

Verifying Compliance with Commitment Protocols

Enabling Open Web-Based Multiagent Systems

Mahadevan Venkatraman and Munindar P. Singh

Department of Computer Science, North Carolina State University, Raleigh, NC 27695, USA

Abstract. Interaction protocols are specific, often standard, constraints on the behaviors of autonomous agents in a multiagent system. Protocols are essential to the functioning of open systems, such as those that arise in most interesting web applications. A variety of common protocols in negotiation and electronic commerce are best treated as *commitment protocols*, which are defined, or at least analyzed, in terms of the creation, satisfaction, or manipulation of the commitments among the participating agents.

When protocols are employed in open environments, such as the Internet, they must be executed by agents that behave more or less autonomously and whose internal designs are not known. In such settings, therefore, there is a risk that the participating agents may fail to comply with the given protocol. Without a rigorous means to verify compliance, the very idea of protocols for interoperation is subverted. We develop an approach for testing whether the behavior of an agent complies with a commitment protocol. Our approach requires the specification of commitment protocols in temporal logic, and involves a novel way of synthesizing and applying ideas from distributed computing and logics of program.

Key words: Commitments; Protocols; Causality; Temporal logic; Formal methods

1. Introduction

Interaction among agents is the distinguishing property of multiagent systems. However, ensuring that only the desirable interactions occur is one of the most challenging aspects of multiagent system analysis and design. This is especially so when the given multiagent system is meant to be used as an open system, for example, in web-based applications.

Because of its ubiquity and ease of use, the web is rapidly becoming the platform of choice for a number of important applications, such as trading, supply-chain management, and in general electronic commerce. However, the web can enforce few constraints on the agents' behavior. Current approaches to security on the web emphasize how the different parties to a transaction may be authenticated or how their data may be encrypted to prevent unauthorized access. Even with authentication and controlled access, the parties would have support beyond conventional protocol techniques (such as finite state machine models) neither to specify the desired interactions nor to detect any violation. However, authentication and access control are conceptually orthogonal to ensuring that the parties behave and interact correctly. Even when the parties are authenticated, they may act undesirably through error or

malice. Conversely, the parties involved may resist going through authentication, but may be willing to be governed by the applicable constraints.

The web provides an excellent infrastructure through which agents can communicate with one another. But the above problems are exacerbated when agents are employed in the web. In contrast with traditional programs and interfaces, neither their behaviors and interactions nor their construction is fixed or under the control of a single authority. In general, in an open system, the member agents are contributed by several sources and serve different interests. Thus, these agents must be treated as

- *autonomous*—with few constraints on behavior, reflecting the independence of their users, and
- *heterogeneous*—with few constraints on construction, reflecting the independence of their designers.

Effectively, the multiagent system is specified as a kind of standard that its member agents must respect. In other words, the multiagent system can be thought of as specifying a protocol that governs how its member agents must act. For our purposes, the standard may be *de jure* as created by a standards body, or *de facto* as may emerge from practice or even because of the arbitrary decisions of a major vendor or user organization. All that matters for us is that a standard imposes some restrictions on the agents. Consider the fish-market protocol as an example of such a standard protocol [14].

Example 1. In the fish-market protocol, we are given agents of two roles: a single auctioneer and one or more potential bidders. The fish-market protocol is designed to sell fish. The seller or auctioneer announces the availability of a bucket of fish at a certain price. The bidders gathered around the auctioneer can scream back *Yes* if they are interested and *No* if they are not; they may also stay quiet, which is interpreted as a lack of interest or *No*. If exactly one bidder says *Yes*, the auctioneer will sell him the fish; if no one says *Yes*, the auctioneer lowers the price; if more than one bidder says *Yes*, the auctioneer raises the price. In either case, if the price changes, the auctioneer announces the revised price and the process iterates. ■

Because of its relationship to protocols in electronic commerce and because it is more general than the popular English and Dutch auctions, the fish-market protocol has become an important one in the recent multiagent systems literature. Accordingly, we use it as our main example in this paper.

Because of the autonomy and heterogeneity requirements of open systems, compliance testing can be based neither on the internal designs of the agents nor on concepts such as beliefs, desires, and intentions that map to internal representations [16]. The only way in which compliance can be tested

is based on the behavior of the participating agents. The testing may be performed by a central authority or by any of the participating agents. However, the requirements for behavior in multiagent systems can be quite subtle. Thus, along with languages for specifying such requirements, we need corresponding techniques to test compliance.

1.1. COMMITMENTS IN AN OPEN ARCHITECTURE

There are three levels of architectural concern in a multiagent system. One deals with individual agents; another deals with the systemic aspects of how different services and brokers are arranged. Both of these have received much attention in the literature. In the middle is the multiagent *execution* architecture, which has not been as intensively studied within the community. An execution architecture must ultimately be based on distributed computing ideas albeit with an open flavor, e.g., [1, 5, 11]. A well-defined execution functionality can be given a principled design, and thus facilitate the construction of robust and reusable systems. Some recent work within multiagent systems, e.g., Ciancarini *et al.* [8, 9] and Singh [18], has begun to address this level.

Much of the work on this broad theme, however, focuses primarily on coordination, which we think of as the lowest level of interaction. Coordination deals with how autonomous agents may align their activities in terms of what they do and when they do it. However, there is more to interaction in general, and compliance in particular. Specifically, interaction must include some consideration of the commitments that the agents enter into with each other. The commitments of the agents are not only base-level commitments dealing with what actions they must or must not perform, but also metacommitments dealing with how they will adjust their base-level commitments [20]. Commitments provide a layer of coherence to the agents' interactions with each other. They are especially important in environments where we need to model any kind of contractual relationships among the agents.

Such environments are crucial wherever open multiagent systems must be composed on the fly, e.g., in electronic commerce of various kinds on the Internet. The addition of commitments as an explicit first-class object results in considerable flexibility of how the protocols can be realized in changing situations. We term such augmented protocols *commitment protocols*.

Example 2. We informally describe the protocol of Example 1 in terms of commitments. When a bidder says *Yes*, he commits to buying the bucket of fish at the advertised price. When the auctioneer advertises a price, he commits that he will sell fish at that price if he gets a unique *Yes*. Neither commitment is irrevocable. For example, if the fish are spoiled, the auctioneer releases the bidder from paying for them. Specifying all possibilities in terms of irrevocable commitments would complicate each commitment, but would still fail to capture the practical meanings of such a protocol. For instance,

the auctioneer may not honor his offering price if a sudden change in weather indicates that fishing will be harder for the next few days. ■

1.2. COMPLIANCE IN OPEN SYSTEMS

The existence of standardized protocols is necessary but not sufficient for the correct functioning of open multiagent systems. We must also ensure that the agents behave according to the protocols. This is the challenge of *compliance*. However, unlike in traditional closed systems, verifying compliance in open systems is practically and even conceptually nontrivial.

Preserving the autonomy and heterogeneity of agents is crucial in an open environment. Otherwise, many applications would become infeasible. Consequently, protocols must be specified as flexibly as possible without making untoward requirements on the participating agents. Similarly, an approach for testing compliance must not require that the agents are homogeneous or impose stringent demands on how they are constructed.

Consequently, in open systems, compliance can be meaningfully expressed only in terms of observable behavior. This leads to two subtle considerations. One, although we talk in terms of behavior, we must still consider the high-level abstractions that differentiate agents from other active objects. The focus on behavior renders approaches based on mental concepts ineffective [16]. However, well-framed social constructs can be used. Two, we must clearly delineate the role of the observer who assesses compliance.

1.3. CONTRIBUTIONS

The approach developed here treats multiagent systems as distributed systems. There is an underlying messaging layer, which delivers messages asynchronously and, for now, reliably. However, the approach assumes for simplicity that the agents are not malicious and do not forge the timestamps on the messages that they send or receive.

The compliance testing is performed by any observer of the system—typically, a participating agent. Our approach is to evaluate temporal logic specifications with respect to locally constructed models for the given observer. The model construction proposed here employs a combination of the notion of potential causality and operations on social commitments (both described below). Our contributions are in

- incorporating potential causality in the construction of local models
- identifying patterns of messages corresponding to different operations on commitments
- showing how to verify compliance based on local information.

Our approach also has important ramifications on agent communication in general, which we discuss in Section 4.

Organization. The rest of this paper is organized as follows. Section 2 presents our technical framework, which combines commitments, potential causality, and temporal logic. Section 3 presents our approach for testing (non-)compliance of agents with respect to a commitment protocol. Section 4 concludes with a discussion of our major themes, the literature, and the important issues that remain outstanding.

2. Technical Framework

Commitment protocols as defined here are a multiagent concept. They are far more flexible and general than commitment protocols in distributed computing and databases, such as *two-phase commit* [12, pp. 562–573]. This is because our underlying notion of commitment is flexible, whereas traditional commitments are rigid and irrevocable. However, because multiagent systems are distributed systems and commitment protocols are protocols, it is natural that techniques developed in classical computer science will apply here. Accordingly, our technical framework integrates approaches from distributed computing, logics of program, and distributed artificial intelligence.

2.1. POTENTIAL CAUSALITY

The key idea behind potential causality is that the ordering of events in a distributed system can be determined only with respect to an observer [13]. If event e precedes event f with respect to an observer, then e may *potentially* cause f . The observed precedence suggests the possibility of an information flow from e to f , but without additional knowledge of the internals of the agents, we cannot be sure that true causation was involved. It is customary to define the local time of an agent as the number of steps it has executed. A *vector clock* is a vector, each of whose elements corresponds to the local time of each communicating agent. A vector v is considered later than a vector u if v is later on some, and not sooner on any, element.

Definition 1. A clock over n agents is an n -ary vector $v = \langle v_1 \dots v_n \rangle$ of natural numbers. The starting clock is $\vec{0} \triangleq \langle 0 \dots 0 \rangle$. ■

Notice that the vector representation is just a convenience. We could just as well use pairs of the form $\langle \text{agent-id}, \text{local-time} \rangle$, which would allow us to model systems of varying membership more easily.

Definition 2. Given n -ary vectors u and v , $u \prec v$ if and only if $(\forall i : 1 \leq i \leq n : u_i \leq v_i)$ and $(\exists i : 1 \leq i \leq n : u_i < v_i)$. ■

Each agent starts at $\vec{0}$. It increments its entry in that vector whenever it performs a local event [15]. It attaches the entire vector as a timestamp to any message it sends out. When an agent receives a message, it updates its vector clock to be the element-wise maximum of its previous vector and the vector timestamp of the message it received. Intuitively, the message brings news of how far the system has progressed; for some agents, the recipient may have better news already. However, any message it sends after this receive event will have a later timestamp than the message just received.

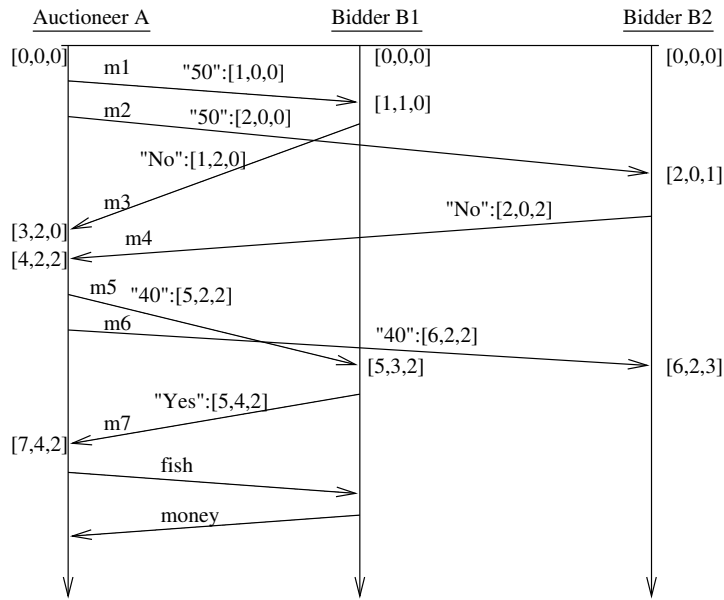


Figure 1. Vector clocks in the fish-market protocol.

Example 3. Figure 1 illustrates the evolution of vector timestamps for one possible run of the fish-market protocol. In the run described here, the auctioneer (A) announces a price of 50 for a certain bucket of fish. Bidders B1 and B2 both decline. A lowers the price to 40 and announces it. This time B1 says *Yes*, leading A to transfer the fish to B1 and B1 to send money to A. For uniformity, the last two steps are also modeled as communications. The messages are labeled m_i to facilitate reference from the text. ■

2.2. TEMPORAL LOGIC

The progression of events, which is inherent in the execution of any protocol, suggests the need for representing and reasoning about time. Temporal logics

provide a well-understood means of doing so, and have been applied in various subareas of computer science. Because of their naturalness in expressing properties of systems that may evolve in more than one possible way and for the efficiency of reasoning that they support, the branching-time logics have been especially popular in this regard [10]. Of these, the best known is Computation Tree Logic (CTL), which we adapt here in our formal language \mathcal{L} . Conventionally, a model of CTL is expressed as a tree. Each node in the tree is associated with a state of the system being considered; the branches of the tree or *paths* thus indicate the possible courses of events or ways in which the system's state may evolve. CTL provides a natural means by which to specify acceptable behaviors of the system.

The following Backus-Naur Form (BNF) grammar with a distinguished start symbol L gives the syntax of \mathcal{L} . \mathcal{L} is based on a set Φ of atomic propositions. Below, *slant* typeface indicates nonterminals; \longrightarrow and $|$ are meta-symbols of BNF specification; \ll and \gg delimit comments; the remaining symbols are terminals. As is customary in formal semantics, we are only concerned with abstract syntax.

L1. $L \longrightarrow Prop \ll\text{atomic propositions: members of } \Phi \gg$

L2. $L \longrightarrow \neg L \ll\text{negation}\gg$

L3. $L \longrightarrow L \wedge L \ll\text{conjunction}\gg$

L4. $L \longrightarrow A P \ll\text{universal quantification over paths}\gg$

L5. $L \longrightarrow E P \ll\text{existential quantification over paths}\gg$

L6. $P \longrightarrow L U L \ll\text{until: operator over a single path}\gg$

The meanings of formulas generated from L are given relative to a model and a state in the model. The meanings of formulas generated from P are given relative to a path and a state on the path. The boolean operators are standard. Useful abbreviations include $\text{false} \equiv (p \wedge \neg p)$, for any $p \in \Phi$, $\text{true} \equiv \neg \text{false}$, $p \vee q \equiv \neg p \wedge \neg q$ and $p \rightarrow q \equiv \neg p \vee q$. The temporal operators A and E are quantifiers over paths. Informally, pUq means that on a given path from the given state, q will eventually hold and p will hold until q holds. Fq means “eventually q ” and abbreviates $\text{true}Uq$. Gq means “always q ” and abbreviates $\neg F\neg q$. Therefore, $EpUq$ means that on some future path from the given state, q will eventually hold and p will hold until q holds.

Definition 3. $M = \langle \mathbf{S}, <, \mathbf{I} \rangle$ is a formal model for \mathcal{L} . \mathbf{S} is a set of states; $< \subseteq S \times S$ is a partial order indicating branching time, and $\mathbf{I} : \mathbf{S} \mapsto \mathcal{P}(\Phi)$ is an interpretation, which tells us which atomic propositions are true in a given state. For $t \in \mathbf{S}$, \mathbf{P}_t is the set of paths emanating from t . ■

$M \models_t p$ expresses “ M satisfies p at t ” and $M \models_{P,t} p$ expresses “ M satisfies p at t along path P .”

M1. $M \models_t \psi$ iff $\psi \in \mathbf{I}(t)$, where $\psi \in \Phi$

M2. $M \models_t p \wedge q$ iff $M \models_t p$ and $M \models_t q$

M3. $M \models_t \neg p$ iff $M \not\models_t p$

M4. $M \models_t Ap$ iff $(\forall P : P \in \mathbf{P}_t \Rightarrow M \models_{P,t} p)$

M5. $M \models_t Ep$ iff $(\exists P : P \in \mathbf{P}_t \text{ and } M \models_{P,t} p)$

M6. $M \models_{P,t} pUq$ iff $(\exists t' : t \leq t' \text{ and } M \models_{P,t'} q \text{ and } (\forall t'' : t \leq t'' \leq t' \Rightarrow M \models_{P,t''} p))$

The above is an abstract semantics. In Section 3.3, we specify the concrete form of Φ , \mathbf{S} , $<$, and \mathbf{I} , so the semantics can be exercised in our computations.

3. Approach

In their generic forms, both causality and temporal logic are well-known. However, applying them in combination and in the particular manner suggested here is novel to this paper.

Temporal logic model checking is usually applied for design-time reasoning [10, pp. 1042–1046]. We are given a specification and an implementation, i.e., program, that is supposed to meet it. A model is generated from the program. A model checking algorithm determines whether the specification is true in the generated model. However, in an open, heterogeneous environment, a design may not be available at all. For example, the vendors who supply the agents may consider their designs to be trade secrets.

By contrast, ours is a run-time approach, and can meaningfully apply model checking even in open settings. This is because it uses a model generated from the joint executions of the agents involved. Model checking in this setting simply determines whether the present execution satisfies the specification. If an execution respects the given protocol, that does not entail that all executions will, because an agent act inappropriately in other circumstances. However, if an execution is inappropriate, that does entail that the system does not satisfy the protocol. Consequently, although we are verifying specific executions of the multiagent system, we can only falsify (but not verify) the correctness of the construction of the agents in the system.

Model checking of the form introduced above may be applied by any observer in the multiagent system. A useful case is when the observer is one of the participating agents. Another useful case is when the observer is some

agent dedicated to the task of managing or auditing the interactions of some of the agents in the multiagent system.

Potential causality is most often applied in distributed systems to ensure that the messages being sent in a system satisfy causal ordering [3]. Causality motivates vector clocks and vector timestamps on messages, which help ensure correct ordering by having the messaging subsystem reorder and retransmit messages as needed. This application of causality can be important, but is controversial [4, 6], because its overhead may not always be justifiable.

In our approach, the delivery of messages may be noncausal. However, causality serves the important purpose of yielding accurate models of the observations of each agent. These are needed, because in a distributed system, the global model is not appropriate. Creating a monolithic model of the execution of the entire system requires imposing a central authority through which all messages are routed. Adding such an authority would take away many of the advantages that make distributed systems attractive in the first place. Consequently, our method of constructing and reasoning with models should

- not require a centralized message router
- work from a single vantage of observation, but be able to handle situations where some agents pool their evidence.

Such a method turns out to naturally employ the notion of potential causality.

3.1. MODELS FROM OBSERVATIONS

The observations made by each agent are essentially a record of the messages it has sent or received. Since each message is given a vector timestamp, the observations can be partially ordered. In general, this order is not total, because messages received from different agents may be mutually unordered.

Example 4. Figure 2 shows the models constructed locally from the observations of the auctioneer and a bidder in the run of Example 3. ■

Although a straightforward application of causality, the above example shows how local models may be constructed. Some subtleties are discussed next.

As remarked above, commitments give the core meaning of a protocol. Our approach builds on a flexible and powerful variety of social commitments, which are the commitments of one agent to another [20]. These commitments are defined relative to a *context*, which is typically the multiagent system itself. The *debtor* refers to the agent that makes a commitment, and the *creditor* to the agent who receives the commitment. Thus we have the following logical form.

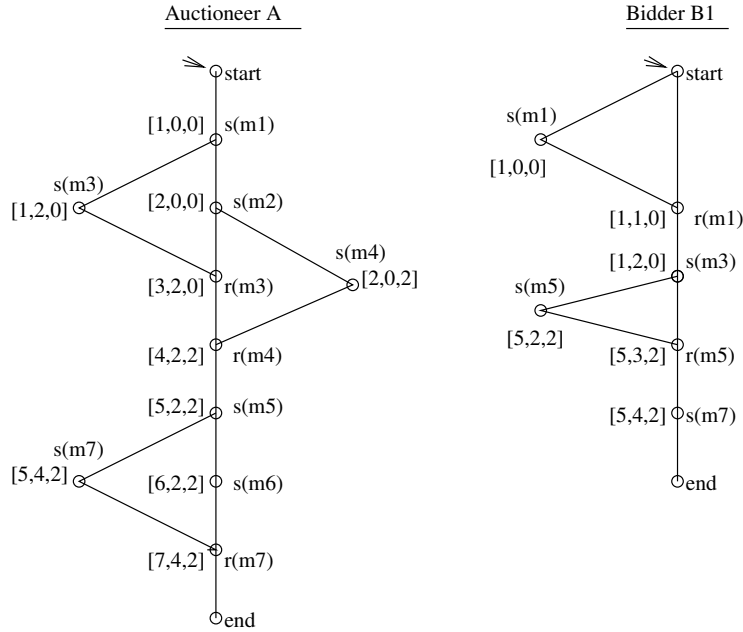


Figure 2. Observations for auctioneer and a bidder in the fish-market protocol.

Definition 4. A commitment is an expression $C(x, y, G, p)$, where x is the debtor, y the creditor, G the context, and p the condition committed to. ■

The expression c is considered true in states where the corresponding commitment exists.

Definition 5. A commitment $c = C(x, y, G, p)$ is *base-level* if p does not refer to any other commitments; c is a *metacommitment* if p refers to a base-level commitment (we do not consider higher-order commitments here). ■

Intuitively, a protocol definition is a set of metacommitments for the different roles (along with a mapping of the message tokens to operations on commitments). In combination with what the agents communicate, these lead to base-level commitments being created or manipulated, which is primarily how a commitment may be referred to within a protocol. The violation of a base-level commitment can give us proof or the “smoking gun” that an agent is noncompliant.

The following *operations* on commitments define how they may be created or manipulated. When we view commitments as an abstract data type, the operations are methods of that data type.

Each operation is realized through a simple message pattern, which states what messages must be communicated among which of the participants and in what order. For the operations on commitments we consider, the patterns

are simple. As described below, most patterns require only a single message, but some require three messages. Obeying the specified patterns ensures that the local models have the information necessary for testing compliance. That the given operation can be performed at all depends on whether the protocol, through its metacommitments, allows that operation. However, when an operation is allowed, it affects the agents' commitments. For simplicity, we assume that the operations on commitments are given a deterministic interpretation. Here z is an agent and $c = C(x, y, G, p)$ is a commitment.

- O1. *Create*(x, c) instantiates a commitment c . *Create* is typically performed as a consequence of the commitment's debtor promising something contractually or by the creditor exercising a metacommitment previously made by the debtor. *Create* usually requires a message from the debtor to the creditor.
- O2. *Discharge*(x, c) satisfies the commitment c . It is performed by the debtor concurrently with the actions that lead to the given condition being satisfied, e.g., the delivery of promised goods or funds. For simplicity, we treat the *discharge* actions as performed only when the proposition p is true. Thus the *discharge* actions are *detached*, meaning that p can be treated as true in the given moment. We model the *discharge* as a single message from the debtor to the creditor.
- O3. *Cancel*(x, c) revokes the commitment c . It can be performed by the debtor as a single message. At the end of this action, $\neg c$ usually holds. However, depending on the existing metacommitments, the *cancel* of one commitment may lead to the *create* of other commitments.
- O4. *Release*(G, c) or *release*(y, c) essentially eliminates the commitment c . This is distinguished from both *discharge* and *cancel*, because *release* does not mean success or failure, although it lets the debtor off the hook. At the end of this action, $\neg c$ usually holds. The *release* action may be performed by the context or the creditor of the given commitment, also as a single message. Because *release* is not performed by the debtor, different metacommitments apply than for *cancel*.
- O5. *Delegate*(x, z, c) shifts the role of debtor to another agent within the same context, and can be performed by the (old) debtor (or the context). Let $c' = C(z, y, G, p)$. At the end of the *delegate* action, $c' \wedge \neg c$ holds.

To prevent the risk of miscommunication, we require the creditor to also be involved in the message pattern. Figure 3(1) shows the associated pattern. The first message sets up the commitment c from x to y and is not part of the pattern. When x delegates the commitment c to z , x tells both y and z that the commitment is delegated. z is now committed to y . Later

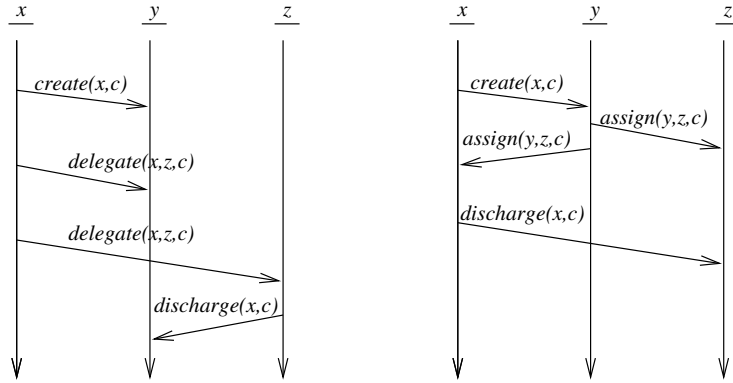


Figure 3. Message pattern for delegate (l) and assign (r).

z may discharge the commitment. The two delegate messages constitute the pattern.

06. $Assign(y, z, c)$ transfers a commitment to another creditor within the same context, and can be performed by the present creditor or the context. Let $c' = C(x, z, G, p)$. At the end of the $assign$ action, $c' \wedge \neg c$ holds.

Here we require that the new creditor and the debtor are also involved as shown in Figure 3(r). The figure shows only the general pattern. Here x is committed to y . When y assigns the commitment to z , y tells both x and z (so z knows it is the new creditor). Eventually, x should discharge the commitment to z . A potentially tricky situation is if x discharges the commitment c even as y is assigning c to z (i.e., the messages cross). In this case, we require y to discharge the commitment to z —essentially by forwarding the contents of the message from x . Thus the worst case requires three messages.

We write the operations as propositions indicating successful execution. Based on the applicable metacommitments, each operation may entail additional operations that take place implicitly.

Definition 6. A commitment c is *resolved* through a *release*, *discharge*, *cancel*, *delegate*, or *assign* performed on c . c ceases to exist when resolved. However, a new commitment is created for *delegate* or *assign*. ■

(New commitments created because of some existing metacommitment are not included in the definition of resolution. Theorem 1 states that the creditor knows the disposition of any commitments due to it. This result helps establish that the creditor can always determine compliance of others relative to what was committed to it.

Theorem 1. If message m_i creates commitment c and message m_j resolves c , then the creditor of c sees both m_i and m_j .

Proof. By inspection of the message patterns constructed for the various operations on commitments. ■

Definition 7. A commitment c is *ultimately resolved* through a *release*, *discharge*, or *cancel* performed on c , or through the ultimate resolution of any commitments created by the *delegate* or *assign* of c . ■

Theorem 2 essentially states that the creation and ultimate resolution of a commitment occur along the same causal path. This is important, because it legitimizes a significant optimization below. Indeed, we defined the above message patterns so we would obtain Theorem 2.

Theorem 2. If message m_i creates commitment c and message m_j ultimately resolves c , then $m_i \prec m_j$.

Proof. By inspection of the message patterns constructed for the various operations on commitments. ■

3.2. SPECIFYING PROTOCOLS

We first consider the coordination and then the commitment aspects of compliance. A *skeleton* is a coarse description of how an agent may behave [18]. A skeleton is associated with each role in the given multiagent system to specify how an agent playing that role may behave in order to coordinate with others. Coordination includes the simpler aspects of interaction, e.g., turn-taking. Coordination is required so that the agents' commitments make sense. For instance, a bidder should not make a bid prior to the advertisement; otherwise, the commitment content of the bid would not even be fully defined.

The skeletons may be constructed by introspection or through the use of a suitable methodology [19]. No matter how they are created, the skeletons are the first line of compliance testing, because an agent that does not comply with the skeleton for its role is automatically in violation. So as to concentrate on commitments in this paper, we postulate that a "proxy" object is interposed between an agent and the rest of the system and ensures that the agent follows the dictates of the skeleton of its role.

We now define the syntax of the specification language through the following grammar whose start symbol is *Protocol*. The braces { and } indicate that the enclosed item is repeated 0 or more times.

L7. $Protocol \longrightarrow \{Meta\} \{Message\}$

- L8. *Message* \rightarrow *Token: Commitment* \ll messages correspond to commitments \gg
- L9. *Meta* \rightarrow $C(\text{Debtor}, \text{Creditor}, \text{Context}, \text{MetaProp})$
- L10. *MetaProp* \rightarrow $AG[\text{Bool} \rightarrow \text{AFAct}] \mid AG[\text{Act} \rightarrow \text{Bool}]$
- L11. *Bool* \rightarrow \ll Boolean combinations of \gg *Act* \mid *Commitment* \mid *Dom*
- L12. *Act* \rightarrow $\text{Operation}(\text{Agent}, \text{Commitment})$
- L13. *Operation* \rightarrow \ll the six operations of Section 3.1 \gg
- L14. *Commitment* \rightarrow $\text{Meta} \mid C(\text{Debtor}, \text{Creditor}, \text{Context}, \text{AFDom})$
- L15. *Dom* \rightarrow \ll domain-specific concepts \gg

The above language embeds a subset of \mathcal{L} . Our approach is to detach the outer actions and commitments, so we can process the inner \mathcal{L} part as a temporal logic. By using commitments and actions on them, instead of simple domain propositions, we can capture a variety of subtle situations, e.g., to distinguish between *release* and *cancel* both of which result in the given commitment being removed.

Example 5 applies the above language on the fish-market protocol.

Example 5. The messages in Figure 1 can be given a content based on the following definitions. Here FM is the fish-market context.

- *fish*: a domain proposition meaning the fish is delivered
- *money_i*: a domain proposition meaning that the appropriate money is paid (subscripted to allow different prices)
- *Bid_i(B_j)*: an abbreviation for $C(B_j, A, FM, AG[\text{fish} \rightarrow \text{AFcreate}(B_j, C(B_j, A, FM, \text{AFmoney}_i))])$ —meaning the bidder promises to pay *money_i* if given the fish
- *Ad_i(B_j)*: an abbreviation for $C(A, B_j, FM, AG[\text{Bid}_i(B_j) \rightarrow \text{AFcreate}(A, C(A, B_j, FM, \text{AFfish}))])$ —meaning the auctioneer offers to deliver the fish if he gets a bid for *money_i*
- *demand_i*: an abbreviation for $(\exists j, k : j \neq k \wedge \text{Bid}_i(B_j) \wedge \text{AFBid}_i(B_k))$ —meaning that at least two bidders have bid for the fish at price *i*
- *bad*: a domain proposition meaning the fish is spoiled

Armed with the above, we can now state the commitments associated with the different messages in the fish market protocol.

- Payment of i from B_j : $discharge(B_j, C(B_j, A, FM, AFmoney_i))$
- Delivering fish to B_j : $discharge(A, C(A, B_j, FM, AFfish))$
- Yes from B_j (for price i): $create(B_j, Bid_i(B_j))$
- No from B_j (for price i): true
- Advertise to B_j (for price i): $create(A, Ad_i(B_j))$

Further, the protocol includes metacommitments that are not associated with any single message. In the present protocol, these metacommitments are of the context itself to release a committing party under certain circumstances. For practical purposes, we could treat these as metacommitments of the creditor.

- High demand: $C(FM, A, FM, AG[demand_i \rightarrow AFrelease(FM, C(A, B_j, FM, AFfish))])$
- Bad fish: $C(FM, B_j, FM, AG[bad \rightarrow AFrelease(FM, C(B_j, A, FM, AFmoney_i))])$

In addition, in a monotonic framework, we would also need to state the completion requirements to ensure that only the above actions are performed.

The auctioneer does not commit to a price if no bid is received. If more than one bid is received, the auctioneer is released from the commitment. Notice that the auctioneer can exit the market or adjust the price in any direction if a unique *Yes* is not received for the current price $money_i$. It would neither be rational for the auctioneer to raise the price if there are no takers at the present price, nor to lower the price if takers are available. However, the protocol *per se* does not legislate against either behavior. ■

The *No* messages have no significance on commitments. They serve only to assist in the coordination so the context can determine if enough bids are received. The lower-level aspects of coordination are not being studied in this paper. Now we can see how the reasoning takes place in a successful run of the protocol.

Example 6. The auctioneer sends out an advertisement, which commits the auctioneer to supplying the fish if he receives a suitable bid. This commitment will be discharged if $AG[Bid_i(B_j) \rightarrow AFcreate(A, C(A, B_j, FM, AFfish))]$ holds. When $Bid_i(B_j)$ is sent by B_j , the bidder is committed to the bid, which is discharged if $AG[fish \rightarrow AFcreate(B_j, C(B_j, A, FM, AFmoney_i))]$ holds. To discharge the advertisement, the auctioneer must eventually create a commitment to eventually supply the fish. If he does not create this commitment, he is in violation. If he

creates it, but does not supply the fish, he is still in violation. If he supplies the fish, the bidder is then committed to eventually forming a commitment to supply the money. If the bidder does so, the protocol is executed successfully. ■

3.3. REASONING WITH THE CONCRETE MODEL

Now we explain the main reasoning steps in our approach and show that they are sound. The main reasoning with models applies the CTL model-checking algorithm on a model and a formula denoting the conjunction of the specifications. The algorithm evaluates whether the formula holds in the initial state of the model. Thus a concrete version of the model M (see Section 2.2) is essential. For the purposes of the semantics, we must define a global model with respect to which commitment protocols may be specified. Intuitively, a protocol specification tells us which behaviors of the entire system are correct. Thus, it corresponds naturally to a global model in which those behaviors can be defined.

Our specific concrete model identifies states with messages. Recall that the timestamp of a message is the clock vector attached to it. The states are ordered according to the timestamps of the messages. The proposition true in a state is the one corresponding to the operation that is performed by the message.

Definition 8. $\mathbf{Q} = \{m : m \text{ is a message}\} \cup \{\vec{0}\}$ ■

Definition 9. For $s, t \in \mathbf{Q}$, $s < t$ iff $\text{timestamp}(s) \prec \text{timestamp}(t)$ ■

Definition 10. For $s \in \mathbf{Q}$, $\mathbf{I}(s) = \{\text{the operations executed by message } s\}$ ■

The structure $M_Q = \langle \mathbf{Q}, <, \mathbf{I} \rangle$ is a *quasimodel*. (Here and below, we assume that $<$ and \mathbf{I} are appropriately projected to the available states.) M_Q is structurally a model, because it matches the requirements of Definition 3. However, M_Q is not a model of the computations that may take place, because the branches in M_Q are concurrent events and do not individually correspond to a single path. A quasimodel can be mapped to a model, $M_S = \langle \mathbf{S}, <, \mathbf{I} \rangle$ with an initial state $\vec{0}$, by including all possible interleavings of the transitions. That is, \mathbf{S} would include a distinct state for every message in each possible ordering of the messages in \mathbf{Q} that is consistent with the temporal order $<$ of M_Q . The relation $<$ can be suitably defined for M_S . However, there is potentially an exponential blowup in that the size of \mathbf{S} may be exponentially greater than the size of \mathbf{Q} .

Theorem 3 shows that naively treating a quasimodel as if it were a model is correct. Thus, the above blowup can be eliminated entirely. Our construction

ensures that all the events relevant to another event are totally ordered with respect to each other. Notice that, as showing in Figure 3, the construction may appear to require one more message than necessary for the *assign* and *delegate* operations. This linear amount of extra work (for the entire set of messages), however, pays off in reducing the complexity of our reasoning algorithm. In the following, p refers to the proposition (of the form $\text{AG}[q \rightarrow \text{AF}r]$) of a metacommitment, which becomes true when the metacommitment is discharged.

Definition 11. For a proposition p , p^T is the proposition obtained by substituting EF for AF in p . ■

Theorem 3. $M_Q \models_{\bar{0}} p^T$ iff $M_S \models_{\bar{0}} p$.

Proof. From Theorem 2 and the restricted structure of M_Q . ■

The above results show that compliance can be tested and without blowing up the model unnecessarily. However, we would like to test for compliance based on local information—so that any agent can decide for itself whether it has been wronged by another. For this reason, we would like to be able to project the global model onto local models for each agent, while ensuring that the local models carry enough information that they are indeed usable in isolation from other local models. Accordingly, we can define the construction of local models corresponding to an agent’s observations. This is simply by defining a subset of \mathbf{S} for a given agent a .

Definition 12. $\mathbf{S}_a = \{m : m \text{ is a message from or to } a\}$. $M_a = \langle \mathbf{S}_a, <, \mathbf{I} \rangle$. ■

Theorem 4 shows that if we restrict attention to commitments that the given agent can observe, then the projected quasimodel yields all and only the correct conclusions relative to the global quasimodel. Thus, if the interested party is vigilant, it can check if anyone else violated the protocol.

Theorem 4. $M_a \models_{\bar{0}} p^T$ if and only if $M_Q \models_{\bar{0}} p^T$, provided that a sees all the commitments mentioned in p .

Proof. From Theorem 2 and the construction of M_a . ■

Example 7. If one of the bidders backs down from a successful bid, the auctioneer immediately can establish that he is cheating, because the auctioneer is the creditor for the bidder’s commitment. However, a bidder cannot ordinarily decide whether the auctioneer is noncompliant, because the bidder does not see all relevant commitments based on which the auctioneer may be released from a commitment to the bidder. ■

Theorem 5 lifts the above results to sets of agents. Thus, a set of agents may pool their evidence in order to establish whether a third party is noncompliant.

Thus, in a setting with two bidders, a model that includes all their evidence can be used to determine whether the auctioneer is noncompliant. Ordinarily, the bidders would have to explicitly pool their information to do so. However, in a broadcast-based or outcry protocol (like a traditional fish market in which everyone is screaming), the larger model can be built by anyone who hears all the messages. Let A be a set of agents.

Definition 13. $\mathbf{S}_A = \bigcup_{a \in A} \mathbf{S}_a$. $M_A = \langle \mathbf{S}_A, <, \mathbf{I} \rangle$. ■

Theorem 5. Let the commitments observed by agents in A include all the commitments in p . Then $M_A \models_{\bar{0}} p^T$ iff $M_Q \models_{\bar{0}} p^T$.

Proof. From Theorem 2 and the construction of M_A . ■

Information about commitments that have been resolved, i.e., are not pending, is not needed in the algorithm, and can be safely deleted from each observer's model. This is accomplished by searching backward in time whenever something is added to the model. Pruning extraneous messages from each observer's model reduces the size of the model and facilitates reasoning about it. This simplification is sound, because the CTL specifications do not include nested commitments.

Mapping from an event-based to a state-based representation, we should consider every event as potentially corresponding to a state change. This approach would lead to a large model, which accommodates not only the occurrence of public events such as message transmissions, but also local events. Such an approach would thus capture the evolution of the agent's knowledge about the progress of the system, which would help in accommodating unreliable messaging. Our approach, as described above, loses some of the agents' knowledge by not separating events and states, but has all the details we need to assess compliance assuming reliable messaging.

4. Discussion

Given the autonomy and heterogeneity of agents, the most natural way to treat interactions is as communications. A communication protocol involves the exchange of messages with a streamlined set of tokens. Traditionally, these tokens are not given any meaning except through reference to the beliefs or intentions of the communicating agents. By contrast, our approach assigns *public*, i.e., observable, meanings in terms of social commitments. Viewed in this light, *every communication protocol is a commitment protocol*.

Formulating and testing compliance of autonomous and heterogeneous agents is a key prerequisite for the effective application of multiagent systems in open environments. As asserted by Chiariglione, minimal specifications based on external behavior will maximize interoperability [7]. The research

community has not paid sufficient attention to this important requirement. A glaring shortcoming of most existing semantics for agent communication languages is their fundamental inability to allow testing for the compliance of an agent [16, 22]. The present approach shows how that might be carried out.

While the purpose of the protocols is to specify legal behavior, they should not specify rational behavior. Rational behavior may result as an indirect consequence of obeying the protocols. However, not adding rationality requirements leads to more succinct specifications and also allows agents to participate even if their rationality cannot be established by their designers.

The compliance checking procedure can be used by any agent who participates in, or observes, a commitment protocol. There are two obvious uses. One, the agent can track which of the commitments made by others are pending or have been violated. Two, it might track which of its own commitments are pending or whose satisfaction has not been acknowledged by others. The agent can thus use the compliance checking procedure as an input to its normal processes of deliberation to guide its interactions with other agents.

We have so far discussed how to detect violations. Once an agent detects a violation, as far as the above method is concerned, it may proceed in any way. However, some likely candidates are the following. The wronged agent may

- inform the agents who appeared to have violated their commitments and ask them to respect the applicable metacommitments
- inform the context, who might penalize the guilty parties, if any; the context may require additional information, e.g., certified logs of the messages sent by the different agents, to establish that some agents are in violation.
- inform other agents in an attempt to spoil the reputation of the guilty parties.

4.1. LITERATURE

Some of the important strands of research of relevance to commitment protocols have been carried out before. However, the synthesis, enhancement, and application of these techniques on multiagent commitment protocols is a novel contribution of this paper. Interaction (rightly) continues to draw much attention from researchers. Still, most current approaches do not consider an explicit execution architecture (however, there are some notable exceptions, e.g., [8, 9, 18]). Other approaches lack a formal underpinning; still others focus primarily on monolithic finite-state machine representations for protocols. Such representations can capture only the lowest levels of a multiagent

interaction, and their monolithicity does not accord well with distributed execution and compliance testing. Model checking has recently drawn much attention in the multiagent community, e.g., [2, 17]. However, these approaches consider knowledge and related concepts and are thus not directly applicable for behavior-based compliance.

4.2. FUTURE DIRECTIONS

The present approach highlights the synergies between distributed computing and multiagent systems. Since both fields have advanced in different directions, a number of important technical problems can be addressed by their proper synthesis. One aspect relates to situations where the agents may suffer a Byzantine failure or act maliciously. Such agents may fake messages or deny receiving them. How can they be detected by the other agents? Another aspect is to capture additional structural properties of the interactions so that noncompliant agents can be more readily detected. Alternatively, we might offer an assistance to designers by synthesizing skeletons of agents who participate properly in commitment protocols. Lastly, it is well-known that there can be far more potential causes than real causes [15]. Can we analyze conversations or place additional, but reasonable, restrictions on the agents that would help focus their interactions on the true relationships between their respective computations? We defer these topics to future research.

Acknowledgements

This work is supported by the National Science Foundation under grants IIS-9529179 and IIS-9624425, and IBM corporation. We are indebted to Feng Wan and Sudhir Rustogi for useful discussions and to the anonymous reviewers for helpful comments.

References

1. Gul A. Agha and Nadeem Jamali. Concurrent programming for distributed artificial intelligence. In [21], chapter 12, pages 505–534. 1998.
2. Massimo Benerecetti, Fausto Giunchiglia, and Luciano Serafini. Model checking multi-agent systems. *Journal of Logic and Computation*, 8(3):401–423, June 1998.
3. Kenneth P. Birman. The process group approach to reliable distributed computing. *Communications of the ACM*, 36(12):36–53, December 1993.
4. Kenneth P. Birman. A response to Cheriton and Skeen’s criticism of causal and totally ordered communication. *Operating Systems Review*, 28(1):11–21, 1994.
5. Nicholas Carriero and David Gelernter. Coordination languages and their significance. *Communications of the ACM*, 35(2):97–107, February 1992.
6. David R. Cheriton and Dale Skeen. Understanding the limitations of causally and totally ordered communication. In *Proceedings of the 14th ACM Symposium on Operating System Principles (SOSP)*, pages 44–57. ACM Press, December 1993.

7. Leonardo Chiariglione. Foundation for intelligent physical agents (FIPA) scope, 1998. www.fipa.org/library/scope.html.
8. Paolo Ciancarini, Andreas Knoche, Robert Tolksdorf, and Fabio Vitali. PageSpace: An architecture to coordinate distributed applications on the web. *Computer Networks and ISDN System*, 28(7–11):941–952, 1996. Proceedings of the 5th International World Wide Web Conference.
9. Paolo Ciancarini, Robert Tolksdorf, Fabio Vitali, Davide Rossi, and Andreas Knoche. Coordinating multiagent applications on the WWW: A reference architecture. *IEEE Transactions on Software Engineering*, 24(5):362–375, May 1998.
10. E. Allen Emerson. Temporal and modal logic. In Jan van Leeuwen, editor, *Handbook of Theoretical Computer Science*, volume B, pages 995–1072. North-Holland, Amsterdam, 1990.
11. Nissim Francez and Ira R. Forman. *Interacting Processes: A Multiparty Approach to Coordinated Distributed Programming*. ACM Press and Addison-Wesley, New York, 1996.
12. Jim Gray and Andreas Reuter. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, San Mateo, 1993.
13. Leslie Lamport. Time, clocks, and the ordering of events in a distributed system. *Communications of the ACM*, 21(7):558–565, July 1978.
14. Juan A. Rodríguez-Aguilar, Francisco J. Martín, Pablo Noriega, Pere Garcia, and Carles Sierra. Towards a test-bed for trading agents in electronic auction markets. *AI Communications*, 11(1):5–19, 1998.
15. Reinhard Schwarz and Friedemann Mattern. Detecting causal relationships in distributed computations: In search of the holy grail. *Distributed Computing*, 7(3):149–174, 1994.
16. Munindar P. Singh. Agent communication languages: Rethinking the principles. *IEEE Computer*, 31(12):40–47, December 1998.
17. Munindar P. Singh. Applying the mu-calculus in planning and reasoning about action. *Journal of Logic and Computation*, 8(3):425–445, June 1998.
18. Munindar P. Singh. A customizable coordination service for autonomous agents. In *Intelligent Agents IV: Proceedings of the 4th International Workshop on Agent Theories, Architectures, and Languages (ATAL-97)*, pages 93–106. Springer-Verlag, 1998.
19. Munindar P. Singh. Developing formal specifications to coordinate heterogeneous autonomous agents. In *Proceedings of the 3rd International Conference on Multiagent Systems (ICMAS)*, pages 261–268. IEEE Computer Society Press, July 1998.
20. Munindar P. Singh. An ontology for commitments in multiagent systems: Toward a unification of normative concepts. *Artificial Intelligence and Law*, 1999. In press.
21. Gerhard Weiß, editor. *Multiagent Systems: A Modern Approach to Distributed Artificial Intelligence*. MIT Press, Cambridge, MA, 1998.
22. Michael J. Wooldridge. Verifiable semantics for agent communication languages. In *Proceedings of the 3rd International Conference on Multiagent Systems (ICMAS)*, pages 349–356. IEEE Computer Society Press, July 1998.

Address for correspondence:

Department of Computer Science
Box 7534
North Carolina State University
Raleigh, NC 27695-7534, USA

`mvenkat@eos.ncsu.edu`, `singh@ncsu.edu`