agent  $i$  (respectively for the environment).*

The terms ‘green’ and ‘red’ are chosen as neutral terms, to avoid overloading them with unintended readings and connotations. The term ‘green’ can be read as ‘legal’, ‘acceptable’, ‘desirable’, ‘correct’, depending on the context of a given application.

Note that global states can also be classified as green/red. Let  $\mathcal{G}_i$  denote the set of global states in which agent  $i$  is in a green local state:  $g \in \mathcal{G}_i$  iff  $l_i(g) \in G_i$ .

Deontic systems of global states are used to interpret modalities such as the following

$$(IS, g) \models O_i \varphi \quad \text{if for all } g' \text{ we have that } l_i(g') \in G_i \text{ implies } (IS, g') \models \varphi \\ \text{(or equivalently: if for all } g' \in \mathcal{G}_i \text{ we have } (IS, g') \models \varphi).$$

$O_i \varphi$  is used to represent that  $\varphi$  holds in all (global) states in which agent  $i$  is functioning correctly. Again, one can consider generated models of the standard form  $(S, \sim_1, \dots, \sim_n, R_1^O, \dots, R_n^O, h)$  where the equivalence relations  $\sim_i$  are defined as above and the relations  $R_i^O$  are defined by  $g R_i^O g'$  if  $l_i(g') \in G_i$ . In particular, the semantics of  $O_i$  is essentially that of an obligation operator in (a strengthened form of) standard deontic logic: each  $O_i$  is a normal modality of type  $KD45$ ; the multi-modal case can be axiomatised by the system  $KD45^{i-j}$  [LS03]. However, it would not be appropriate to read the expression  $O_i \varphi$  as ‘there is an obligation on agent  $i$  that  $\varphi$ ’.

Knowledge can be modelled on deontic interpreted systems in the same way as on interpreted systems, and one can study various combinations of the modalities such as  $K_i O_j$ ,  $O_j K_i$ , and others. It is particularly important when reading these expressions, however, to remember that they express the external “bird’s eye” view of the system:  $O_j K_i \varphi$  says that in all states of the system in which agent  $j$  is functioning correctly, from an external perspective agent  $i$  has sufficient information to know that  $\varphi$ . There are many other senses of ‘ought to know’ that are not captured by this construction.

Another concept of particular interest is knowledge that an agent  $i$  has *on the assumption that the system (the environment, agent  $j$ , group of agents  $X$ )*

is functioning correctly. We employ the (doubly relativised) modal operator  $\widehat{K}_i^j$  for this notion, interpreted as follows:

$$(IS, g) \models \widehat{K}_i^j \varphi \quad \text{if for all } g' \text{ such that } l_i(g) = l_i(g') \text{ and } l_j(g') \in G_j \\ \text{we have that } (IS, g') \models \varphi.$$

$K_i$ ,  $O_i$ , and  $\widehat{K}_i^j$  are all normal modalities (of type S5,  $KD45$ , and  $KD45$ , respectively). Various relationships hold between them, such as  $K_i \varphi \rightarrow \widehat{K}_i^j \varphi$ ,  $K_i O_j \varphi \rightarrow \widehat{K}_i^j \varphi$ ,  $O_j \varphi \rightarrow \widehat{K}_i^j \varphi$ , and others. Further details, and metalogical properties such as frame correspondence and completeness, may be found in [LS03]. Our aim in this present paper is to illustrate the use of these modalities in the examples to follow.

Finally, interpreted systems can be extended to deal with temporal evolution. Consider a set of runs over global states  $R = \{r: NS\}$ , representing flows of time for the system. When this structure is in place one can interpret the usual temporal connectives on it (see e.g. [GHR93]). In this paper we limit our analysis to the static case.

### 3 The bit transmission problem

The bit-transmission problem [FHMV95] involves two agents, a *sender*  $S$ , and a *receiver*  $R$ , communicating over a possibly faulty communication channel.  $S$  wants to communicate some information—the value of a bit for the sake of the example—to  $R$ . We would like to design a protocol that accomplishes this objective while minimising the use of the communication channel.

Two (trivial) special cases of the scenario can immediately be solved. The first is the one in which the channel is actually working, or is at least operative for the first few rounds of computation. In this case we would require  $S$  to send the value of the bit once as the system comes alive. The other special case arises when the channel is constantly non-operative. There is obviously no protocol that can ensure that  $R$  receives the information in that case.

The interesting scenario arises when the channel is working correctly at certain times while failing at others. There are several ways in which this can be approached, for instance by stipulating that the channel delivers messages with a fixed probability  $P > 0$  at any given round. In this paper we do not make any of these assumptions; instead, we analyse the most general case with respect to a most simple protocol. The protocol is as follows.  $S$  immediately starts sending the bit to  $R$ , and continues to do so until it receives an acknowledgement from  $R$ .  $R$  does nothing until it receives the bit; from then on it sends acknowledgements of receipt to  $S$ .  $S$  stops sending the bit to  $R$  when it receives an acknowledgement. Note that  $R$  will continue sending acknowledgements even after  $S$  has received its acknowledgement. Intuitively  $S$  will know for sure that the bit has been received by  $R$  when it gets an acknowledgement from  $R$ .  $R$ , on the other hand, will never be able to know whether its acknowledgement has been received since  $S$  does not answer the acknowledgement<sup>1</sup>. We refer to [HZ92, FHMV95] for further discussion.

<sup>1</sup>One might think that this problem can be solved by insisting that  $S$  sends an acknowledgement of the acknowledgement, but by doing so we simply push the problem one level deeper, and  $S$  would never know whether its acknowledgement of the acknowledgement has been received by  $R$ .

## 4 Analysis

The bit-transmission problem can be analysed using the formalism of interpreted systems. To do this we follow the approach taken by Halpern and colleagues in [FHMV95].

There are three active components in the scenario: a sender, a receiver, and a communication channel. In line with the spirit of the formalism of interpreted systems, it is convenient to see sender and receiver as agents, and the communication channel as the environment. Each of these can be modelled by considering their local states. For the sender  $S$ , it is enough to consider four possible local states. They represent the value of the bit that  $S$  is attempting to transmit, and whether or not  $S$  has received an acknowledgement from  $R$ . Three different local states are enough to capture the state of  $R$ : the value of the received bit, if any, and the circumstance in which no bit has been received yet, represented by  $\epsilon$ . So we have

$$L_S = \{0, 1, 0\text{-ack}, 1\text{-ack}\}, \quad L_R = \{0, 1, \epsilon\}.$$

To model the environment it is often convenient to represent in the local state (amongst other things) the combinations of messages that were sent in the previous round, by  $S$  and  $R$ , respectively. The four local states would then be:

$$L_E = \{(\cdot, \cdot), (\text{sendbit}, \cdot), (\cdot, \text{sendack}), (\text{sendbit}, \text{sendack})\},$$

where ‘ $\cdot$ ’ represents configurations in which no message was sent by the corresponding agent. However, for the examples we will be analysing in this paper this level of detail is unnecessary, and so for simplicity we shall take the local states of the environment to be just a singleton:

$$L_E = \{\cdot\}.$$

This is to simplify (and shorten) the presentation. It should be obvious how the formulas of later sections can be adjusted for examples requiring a more complicated representation of the environment’s local states.

Global states  $\mathcal{S}$  for the system are defined as  $\mathcal{S} \subseteq L_S \times L_R \times L_E$ . A global state  $g = (l_S, l_R, l_E)$  gives a snapshot of the system at a given time. Note that not all triples of the product are admissible in principle, but only those that can be reached in a run of the protocol from given initial configurations, as will be explained below. Although in this simple example the  $L_E$  component is fixed, it is not redundant, and we will retain it in the global states as a reminder that reference to it would be required in more complicated settings.

The interpreted systems framework has the advantage of being suitable for integration with finer semantics representing actions and protocols. To do so consider a set of actions  $Act_i$  for every agent in the system and the environment. The actions for the agents  $S$  and  $R$  are as follows.

$$Act_S = \{\text{sendbit}, \lambda\}, \quad Act_R = \{\text{sendack}, \lambda\}$$

Here  $\lambda$  stands for no action (‘no-op’).

The actions  $Act_E$  for the environment correspond to the transmission of messages between  $S$  and  $R$  on the unreliable communication channel. To make the example sufficiently rich, we will assume that the communication channel can

transmit messages in both directions simultaneously, and that a message travelling in one direction can get through while a message travelling in the opposite direction is lost. (Alternatively, think of a pair of unidirectional communication channels whose faults are independent of one other.) The set of actions for the environment is

$$Act_E = \{\leftrightarrow, \rightarrow, \leftarrow, -\}$$

$\leftrightarrow$  represents the action in which the channel transmits any message successfully in both directions,  $\rightarrow$  that it transmits successfully from  $S$  to  $R$  but loses any message from  $R$  to  $S$ ,  $\leftarrow$  that it transmits successfully from  $R$  to  $S$  but loses any message from  $S$  to  $R$ , and  $-$  that it loses any messages sent in either direction. Alternatively, we could model the communications channel by having a single action for the environment, *transmit*, and then representing the unreliability of the channel by four possible local states of the environment. It comes to much the same thing, except that this alternative treatment makes the specification of the rest of the system slightly more cumbersome.

We can model the evolution of the system by means of a transition function  $\pi: \mathcal{S} \times Act \rightarrow 2^{\mathcal{S}}$ , where  $Act = Act_S \times Act_R \times Act_E$  is the set of joint actions for the system:  $\pi(g, (\alpha_S, \alpha_R, \alpha_E))$  is the set of global states resulting (non-deterministically) from performance in global state  $g$  of actions  $\alpha_S$ ,  $\alpha_R$ ,  $\alpha_E$  by sender  $S$ , receiver  $R$ , and the environment, respectively.  $\pi$  codes the fact that the state of the communication channel determines whether the actions performed by the agents (i.e., the messages they send on the channel) get through or not, and what their effects are. For example, the definition of  $\pi$  contains the following<sup>2</sup>:

$$\begin{aligned} \pi((0, \epsilon, \cdot), (sendbit, \lambda, \leftrightarrow)) &= (0, 0, \cdot), \\ \pi((0, \epsilon, \cdot), (sendbit, \lambda, \rightarrow)) &= (0, 0, \cdot), \\ \pi((0, \epsilon, \cdot), (sendbit, \lambda, \leftarrow)) &= (0, \epsilon, \cdot), \\ \pi((0, \epsilon, \cdot), (sendbit, \lambda, -)) &= (0, \epsilon, \cdot), \end{aligned}$$

to capture that when the channel works properly the message *sendbit* from  $S$  gets through and gets processed accordingly by  $R$ .

Some other cases are as follows:

$$\begin{aligned} \pi((0, 0, \cdot), (sendbit, sendack, \alpha_E)) &= (0-ack, 0, \cdot) & (\alpha_E \in \{\leftarrow, \leftrightarrow\}) \\ \pi((0, 0, \cdot), (sendbit, sendack, \alpha_E)) &= (0, 0, \cdot) & (\alpha_E \in \{\rightarrow, -\}) \end{aligned}$$

The remaining cases can be similarly expressed and are added straightforwardly. We leave the details to the reader. The transition function  $\pi$  is depicted, in simplified form, in Figure 1.

For compliance with a given protocol, only certain actions are performable at a given time for an agent. For example if sender  $S$  has not yet received an acknowledgement from receiver  $R$ , i.e., in the model, when  $S$  is in the local state 0 or 1, then according to the simple protocol under consideration,  $S$  should send the value of the bit over the channel to  $R$ , i.e., perform the action *sendbit*. To capture such requirements the concept of *protocol* can be used. A protocol for agent  $i$  is a function  $P_i: L_i \rightarrow 2^{Act_i}$  mapping local states to sets of actions.  $P_i(l_i)$  is the set of actions performable according to the protocol by agent  $i$

<sup>2</sup>We omit brackets when writing singleton sets to reduce clutter.

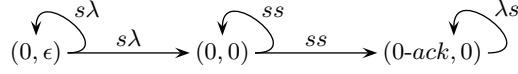


Figure 1: The transition function  $\pi$  for the bit-transmission system. The case  $\text{bit}=1$  is similar and is omitted. Labels  $s\lambda$ ,  $\lambda s$  and  $ss$  stand for the joint actions (*sendbit*,  $\lambda$ ), ( $\lambda$ , *sendack*) and (*sendbit*, *sendack*), respectively. Actions of the environment  $Act_E$  are omitted for simplicity.

when its local state is  $l_i$ . For the example under consideration the protocol can be defined as follows:

$$\begin{aligned} P_S(0) &= P_S(1) = \textit{sendbit}, & P_S(0\text{-ack}) &= P_S(1\text{-ack}) = \lambda, \\ P_R(\epsilon) &= \lambda, & P_R(0) &= P_R(1) = \textit{sendack}. \end{aligned}$$

(As usual, singleton sets are written here without brackets.) For the environment, for this simple example, we use the constant function:

$$P_E(l_E) = Act_E = \{\leftrightarrow, \rightarrow, \leftarrow, -\}, \quad \text{for all } l_E \in L_E.$$

If we assume that the system starts from an initial state  $(0, \epsilon, \cdot)$  or  $(1, \epsilon, \cdot)$ , it is possible to show that  $S$  will start sending the bit and will only stop after having received an acknowledgement from  $R$ . In turn  $R$  will remain silent until it receives the bit, and it will never stop sending acknowledgements from then on. The analysis can be made formally by using the mechanism of *contexts* [FHMV97] and implemented in model checking languages such as in [LRS02]; for our purposes it is not necessary to pursue that analysis here.

**Definition 3 (Run, reachable global state,  $P_i$ -compliant run)** Let  $\mathcal{S}$  be a set of global states,  $\pi: \mathcal{S} \times Act \rightarrow 2^{\mathcal{S}}$  a transition function,  $P_i$  a protocol function for the agent  $i$ , and  $\mathcal{I}$  a set of initial configurations  $\mathcal{I} \subseteq \mathcal{S}$ . A run is any sequence

$$g_0 \gamma_1 g_1 \dots g_{k-1} \gamma_k g_k \dots g_{n-1} \gamma_n g_n \quad (n \geq 0)$$

where  $g_0 \in \mathcal{I}$ , each  $\gamma_k$  ( $k = 1 \dots n$ ) is a joint action in  $Act$ , and  $g_k \in \pi(g_{k-1}, \gamma_k)$ . A global state  $g$  is  $\pi$ -reachable if  $g = g_n$  for some run  $g_0 \gamma_1 g_1 \dots g_{n-1} \gamma_n g_n$ .

Let  $[\gamma]_i$  denote the action performed by agent  $i$  in global transition  $\gamma$ , i.e., when  $\gamma = (\alpha_1, \dots, \alpha_n, \alpha_E)$ ,  $[\gamma]_i = \alpha_i$ .

A run  $g_0 \gamma_1 g_1 \dots g_{n-1} \gamma_n g_n$  is  $P_i$ -compliant when  $[\gamma_k]_i \in P_i(l_i(g_{k-1}))$  for each  $k = 1 \dots n$ . A global state  $g$  is  $P_i$ -reachable if  $g = g_n$  for some  $P_i$ -compliant run  $g_0 \gamma_1 g_1 \dots g_{n-1} \gamma_n g_n$ .

The set of global states reachable from the initial configurations  $(0, \epsilon, \cdot)$  and  $(1, \epsilon, \cdot)$  as defined by the transition function  $\pi$  and the protocol functions  $P_S$ ,  $P_R$  and  $P_E$ , is summarised in Figure 2. The environment component is omitted for clarity.

Having defined the set of reachable global states, we can apply the tools of formal logic to the analysis of the scenario by considering an interpretation of a suitably chosen set of propositional variables. We shall use the set  $\mathcal{P} = \{\mathbf{bit} = \mathbf{0}, \mathbf{bit} = \mathbf{1}, \mathbf{recbit}, \mathbf{recack}\}$ . We interpret these on the interpreted system  $IS_b = (\mathcal{S}_\pi, \sim_S, \sim_R, h)$ , where  $\mathcal{S}_\pi$  is the set of *reachable global states* as



(0, 0)		(0-ack, 0)	
	(1, 1)		(1-ack, 1)
(0, $\epsilon$ )	(1, $\epsilon$ )		

Figure 2: The state space of the bit-transmission system. Columns (respectively rows) represent global states epistemically equivalent for  $S$  (respectively for  $R$ ). The environment component is omitted for simplicity.

defined by  $\pi$  and the protocol functions  $P_S$ ,  $P_R$  and  $P_E$ , where  $\sim_R$  and  $\sim_S$  are equivalence relations on global states as defined in Section 2, and where  $h$  is an interpretation for the atoms in  $\mathcal{P}$  such that the following holds:

$$\begin{aligned}
(IS_b, g) &\models \mathbf{bit} = \mathbf{0} && \text{if } l_S(g) = 0, \text{ or } l_S(g) = 0\text{-ack} \\
(IS_b, g) &\models \mathbf{bit} = \mathbf{1} && \text{if } l_S(g) = 1, \text{ or } l_S(g) = 1\text{-ack} \\
(IS_b, g) &\models \mathbf{recbit} && \text{if } l_R(g) = 1, \text{ or } l_R(g) = 0 \\
(IS_b, g) &\models \mathbf{recack} && \text{if } l_S(g) = 1\text{-ack}, \text{ or } l_S(g) = 0\text{-ack}.
\end{aligned}$$

This permits us to represent and check properties of the system directly on the semantical models. For example, by employing standard temporal connectives on the runs constructed above one can check that:

$$IS_b \not\models \mathbf{recbit} \rightarrow \text{AF } \mathbf{recack}$$

which represents the intrinsic unreliability of the channel (AF here means truth in some global state in all future branching runs).

Irrespective of the analysis of the dynamic properties of the system, there is one interesting static epistemic property that is worth observing. By ascribing knowledge to the agents using the standard [FHMV95] approach (see Section 2), it can be checked from Figure 2 that

$$IS_b \models \mathbf{recbit} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \quad (1)$$

which confirms our intuition about the model. Furthermore:

$$IS_b \models \mathbf{recack} \rightarrow \mathbf{recbit} \quad (2)$$

$$IS_b \models \mathbf{recack} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \quad (3)$$

and perhaps most interestingly that

$$IS_b \models \mathbf{recack} \rightarrow K_S(K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \quad (4)$$

$$IS_b \models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_S K_R(\mathbf{bit} = \mathbf{0}) \quad (5)$$

(and similarly for the case  $(\mathbf{bit} = \mathbf{1})$ ). So, if an *ack* is received by  $S$ , then  $S$  is sure that  $R$  knows the value of the bit. Intuitively this represents the fact that although the channel is potentially faulty, if messages do manage to travel back and forth the protocol is strong enough to eliminate any uncertainty in the communication. Whereas, for example,

$$IS_b \not\models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_R K_S K_R(\mathbf{bit} = \mathbf{0})$$

The properties (1–5) above can be checked by direct calculation on the table in Figure 2. Properties (3), (4) and (5) can also be derived very straightforwardly from (1) and (2) in the logic  $S5_n$ . (3) follows from (1) and (2) by propositional logic. (4) follows from (3) and the property, easy to check from Figure 2, that  $IS_b \models \mathbf{recack} \rightarrow K_S \mathbf{recack}$ : indeed, suppose  $IS_b \models \mathbf{recbit} \rightarrow \varphi$  for any formula  $\varphi$ . Then  $IS_b \models \mathbf{recack} \rightarrow \varphi$  by propositional logic from (2), and furthermore, since  $K_S$  is normal, we have both  $IS_b \models K_S(\mathbf{recack} \rightarrow \varphi)$  and  $IS_b \models K_S \mathbf{recack} \rightarrow K_S \varphi$ . With  $IS_b \models \mathbf{recack} \rightarrow K_S \mathbf{recack}$ , the property  $IS_b \models \mathbf{recack} \rightarrow K_S \varphi$  follows immediately. Property (4) is a special case. Property (5) can be derived in similar fashion. Again, the form of these derivations surely accords with our intuition about the example. The derivations will be useful, moreover, when we examine more complicated examples later.

It may be felt that the above is an unnecessarily complicated formalisation of such a simple example. Indeed, it is possible to construct much simpler formalisations that will still support the analysis conducted above. For example, instead of having four different actions for the environment, the unreliability of the communications channel can be modelled simply by non-determinism of the transition function  $\pi$ . In this simplified model, the component for the environment action is redundant (it is always ‘transmit’), and so the simplified transition function is defined as follows:

$$\begin{aligned} \pi((0, \epsilon, \cdot), (\mathit{sendbit}, \lambda)) &= \{(0, \epsilon, \cdot), (0, 0, \cdot)\}, \\ \pi((0, 0, \cdot), (\mathit{sendbit}, \mathit{sendack})) &= \{(0, 0, \cdot), (0\text{-ack}, 0, \cdot)\}, \\ \pi((0\text{-ack}, 0, \cdot), (\lambda, \mathit{sendack})) &= \{(0\text{-ack}, 0, \cdot)\} \end{aligned}$$

It is easy to check that this simplified formalisation (with the same protocol functions as before) generates exactly the same space of reachable global states as shown in Figure 2, and so produces the same analysis of the example also. The more complicated formalisation, however, with four explicit transmission actions for the environment and a correspondingly more complicated definition of the transition function  $\pi$ , will be necessary when we consider a more complicated variant of the bit-transmission problem in Section 7 where an additional controller agent is introduced to monitor the messages that pass along the communication channel. It is for this reason that we show the more complicated formalisation in this and subsequent sections.

## 5 Violation of specifications

In the previous section we have assumed that both agents and environment follow their respective protocols. We now apply the machinery of deontic interpreted systems [LS03] to analyse what happens when the specified protocols are violated in the bit-transmission problem. We examine in detail only the possibility that  $R$  is faulty. The possibility that  $S$  is faulty, and other combinations of faulty  $R$ ,  $S$  and  $E$ , can be treated in similar fashion. Specifically, we shall consider in this section the possibility that  $R$  may send acknowledgements without having received the bit. This is a simple example of an agent doing something that it should not do. In the section that follows we shall consider an example where an agent does not do something that it should do.

In order to apply the machinery we modify the framework of the previous section so that the set of local states of agent  $i$  is composed of two disjoint

sets of green ( $G_i$ ) and red ( $R_i$ ) local states, representing correct and incorrect functioning behaviour respectively. For the sender  $S$ , since we are not admitting (for the purposes of the example) the possibility of faults, its local states are all green, that is to say, allowed by the protocol. We thus have:

$$L_S = G_S = \{0, 1, 0\text{-ack}, 1\text{-ack}\}, \quad R_S = \emptyset.$$

For the case of the environment, we have admitted the possibility of faulty, or unreliable, behaviour but these ‘faults’ are not *violations* of the protocol under examination, and are not therefore regarded as ‘red’ states. Accordingly, all local states of the environment are also green, and we have:

$$L_E = G_E, \quad R_E = \emptyset.$$

It remains to model the local states of the, now potentially faulty, receiver  $R$ . We can do so by extending the set of local states  $\{0, 1, \epsilon\}$  with at least one additional element  $\epsilon\text{-ack}$  representing the (red) local state in which  $R$  has sent an acknowledgement without having received the value of the bit. This is a reasonable model if we assume that receiver  $R$  has at least some rudimentary recall or memory capability that allows its local state to record whether a *sendack* message has been sent (not necessarily successfully). In that case, it is also reasonable to extend the set of  $R$ ’s local states to include elements  $0\text{-ack}$  and  $1\text{-ack}$  as well. It remains to determine whether these two elements should be classified as red or green local states of  $R$ . We will discuss that presently. In summary, the local states of  $R$  for this version of the problem are defined as follows:

$$L'_R = \{0, 1, \epsilon, 0\text{-ack}, 1\text{-ack}, \epsilon\text{-ack}\}$$

with  $\{0, 1, \epsilon\} \subseteq G'_R$ ,  $\{\epsilon\text{-ack}\} \subseteq R'_R$ , and states  $0\text{-ack}$  and  $1\text{-ack}$  still to be classified into  $G'_R$  or  $R'_R$ .

How we classify these two remaining local states of  $R$  depends on how we choose to interpret the *ack* component. One possibility is to view the *ack* as indicating that at least one faulty acknowledgement was sent at some time before the value of the bit had been received. On this reading, the local states  $0\text{-ack}$  and  $1\text{-ack}$  would both be classified as red local states for  $R$ .

In this paper, however, we will do it differently, and in a more general way. Here, we will interpret the *ack* in local states  $0\text{-ack}$  and  $1\text{-ack}$  for  $R$  as indicating merely that a *sendack* action has been performed by  $R$ —*ack* in these states does not signify that this acknowledgement was a faulty one, just that at least one acknowledgement has been sent. We will let the protocol function  $P'_R$  for this extended system determine whether local states for  $R$  are green or red. Specifically, we will say that a local state  $l_R$  for receiver  $R$  is ‘green’ (belongs to  $G'_R$ ) iff there is a  $P'_R$ -reachable global state  $(l_S, l_R, l_E)$ , i.e. by Definition 3, iff  $(l_S, l_R, l_E)$  is the final global state of some  $P'_R$ -compliant run of the system starting at one of the initial states  $(0, \epsilon)$  or  $(1, \epsilon)$ . The green local states  $G'_R$  for  $R$  will thus represent those local states of  $R$  that could have been reached in a run in which  $R$  behaved according to its protocol at each step; this is the appropriate notion of ‘green’ for the analysis we want to conduct. The red local states  $R'_R$  are those local states of  $R$  which could only be reached by incorrect functioning behaviour by  $R$ .

We turn now to defining the protocol functions  $P'_R$ ,  $P'_S$ ,  $P'_E$ , and the transition function  $\pi'$ , of the deontic interpreted system for this version of the problem.

Since the two sets of local states for  $S$  and  $E$  have not changed we can keep the function  $P_S$  and  $P_E$  described in Section 4, i.e., we take  $P'_S = P_S$  and  $P'_E = P_E$ . But we need to extend the protocol function  $P_R$  of Section 4 so that  $P'_R$  is defined also on the local states  $L'_R$  of  $R$  that were not in  $L_R$ . We want to define what we might call a ‘conservative’ extension of the protocol function  $P_R$ : the protocol  $P'_R$  should be the same as  $P_R$  when evaluated on the local states  $L_R$ , leaving it only to define its values for the new local states  $L'_R - L_R$ . The projection of the new interpreted system onto the green local states should result in the ‘old’ system  $IS_b$  of Section 4. So we have:

$$P'_R(\epsilon) = P_R(\epsilon) = \lambda, \quad P'_R(0) = P'_R(1) = P_R(0) = P_R(1) = \text{sendack}$$

How shall we define  $P'_R$  for the new local states, and in particular for the (red) local state  $\epsilon\text{-ack}$ ? Of course, the protocol described in Section 3 *does not say*: as presented, it does not cover the possibility of violation, and does not specify what actions are to be taken if errors (here, the sending of premature acknowledgements) arise. For the sake of concreteness, let us define

$$P'_R(0\text{-ack}) = P'_R(1\text{-ack}) = \text{sendack}, \quad P'_R(\epsilon\text{-ack}) = \lambda$$

In fact, in this example, it makes little difference how we define  $P'_R$  for these new local states. As is perhaps obvious, there is no extension of the bit-transmission protocol that will recover effectively from the erroneous sending of an acknowledgement by  $R$ . The point, however, is that in general it is meaningful to define protocols on red states as well as green, whenever the protocol makes provision for remedial or error recovery actions.

The question that we would now like to ask is how the runs of the system change following the introduction of the new local states for  $R$ . Suppose the system starts as before from the global state  $g_0 = (0, \epsilon, \cdot)$  or  $g_0 = (1, \epsilon, \cdot)$ . It can either produce an error-free run or  $R$  can act faultily at any time during the evolution. Notice first that  $R$ 's faults may indeed inhibit the communication of the bit. For example, if at any point  $R$  sends an *ack* without having received the bit, this, if received by  $S$ , will make  $S$  stop sending messages, preventing any communication between the agents. This simple observation indicates that some of the properties that hold in the case of no incorrect behaviour will no longer be valid here. This should be reflected in the analysis.

Let us explore then how the analysis turns out by applying the same multi-modal language based on the set of atoms  $\mathcal{P} = \{\mathbf{bit} = \mathbf{0}, \mathbf{bit} = \mathbf{1}, \mathbf{recbit}, \mathbf{recack}\}$ , and augmented by the modal operators  $O_i, K_i, \widehat{K}_i^j$  described in Section 2. We interpret formulas on the deontic interpreted system  $IS'_b = (\mathcal{S}'_\pi, \sim'_R, \sim'_S, R'^{O'}, R'^{O'}, h')$  resulting from defining the relations  $\sim'_S, \sim'_R, R'^{O'}$  on the set of the  $\pi'$ -reachable global states  $\mathcal{S}'_\pi \subset L_S \times L'_R \times L_E$ , and with the relation  $R'^{O'}$  defined by taking the green local states for  $R$  to be those in any  $P'_R$ -reachable global state.

For the computation of reachable global states  $\mathcal{S}'_\pi$ , it remains to define the transition function  $\pi'$  for system  $IS'_b$ , by extending the definition of  $\pi$  in system  $IS_b$ . First, we specify the effects of protocol-violating actions in the states  $\mathcal{S}_\pi$  the original system  $IS_b$ . For the case where the bit is 0 (the other can be done

similarly) we shall impose:

$$\begin{aligned}
\pi'((0, \epsilon, \cdot), (\text{sendbit}, \text{sendack}, \leftrightarrow)) &= (0\text{-ack}, 0\text{-ack}, \cdot) \\
\pi'((0, \epsilon, \cdot), (\text{sendbit}, \text{sendack}, \rightarrow)) &= (0, 0\text{-ack}, \cdot) \\
\pi'((0, \epsilon, \cdot), (\text{sendbit}, \text{sendack}, \leftarrow)) &= (0\text{-ack}, \epsilon\text{-ack}, \cdot) \\
\pi'((0, \epsilon, \cdot), (\text{sendbit}, \text{sendack}, -)) &= (0, \epsilon\text{-ack}, \cdot)
\end{aligned}$$

These definitions can be expressed more concisely in various ways<sup>3</sup> but we present them in this long-winded style for ease of reading. Note that in the second and fourth cases the result state is a faulty (red) state even though the erroneous acknowledgement sent by  $R$  was not received by  $S$ .

Now we extend the definition of  $\pi$  so that  $\pi'$  is defined also on the new states corresponding to the new local states  $L'_R - L_R$ .

$$\begin{aligned}
\pi'((0, \epsilon\text{-ack}), \cdot), (\text{sendbit}, \text{sendack}, \leftrightarrow) &= (0\text{-ack}, 0\text{-ack}, \cdot) \\
\pi'((0, \epsilon\text{-ack}), \cdot), (\text{sendbit}, \text{sendack}, \rightarrow) &= (0, 0\text{-ack}, \cdot) \\
\pi'((0, \epsilon\text{-ack}), \cdot), (\text{sendbit}, \text{sendack}, \leftarrow) &= (0\text{-ack}, \epsilon\text{-ack}, \cdot) \\
\pi'((0, \epsilon\text{-ack}), \cdot), (\text{sendbit}, \text{sendack}, -) &= (0, \epsilon\text{-ack}, \cdot)
\end{aligned}$$

Similarly, for further illustration:

$$\pi'((0, 0\text{-ack}), \cdot), (\text{sendbit}, \text{sendack}, \alpha_E) = (l'_S, l'_R, \cdot)$$

where  $l'_S = 0\text{-ack}$ ,  $l'_R = 0\text{-ack}$  for  $\alpha_E = \leftrightarrow$  or  $\alpha_E = \leftarrow$ ; and  $l'_S = 0$ ,  $l'_R = 0\text{-ack}$  for  $\alpha_E = \rightarrow$  or  $\alpha_E = -$ .

The remaining cases are expressed similarly and are left to the reader. For the rest of the system we define  $\pi'$  to behave in the same way as  $\pi$  in  $IS_b$ , i.e., the projection of  $\pi'$  onto the components of  $S$  and  $E$  coincides with  $\pi$ . The transition function  $\pi'$  is depicted, in simplified form, in Figure 3.

Given these definitions we can compute the set of  $\pi'$ -reachable states as before, and then pick out the set of  $P'_R$ -reachable states (the green states for  $R$ ) by identifying the  $P'_R$ -compliant runs. The state space for this version of the problem is shown in Figure 4, again omitting the environment component for clarity.

The interpretation  $h'$  for the new system is the interpretation  $h$  from the previous section with the interpretation of **recbit** updated, thus:

$$\begin{aligned}
(IS'_b, g) \models \mathbf{recbit} \quad \text{if} \quad & l_R(g) = 1, \text{ or } l_R(g) = 0, \text{ or} \\
& l_R(g) = 0\text{-ack}, \text{ or } l_R(g) = 1\text{-ack}.
\end{aligned}$$

We can now investigate whether or not the formulas that held true in the scenario with no fault remain true here. It is easily checked that

$$IS'_b \models \mathbf{recbit} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \quad (6)$$

<sup>3</sup>In practice we formulate the definition of a transition function either as a (small) Prolog program or by means of model checking languages such as SMV, and then use those to generate the sets of reachable global states shown in Figures in this paper. Once again, a much simpler formalisation of the example can be obtained by modelling the unreliability of the communication channel by non-determinism of the transition function  $\pi'$  instead of dealing explicitly with the four different kinds of environment actions. We will need the more complicated version when we consider the variant of the example in Section 7.

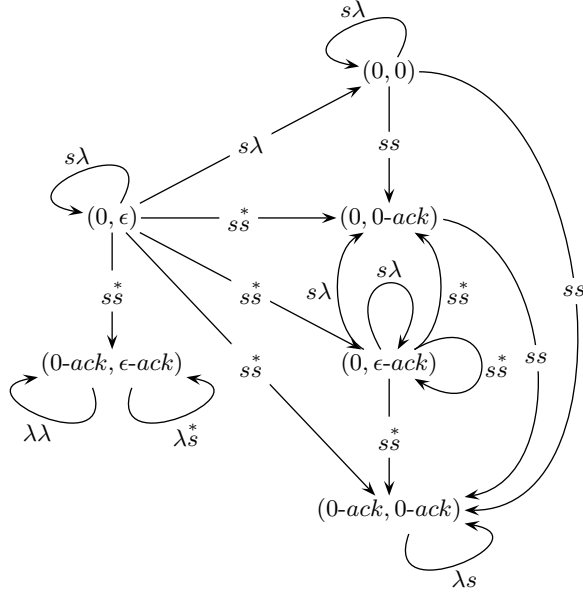


Figure 3: The transition function  $\pi'$  for the bit-transmission system in which  $R$  may send incorrect acknowledgements. The case  $\text{bit}=1$  is similar and is omitted. Labels  $s\lambda$ ,  $\lambda s$  and  $ss$  stand for the joint actions ( $\text{sendbit}, \lambda$ ),  $(\lambda, \text{sendack})$  and ( $\text{sendbit}, \text{sendack}$ ), respectively. Actions of the environment  $Act_E$  are omitted for simplicity. The  $\bar{s}$  annotation indicates actions by  $R$  that do not conform to the protocol  $P_R$ .

$(0, 0-ack)$		$(0-ack, 0-ack)$	
	$(1, 1-ack)$		$(1-ack, 1-ack)$
$(0, \epsilon-ack)$	$(1, \epsilon-ack)$	$(0-ack, \epsilon-ack)$	$(1-ack, \epsilon-ack)$
$(0, 0)$			
	$(1, 1)$		
$(0, \epsilon)$	$(1, \epsilon)$		

Figure 4: The state space of the bit-transmission system in case  $R$  may send incorrect acknowledgements. Columns (respectively rows) represent global states epistemically equivalent for  $S$  (respectively for  $R$ ). The shaded entries indicate the states that are not  $P'_R$ -reachable.

which confirms our intuition about the model. Even if faults occur, if the bit has been received,  $R$  will know its value. It is the mechanism of acknowledgements that is no longer reliable. Indeed we find in our model that:

$$\begin{aligned} IS'_b &\not\models \mathbf{recack} \rightarrow \mathbf{recbit} \\ IS'_b &\not\models \mathbf{recack} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \end{aligned}$$

and as expected:

$$\begin{aligned} IS'_b &\not\models \mathbf{recack} \rightarrow K_S(K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\ IS'_b &\not\models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_S K_R(\mathbf{bit} = \mathbf{0}) \end{aligned}$$

Notice however that using the operator  $O_R$  introduced in Section 2, which represents what holds in states where  $R$  is operating correctly, we have the following:

$$IS'_b \models O_R(\mathbf{recack} \rightarrow \mathbf{recbit}) \quad (7)$$

$$IS'_b \models O_R(\mathbf{recack} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1}))) \quad (8)$$

More interesting is that a particular form of knowledge also still holds. Intuitively if  $S$  makes the assumption of  $R$ 's correct functioning behaviour, then, upon receipt of an acknowledgement, it would make sense for it to assume that  $R$  does know the value of the bit. To model this intuition we use the operator  $\widehat{K}_i^j$  “knowledge under the assumption of correct behaviour” as presented in Section 2. This describes the knowledge agent  $i$  has if attention is restricted to states in which  $j$  is performing as intended. We refer to [LS03] for more details, but note that unlike the usual epistemic operators associated with interpreted systems,  $\widehat{K}_i^j$  is not an S5 operator, and in particular it does not have axiom T, i.e., knowledge under assumptions of correct functioning behaviour does not imply truth. This operator is of particular interest here because it captures precisely our intuition about the example. We have:

$$IS'_b \models \mathbf{recack} \rightarrow \widehat{K}_S^R(K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \quad (9)$$

$$IS'_b \models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow \widehat{K}_S^R K_R(\mathbf{bit} = \mathbf{0}) \quad (10)$$

Again, as in Section 4, we can check these properties by calculation from the table in Figure 4, and we can derive them syntactically from properties simpler to evaluate. Check first from Figure 4 that properties (6) and (7) hold. These are the key properties. From (6) and (7), and  $O_R$  normal, we easily derive (8). To derive (9) and (10) we modify the argument used in Section 4 allowing now for the presence of  $O_R$ . Suppose that the following generalised form of property (8) is valid in the model:  $IS'_b \models O_R(\mathbf{recack} \rightarrow \varphi)$ . Since  $K_S$  is normal, we get  $IS'_b \models K_S O_R(\mathbf{recack} \rightarrow \varphi)$ . The general property  $K_S O_R \varphi \rightarrow \widehat{K}_S^R \varphi$  is valid in the class of all deontic interpreted systems (see the summary of the logic in Section 2 and [LS03] for details), so we have  $IS'_b \models \widehat{K}_S^R(\mathbf{recack} \rightarrow \varphi)$ , and hence also  $IS'_b \models \widehat{K}_S^R \mathbf{recack} \rightarrow \widehat{K}_S^R \varphi$  because  $\widehat{K}_S^R$  is normal. Now  $IS'_b \models \mathbf{recack} \rightarrow K_S \mathbf{recack}$  still holds, as is easily checked from Figure 4, and the general property  $K_S \varphi \rightarrow \widehat{K}_S^R \varphi$  of the logic (see Section 2) gives  $IS'_b \models \mathbf{recack} \rightarrow \widehat{K}_S^R \mathbf{recack}$ . So we get  $IS'_b \models \mathbf{recack} \rightarrow \widehat{K}_S^R \varphi$ , of which (9) is a special case. The derivation of property (10) can be obtained in similar fashion.

To summarise: we have modified the scenario of the bit-transmission problem, relaxing slightly the assumptions of correct functioning behaviour that hold true for it. In particular, we have allowed for the possibility that one of the agents, the receiver  $R$ , may violate the protocol by performing an action it should not perform according to the protocol. We have seen that some key properties of the system then no longer hold. In particular under the assumptions we have studied,  $S$  will never know whether or not  $R$  knows the value of the bit; at best  $S$  can know this only under the assumption of correct functioning behaviour for  $R$ <sup>4</sup>. This is an intuitive result that was validated semantically on the model, and derived syntactically from some simple key properties of the model.

In specifying the extended protocol, a value of  $P'_R$  for the new local states of  $R$  had to be chosen. It seems intuitively obvious that no specification of  $P'_R$  for these states would allow us to recover the key property we require of the protocol. In other words, it should be possible to show that  $\mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_S K_R (\mathbf{bit} = \mathbf{0})$  is not valid in  $IS'_b$ , no matter how we extend the definition of  $P_R$  to cover the new local states of  $R$ —there is no protocol that will ensure recovery of communication once a violation has occurred.

How could the system recover after  $R$  has sent an incorrect *ack*? It seems as though the only way to proceed is for the two agents to re-synchronise, perhaps with a message from  $R$  signalling the failure of all communication so far with a request to restart the protocol from scratch. Clearly this request could *per se* fail to reach  $S$ , so we would need to require  $S$  to acknowledge the request to  $R$ . So the roles of  $S$  and  $R$  would need to be swapped before communication can be resumed. We see no conceptual difficulty in modelling this setting with the tools presented so far.

## 6 An error correcting protocol

Consider again the bit-transmission problem as modelled in Section 4, but assume now that  $R$  can be faulty in a different way, in that it may fail to send acknowledgements when in fact it has received the bit. To what extent would this second kind of fault compromise communication, and are there ways of recovering from it? In line with the development of the previous section, let us define the elements of the model for this new setting: sets local states, protocols, the transition function, reachable global states, and the classification of states into green (protocol-compliant) and red.

The set of local red and green states for  $S$  and  $E$  are the same as in the models of Sections 4 and 5. The local states for  $R$  are as follows:

$$L''_R = \{0, 1, \epsilon, 0-f, 1-f\}, \quad G''_R \supseteq \{0, 1, \epsilon\}$$

Here  $f$  is intended to indicate that the bit has been received but  $R$  has failed (at least once) to send an acknowledgement. It seems obvious that the red local

---

<sup>4</sup>In the case where the communication channel operates symmetrically in both directions, i.e., in the case where either the channel loses messages in both directions or both messages get through, the sender  $S$  knows that receiver  $R$  knows the value of the bit when  $S$  receives an acknowledgement, even without assuming correct behaviour by  $R$ . This is because if  $S$  receives an acknowledgement, its transmission of the bit must also have got through to  $R$ . The reader may care to re-work the analysis for the simpler case where  $Act_E = \{\leftrightarrow, -\}$  to check that the formal model confirms this informal argument.



states for  $R$  are then  $\{0-f, 1-f\}$  but for the sake of following a general method we will let the protocol function and the transition function pick these out when we compute the protocol-compliant runs.

The protocol functions we use for  $S$  and  $E$  are again the ones we employed in Sections 4 and 5; what changes is the protocol function for agent  $R$ . Again, we wish to define an extension  $P_R''$  of  $P_R$  that has the same values as  $P_R$  for the original (green) local states  $L_R$  and specifies in addition how  $R$  should behave when in one of the new (red) local states  $L_R'' - L_R$ . The transition function  $\pi''$  will specify the results of actions by  $R$  in these new states and of protocol-violating actions by  $R$  in the original (green) local states  $L_R$ . In this example, unlike the one of Section 5, there is an obvious way of extending the protocol so that we obtain error-correcting behaviour in case a fault has occurred: if  $R$  has failed to send an acknowledgement, we simply require that  $R$  does so at the next round. Formally:

$$P_R''(\epsilon) = \lambda, \quad P_R''(0) = P_R''(1) = \text{sendack}, \quad P_R''(0-f) = P_R''(1-f) = \text{sendack}.$$

So in this case, it is easy to spell out the conditions for recovery from a red state. To check that recovery is indeed obtained we evaluate formulas on the evolution of the deontic interpreted system just constructed. It remains to define the transition function of the system, so that we can compute the set of reachable global states.

As usual, we extend the definition of  $\pi$  for the original model, by considering first the conditions under which we move to one of the new local states for agent  $R$  (i.e., one of the states  $L_R'' - L_R$ ), and then the outcome of transitions originating from these local states. For all  $\alpha_E \in Act_E$ , unless stated otherwise:

$$\begin{aligned} \pi''((0, 0, \cdot), (\text{sendbit}, \lambda, \alpha_E)) &= (0, 0-f, \cdot) \\ \pi''((0-ack, 0, \cdot), (\lambda, \lambda, \alpha_E)) &= (0-ack, 0-f, \cdot) \\ \pi''((0, 0-f, \cdot), (\text{sendbit}, \lambda, \alpha_E)) &= (0, 0-f, \cdot) \\ \pi''((0, 0-f, \cdot), (\text{sendbit}, \text{sendack}, \alpha_E)) &= (0-ack, 0-f, \cdot) && (\alpha_E \in \{\leftrightarrow, \leftarrow\}) \\ \pi''((0, 0-f, \cdot), (\text{sendbit}, \text{sendack}, \alpha_E)) &= (0, 0-f, \cdot) && (\alpha_E \in \{\rightarrow, -\}) \end{aligned}$$

The definitions for bit=1 are analogous. For completeness, the remaining cases are:

$$\begin{aligned} \pi''((0-ack, 0-f, \cdot), (\lambda, \lambda, \alpha_E)) &= (0-ack, 0-f, \cdot) \\ \pi''((0-ack, 0-f, \cdot), (\lambda, \text{sendack}, \alpha_E)) &= (0-ack, 0-f, \cdot) \end{aligned}$$

With this information, we can compute the set of  $\pi''$ -reachable global states and pick out from them the  $P_R''$ -reachable ones. Once again we do not show the entire computation here, but simply summarise the results in Figure 5, with the environment component omitted as usual.

Following what is by now our standard procedure we can determine the corresponding deontic interpreted system  $IS_b''$  and interpret the formulas of interest. The interpretation function  $h''$  for atoms is as before, adjusted to cover the newly introduced global states:  $h''$  is  $h$  with the obvious extension to make atom **recbit** true also in all global states  $g$  with  $l_R(g) = 0-f$  or  $l_R(g) = 1-f$ . It

$(0, 0-f)$		$(0-ack, 0-f)$	
	$(1, 1-f)$		$(1-ack, 1-f)$
$(0, 0)$		$(0-ack, 0)$	
	$(1, 1)$		$(1-ack, 1)$
$(0, \epsilon)$	$(1, \epsilon)$		

Figure 5: The state space of the bit-transmission system in case  $R$  may fail to send acknowledgements when supposed to do so. Columns (respectively rows) represent global states epistemically equivalent for  $S$  (respectively for  $R$ ). The shaded entries indicate the states that are not  $P_R''$ -reachable.

is easily confirmed that we have all of:

$$\begin{aligned}
IS_b'' &\models \mathbf{recack} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\
IS_b'' &\models \mathbf{recack} \rightarrow K_S(K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\
IS_b'' &\models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_S K_R(\mathbf{bit} = \mathbf{0})
\end{aligned}$$

Indeed, these follow immediately by checking in Figure 5 that

$$\begin{aligned}
IS_b'' &\models \mathbf{recbit} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\
IS_b'' &\models \mathbf{recack} \rightarrow \mathbf{recbit}
\end{aligned}$$

The rest then follows by the same syntactic derivation as in Section 4.

Naturally, assuming correctness of  $R$ 's behaviour does not invalidate any of the above, so that we have also, e.g.:

$$IS_b'' \models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow \widehat{K}_S^R K_R(\mathbf{bit} = \mathbf{0})$$

Indeed,  $K_S \varphi \rightarrow \widehat{K}_S^R \varphi$  is valid in the class of all deontic interpreted systems:  $K_i \varphi \rightarrow \widehat{K}_i^j \varphi$  is a property of the logic for all pairs of agents  $i$  and  $j$  (see [LS03] for details).

This example demonstrates how the formalism of deontic interpreted systems deals with two rather intuitive points. First, not all incorrect behaviours by the participating agents necessarily compromise the validity of crucial properties of the system being modelled. Secondly, it is possible and useful to reason about these incorrect error states and devise protocols that attempt recovery from them.

## 7 Control

In the previous section it was possible to devise a simple error-correcting protocol that, if followed, ensures that the system will get back to a state in which all the agents are in a green local state. This is not always possible, of course. Given a system in which agents cannot be assumed to behave in accordance with the specified standards of behaviour or to follow the prescribed protocols, it is natural to seek additional control or enforcement mechanisms that can be introduced to encourage or even force agents to comply. Of these the simplest (but by no means the only) strategy is to look for a way of constraining the agents' behaviour so that the possibility of violation is simply eliminated. In

the present example, for instance, we could add an additional filter on the transmission of messages from  $R$  to block out those that would cause violations of the protocol. If successful, we would have an example of an enforcement strategy called ‘regimentation’ in [JS93]. Discussion of other possible strategies is beyond the scope of this paper.

So let us consider another variation of the bit-transmission problem. This time we will assume that  $R$  may develop faults of *either* of the two kinds analysed in Sections 5 and 6. We also introduce a third agent  $C$  into the system, a *controller* whose function is to constrain the behaviour of  $R$ . The controller agent  $C$  is a kind of ‘regimentation’ device. As we understand it, it is a (simple) example of what is called an ‘inter-agent’ in [RAMN<sup>+</sup>98] and a ‘controller’ in [MU00].

$C$ ’s actions are the following:

$$Act_C = \{allow, block\}.$$

These actions provide an additional filter on the attempted actions of receiver  $R$ . The *block* action will be used to override any transmit action of the environment when  $R$  attempts to send an erroneous acknowledgement; *allow* allows a message from  $R$  to enter the transmission channel to  $S$ .

The controller  $C$  must be able to detect violations of the protocol by  $R$ . However, it is a fundamental assumption of the formalism of interpreted systems that all local states are private and only actions are public. This is to model what actually happens in distributed and multi-agent systems. The local state of an agent  $i$  is invisible to a different agent  $j$ , unless there are actions by means of which agent  $i$  communicates selected elements of its local state to agent  $j$ . And likewise for the local state of the environment. Given this, the controller may block  $R$ ’s messages only on the basis of observations of the actions performed. For this example, the behaviour of  $C$  can be modelled by giving it two local states, *ok* and  $\overline{ok}$ , say. We shall not model (in this paper) the possibility that the controller  $C$  fails to act correctly, so both of its local states are green:

$$G_C = L_C = \{ok, \overline{ok}\}, \quad R_C = \emptyset.$$

We use *ok* for the state where  $C$  is allowing messages from  $R$  to enter the communication channel, and  $\overline{ok}$  for the state where  $C$  has detected an attempted violation by  $R$  and is thus blocking messages from  $R$ . The protocol of  $C$  is simply:

$$P_C'''(ok) = allow, \quad P_C'''(\overline{ok}) = block.$$

It remains to specify how the controller’s local state switches from *ok* to  $\overline{ok}$  bearing in mind that the local state of  $R$  is invisible to  $C$ . That will be done when specifying the transition function of the system.

The protocols for  $S$  and  $E$  are the same as those we have employed so far. As for  $R$ , we are now assuming it can develop both kinds of faults discussed in previous sections. So, in this version of the example its local states are the following:

$$L_R''' = \{0, 1, \epsilon, 0-ack, 1-ack, \epsilon-ack, 0-f, 1-f, 0-ack-f, 1-ack-f\}$$

where the states  $\{0, 1, \epsilon\}$  are green as usual, and where we allow the protocol and transition functions to determine which of the remaining states are green

(reachable by a protocol-compliant run) and which are red. As before, we will read *ack* as signifying that an acknowledgement (not necessarily faulty) has been sent by *R*, and *f* that at least one erroneous failure to send an acknowledgement by *R* has occurred. It is possible to devise models with much simpler representation of *R*'s red local states but we will work with this larger one as it is clearer how it combines elements from the two previous analyses.

For the protocol function for *R*, we want  $P_R'''$  to be the same as  $P_R$  when evaluated on the original (green) local states  $L_R$ , as usual, and to have the error-correcting behaviour discussed in Section 6 for the local states *0-f* and *1-f*. For these two states,  $P_R'''$  has the value of  $P_R''$ . We also want this behaviour for the local states *0-ack-f* and *1-ack-f*. For the remaining states *0-ack*, *1-ack*, *ε-ack*, the value of  $P_R'''$  can be chosen arbitrarily since again nothing of interest will depend on this. For concreteness, we will follow the choice in Section 5. So we have:

$$\begin{aligned} P_R'''(\epsilon) &= \lambda, & P_R'''(0) &= P_R'''(1) = \textit{sendack}, \\ P_R'''(0-f) &= P_R'''(1-f) = P_R'''(0-ack-f) = P_R'''(1-ack-f) = \textit{sendack}, \\ P_R'''(0-ack) &= P_R'''(1-ack) = \textit{sendack}, & P_R'''(\epsilon-ack) &= \lambda \end{aligned}$$

Note that while *C* is expected to block erroneous *sendack* messages when *R* is in the critically faulty states (those with *ack*), it cannot *prevent* *R* from entering these or any other faulty states.

We now consider the transition function of this system. In this case both global states and joint actions are 4-tuples. Let us first consider agent *C*. The controller may block *R*'s messages only on the basis of observations of the joint actions performed. In this example, we get the desired effect by stipulating that *C* starts in state  $\overline{ok}$  (blocking) and remains in state  $\overline{ok}$  until the successful transmission of a bit from *S* to *R* is observed, i.e., given our assumptions about the communication channel, until there occurs a *sendbit* action by *S* when the environment action is either  $\leftrightarrow$  or  $\rightarrow$ . When this happens, *C* switches to local state *ok*, and it remains in that state for the rest of the run. How do we know that this protocol for *C* does indeed give the desired effect? It is precisely the point of the paper to demonstrate how such claims can be expressed and verified formally.

It is clearer (and shorter) to define  $\pi'''$  in terms of another function  $\pi_+'''$  which is a combination of the effects of  $\pi'$  and  $\pi''$  defined in earlier sections, extended to cover the new local states for *R* containing both elements *ack* and *f*.  $\pi_+'''$  corresponds to the transition function of the system without the presence of the controller agent *C*. We sketch its definition presently. The definition of  $\pi'''$  is (for all  $\alpha_R \in Act_R$ ,  $\alpha_S \in Act_S$ ,  $\alpha_E \in Act_E$ ):

$$\begin{aligned} \pi'''((l_S, l_R, \overline{ok}, l_E), (\textit{sendbit}, \alpha_R, \textit{block}, -)) &= \pi_+'''((l_S, l_R, l_E), (\textit{sendbit}, \alpha_R, -)) : \overline{ok} \\ \pi'''((l_S, l_R, \overline{ok}, l_E), (\textit{sendbit}, \alpha_R, \textit{block}, \leftarrow)) &= \pi_+'''((l_S, l_R, l_E), (\textit{sendbit}, \alpha_R, -)) : \overline{ok} \\ \pi'''((l_S, l_R, \overline{ok}, l_E), (\textit{sendbit}, \alpha_R, \textit{block}, \rightarrow)) &= \pi_+'''((l_S, l_R, l_E), (\textit{sendbit}, \alpha_R, \rightarrow)) : ok \\ \pi'''((l_S, l_R, \overline{ok}, l_E), (\textit{sendbit}, \alpha_R, \textit{block}, \leftrightarrow)) &= \pi_+'''((l_S, l_R, l_E), (\textit{sendbit}, \alpha_R, \rightarrow)) : ok \\ \pi'''((l_S, l_R, ok, l_E), (\alpha_S, \alpha_R, \textit{allow}, \alpha_E)) &= \pi_+'''((l_S, l_R, l_E), (\alpha_S, \alpha_R, \alpha_E)) : ok \end{aligned}$$

where  $(l_S, l_R, l_E) : l_C$  represents the 4-tuple  $(l_S, l_R, l_C, l_E)$ .

The key things to note are (1) when the controller *C* is blocking the transmission of messages from the receiver, the effects of environment actions  $\leftrightarrow$

and  $\rightarrow$  are those of  $\leftarrow$  and  $-$ , respectively; and (2) a successful transmission of *sendbit* from the receiver (cases 3 and 4 above) switches the controller  $C$ 's local state from  $\overline{ok}$  to  $ok$ .

It remains to define the function  $\pi_+'''$ . We wish to combine the effects of  $\pi'$  and  $\pi''$  from previous sections. So, for the case where the local state  $l_R \in L_R$  and the action  $\alpha_R$  of  $R$  belongs to  $P_R(l_R)$ , the value of  $\pi_+'''$  is the value of  $\pi$ . For the case where  $l_R \in L'_R$  and  $\alpha_R \in P'_R(l_R)$  (which by construction subsumes the previous case of  $L_R$  also) the value of  $\pi_+'''$  is the value of  $\pi'$ . And for the case where  $l_R \in L''_R$  and  $\alpha_R \in P''_R(l_R)$  (which also subsumes the first case  $L_R$ ) the value of  $\pi_+'''$  is the value of  $\pi''$ . The remaining cases are those dealing with the new local states. Various concise formulations are possible but it is perhaps clearest to present the definition here in the form of a table (with corresponding entries for the case  $\text{bit}=1$ ). The  $*$  in some entries indicates that the state remains unchanged.

	$\lambda, \lambda$ $Act_E$	<i>sendbit</i> , $\lambda$ $\leftarrow -$	<i>sendbit</i> , $\lambda$ $\rightarrow \leftrightarrow$	$\lambda$ , <i>sendack</i> $\leftarrow \leftrightarrow$	$\lambda$ , <i>sendack</i> $\rightarrow -$
$(0, \epsilon)$		*	$(0, 0)$		
$(0, 0)$		$(0, 0-f)$	$(0, 0-f)$		
$(0, \epsilon-ack)$		*	$(0, 0-ack)$		
$(0, 0-ack)$		$(0, 0-ack-f)$	$(0, 0-ack-f)$		
$(0, 0-ack-f)$		*	*		
$(0-ack, \epsilon)$	*			$(0-ack, \epsilon-ack)$	$(0-ack, \epsilon-ack)$
$(0-ack, 0)$	$(0-ack, 0-f)$			$(0-ack, 0-ack)$	$(0-ack, 0-ack)$
$(0-ack, \epsilon-ack)$	*			*	*
$(0-ack, 0-ack)$	$(0-ack, 0-ack-f)$			*	*
$(0-ack, 0-ack-f)$	*			*	*

	<i>sendbit</i> , <i>sendack</i> $-$	<i>sendbit</i> , <i>sendack</i> $\rightarrow$	<i>sendbit</i> , <i>sendack</i> $\leftarrow$	<i>sendbit</i> , <i>sendack</i> $\leftrightarrow$
$(0, \epsilon)$	$(0, \epsilon-ack)$	$(0, 0-ack)$	$(0-ack, \epsilon-ack)$	$(0-ack, 0-ack)$
$(0, 0)$	$(0, 0-ack)$	$(0, 0-ack)$	$(0-ack, 0-ack)$	$(0-ack, 0-ack)$
$(0, \epsilon-ack)$	*	$(0, 0-ack)$	$(0-ack, \epsilon-ack)$	$(0-ack, 0-ack)$
$(0, 0-ack)$	*	*	$(0-ack, 0-ack)$	$(0-ack, 0-ack)$
$(0, 0-ack-f)$	*	*	$(0-ack, 0-ack-f)$	$(0-ack, 0-ack-f)$

We may now move to consider the deontic interpreted system

$$IS_b''' = (\mathcal{S}_\pi''', \sim_S''', \sim_R''', \sim_C''', R_S^{O'''}, R_R^{O'''}, R_C^O, h''')$$

generated by the system above. The set of reachable global states  $\mathcal{S}_\pi''' \subset L_S \times L_R''' \times L_C \times L_E$  can be computed from the initial states  $(0, \epsilon, \overline{ok}, \cdot)$  and  $(1, \epsilon, \overline{ok}, \cdot)$  using the transition function  $\pi'''$ , with the green ( $P_R'''$ -reachable) states picked out as usual. The resulting set of reachable global states is summarised in Figure 6. Given that the local states of  $C$  do not affect the epistemic states of  $S$  and  $R$ , which is our prime object of interest in this analysis, we leave out  $C$ 's

<del>(0, 0-ack-f)</del>		<del>(0-ack, 0-ack-f)</del>	
	(1, 1-ack-f)		(1-ack, 1-ack-f)
(0, 0-f)		(0-ack, 0-f)	
	(1, 1-f)		(1-ack, 1-f)
(0, 0-ack)		(0-ack, 0-ack)	
	(1, 1-ack)		(1-ack, 1-ack)
<del>(0, <math>\epsilon</math>-ack)</del>	<del>(1, <math>\epsilon</math>-ack)</del>	<del>(0-ack, <math>\epsilon</math>-ack)</del>	<del>(1-ack, <math>\epsilon</math>-ack)</del>
(0, 0)		(0-ack, 0)	
	(1, 1)		(1-ack, 1)
(0, $\epsilon$ )	(1, $\epsilon$ )		

Figure 6: The state space of the ‘regimented’ bit-transmission system in case  $R$  may both fail to send acknowledgements when supposed to, and send acknowledgements when supposed not to. Columns (respectively rows) represent global states epistemically equivalent for  $S$  (respectively for  $R$ ). The shaded entries indicate the states that are not  $P_R'''$ -reachable. The two entries struck out are those eliminated by introduction of the controller  $C$ .

component from the table as we have done for  $E$  earlier<sup>5</sup>. The table also shows which global states are eliminated (made unreachable from the initial states) by introduction of the regimentation mechanism  $C$ . Finally, we need to adjust the interpretation function  $h'''$  in the obvious way so that the atom **recbit** is true in all global states  $g$  except those with  $l_R(g) = \epsilon$  or  $l_R(g) = \epsilon\text{-ack}$ .

Let us now go back to the formulas we analysed before and check whether they hold true on  $IS_b'''$ . We should expect that the introduction of the controller  $C$  eliminates the possibility of incorrect acknowledgements reaching  $S$ . This is indeed what we find by analysing the formulas. It is easy to check from Figure 6 that:

$$\begin{aligned} IS_b''' &\models \mathbf{recbit} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\ IS_b''' &\models \mathbf{recack} \rightarrow \mathbf{recbit} \end{aligned}$$

from which follows, by exactly the same syntactic derivations as in Section 4:

$$\begin{aligned} IS_b''' &\models \mathbf{recack} \rightarrow (K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\ IS_b''' &\models \mathbf{recack} \rightarrow K_S(K_R(\mathbf{bit} = \mathbf{0}) \vee K_R(\mathbf{bit} = \mathbf{1})) \\ IS_b''' &\models \mathbf{recack} \wedge (\mathbf{bit} = \mathbf{0}) \rightarrow K_S K_R(\mathbf{bit} = \mathbf{0}) \end{aligned}$$

These can also be confirmed by direct calculation from the table in Figure 6.

Once again, it is possible to construct other formalisations of the example, for instance, corresponding to the cases where the receiver  $R$  has no recall/memory capability. We leave it to the reader to confirm that an analysis in the same style works also with these (simpler) formalisations.

## 8 Conclusions

We have presented three variations of the bit-transmission problem to illustrate, and evaluate, how the machinery of deontic interpreted systems provides

<sup>5</sup>More precisely, the table represents the quotient set of the set of reachable states with respect to an equivalence relation defined by  $(l_S, l_R, l_C, l_E) \sim (l'_S, l'_R, l'_C, l'_E)$  iff  $l_S = l'_S$  and  $l_R = l'_R$ .

a means of analysing violations and (certain) enforcement mechanisms. Clearly the example is trivial compared to the complex multi-agent systems applications referred to in the introductory section. Nevertheless, it does exhibit many of the features we shall have to confront in these complex examples. The fact that we have been able to carry out a detailed formal analysis of different kinds of faults, and of the effects of enforcement and control mechanisms, encourages us to believe that useful tools and methods can be developed on this basis.

Apart from examining further examples, we are pursuing three main lines of development. First, we have identified a number of technical questions concerning what we called ‘conservative’ extensions of protocol functions to include red states. Generally, we can find ways of structuring the definition of protocol functions and system transition functions to make them easier to construct and maintain as the applications become more complex. Space limitations prevented us from discussing alternative styles of definitions in this paper. Second, we are investigating an extended formalism which colours *transitions* red or green and not just states, combining the formalism of interpreted systems with the construction of a dynamic logic of permission reported in [Mey96]. Third, and perhaps most important, is the question of how these methods will scale up to deal with realistic examples with many agents and many kinds of faults.

It is fortunate that in the bit transmission problem it is possible to derive the properties of interest syntactically from very simple properties that are easy to check semantically on the model. This will not always be the case. We do believe, however, that computational support can be provided to enable the analysis of large, realistic problems. We have been experimenting with model checking software for this purpose. Specifically, we use the NuSMV temporal model checker [CCGR99] to compute the set of runs of the system. From these we extract the set of reachable global states, and feed them into AKKA<sup>6</sup>, a software system for testing the validity of multi-modal formulas in a Kripke model, to verify the epistemic properties of interest. Preliminary experiments using this method are reported in [LRS02]. We are also experimenting with the use of action description formalisms of the kind found in AI to make the definition of protocol and system transition functions more concise.

## Acknowledgements

This work was supported by the EU-funded FET project ALFEBIITE (IST-1999-10298).

## References

- [CCGR99] A. Cimatti, E. Clarke, F. Giunchiglia, and M. Roveri. NuSMV: A new symbolic model verifier. In *Proc. 11th International Computer Aided Verification Conference*, pages 495–499, 1999.
- [CJ96] J. Carmo and A. J. I. Jones. Deontic database constraints, violation and recovery. *Studia Logica*, 57(1):139–165, 1996.
- [Coe93] J. Coenen. Top-down development of layered fault-tolerant systems and its problems—A deontic perspective. *Annals of Mathematics and Artificial Intelligence*, 9(1–2):133–150, 1993.

---

<sup>6</sup><http://turing.wins.uva.nl/~lhendrik/>

- [FHMV95] R. Fagin, J. Y. Halpern, Y. Moses, and M. Y. Vardi. *Reasoning about Knowledge*. MIT Press, Cambridge, 1995.
- [FHMV97] R. Fagin, J. Y. Halpern, Y. Moses, and M. Y. Vardi. Knowledge-based programs. *Distributed Computing*, 10(4):199–225, 1997.
- [FM91] J. Fiadeiro and Maibaum. Temporal reasoning over deontic specifications. *Journal of Logic and Computation*, 1(3):357–395, 1991.
- [GHR93] D. M. Gabbay, I. M. Hodkinson, and M. A. Reynolds. *Temporal Logic: Mathematical Foundations and Computational Aspects, Volume 1: Mathematical Foundations*. Oxford University Press, 1993.
- [GMP92] J. Glasgow, G. MacEwen, and P. Panangaden. A logic for reasoning about security. *ACM Transactions on Computer Systems*, 10(3):226–264, August 1992.
- [Hin62] J. Hintikka. *Knowledge and Belief, An Introduction to the Logic of the Two Notions*. Cornell University Press, Ithaca (NY) and London, 1962.
- [HZ92] J. Y. Halpern and L. D. Zuck. A little knowledge goes a long way: Knowledge-based derivations and correctness proofs for a family of protocols. *Journal of the ACM*, 39(3):449–478, 1992.
- [JS93] A. J. I. Jones and M. J. Sergot. On the Characterisation of Law and Computer Systems: The Normative Systems Perspective. In John-Jules Ch. Meyer and Roel J. Wieringa, editors, *Deontic Logic in Computer Science: Normative System Specification*, chapter 12, pages 275–307. John Wiley & Sons, Chichester, England, 1993.
- [KM89] S. Khosla and T. Maibaum. The prescription and description of state-based systems. In B. Banieqbal, H. Barringer, and A. Pnueli, editors, *Temporal Logic in Specification*, number 398 in LNCS, pages 243–294. Springer-Verlag, 1989.
- [LRS02] A. Lomuscio, F. Raimondi, and M. Sergot. Towards model checking interpreted systems. In *Proceedings of Mochart — First International Workshop on Model Checking and Artificial Intelligence*, 2002.
- [LS03] A. Lomuscio and M. Sergot. Deontic interpreted systems. *Studia Logica*, 75, 2003.
- [Mey88] J.-J. Ch. Meyer. A different approach to deontic logic: Deontic logic viewed as a variant of dynamic logic. *Notre Dame Journal of Formal Logic*, 29(1):109–136, 1988.
- [Mey96] R. van der Meyden. The dynamic logic of permission. *Journal of Logic and Computation*, 6(3):465–479, 1996.
- [MH95] J.-J. Ch. Meyer and W. van der Hoek. *Epistemic Logic for AI and Computer Science*, volume 41 of *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press, 1995.
- [ML85] N. H. Minsky and A. D. Lockman. Ensuring integrity by adding obligations to privileges. In *Proc. 8th International Conference on Software Engineering*, pages 92–102. IEEE Computer Society Press, 1985.
- [MU00] N. Minsky and V. Ungureanu. Law-governed interaction: A coordination and control mechanism for heterogeneous distributed systems. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 9(3):273–305, 2000.
- [RAMN<sup>+</sup>98] J. A. Rodriguez-Aguilar, F. J. Martin, P. Noriega, P. Garcia, and C. Sierra. Towards a test-bed for trading agents in electronic auction markets. *AI Communications*, 11:5–19, 1998.
- [Woo00] M. Wooldridge. *Reasoning about Rational Agents*. MIT Press, July 2000.