

# Specification and Analysis of the AER/NCA Active Network Protocol Suite in Real-Time Maude

Peter Csaba Ölveczky<sup>1,2</sup>, José Meseguer<sup>1</sup>, and Carolyn L. Talcott<sup>3</sup>

<sup>1</sup> Department of Computer Science, University of Illinois at Urbana-Champaign

<sup>2</sup> Department of Informatics, University of Oslo

<sup>3</sup> Computer Science Laboratory, SRI International

August 1, 2004

## Abstract

This paper describes the application of the Real-Time Maude tool and the Maude formal methodology to the specification and analysis of the AER/NCA suite of active network multicast protocol components. Because of the time-sensitive and resource-sensitive behavior, the presence of probabilistic algorithms, and the composability of its components, AER/NCA poses challenging new problems for its formal specification and analysis. Real-Time Maude is a natural extension of the Maude rewriting logic language and tool for the specification and analysis of real-time object-based distributed systems. It supports a wide spectrum of formal methods, including: executable specification; symbolic simulation; breadth-first search for failures of safety properties in infinite-state systems; and linear temporal logic model checking of time-bounded temporal logic formulas. These methods complement those offered by network simulators on the one hand, and timed-automaton-based tools and general-purpose theorem provers on the other. Our experience shows that Real-Time Maude is well-suited to meet the AER/NCA modeling challenges, and that its methods have proved effective in uncovering subtle and important errors in the informal use case specification.

## 1 Introduction

Formal system specification and analysis requires exercising good judgment in making decisions that are not themselves amenable to full formalization. Questions such as: what should be formalized? at what level of abstraction? what are the relevant, perhaps informal, properties and how should they be formalized? have to be answered. Indeed, the success of the formal modeling enterprise greatly depends on how well they can be answered within a given formal framework, and on how the formal analysis tasks can then be supported by tools. Furthermore, as systems become more complex, their relevant properties tend to also grow in complexity and become more difficult to model and analyze formally, both because the formalization task becomes harder, and because system complexity tends to give rise to combinatorial explosions that make certain kinds of analyses unfeasible. Therefore, case studies involving challenging complex systems are enormously useful for testing the true mettle of a given formal framework and tool, and for extending the range of its applications. They are also one of the best ways of showing by example how a framework and its tools can be used to tackle a fairly wide range of similar problems.

This paper describes in detail our experience using the real-time rewriting logic formal framework [22] and its associated Real-Time Maude tool [24, 23] in analyzing a quite complex and sophisticated system, namely the AER/NCA active network protocol suite [9]. AER/NCA is a suite of four composable active network protocols, each achieving specific subgoals within the overall goal of making network multicast scalable, fault-tolerant, and congestion-avoidant. Challenges involved in formally specifying and analyzing the AER/NCA suite include:

1. Modeling time-sensitive behavior, including transmission delays, delay estimation, timers, and ordering.
2. Modeling resource-sensitive behavior, including link capacity, latency, congestion/cross traffic, and buffering.
3. Modeling probabilistic behavior, since several of the network algorithms involved are in fact probabilistic.
4. Both performance and correctness are critical aspects and are in fact inter-related, since the correct functioning of several protocol components consists precisely in achieving certain performance goals.
5. Composability is a key feature that must be supported to respect AER/NCA’s modular design, to avoid combinatorial explosions whenever possible, to analyze both individual components and their composed behavior, and to facilitate future reuses and extensions.

This is indeed a tall order. The first thing to observe is that a suite of protocols of this nature does not seem amenable to formalization and analysis within the known *decidable* frameworks for real-time systems and their associated tools [12, 8, 27]. There is indeed a tension between analytic power and expressiveness, where more power, sometimes even decidability, is typically purchased at the price of a more restricted formalism and range of applications. In this regard, Real-Time Maude is a quite expressive and general formal specification and analysis tool supporting a fairly wide middle ground between decidable real-time formalisms and tools on the one hand, and simulation tools on the other.

As we explain in detail in this paper, Real-Time Maude’s expressive features have allowed us to meet all the modeling challenges mentioned above in a satisfactory way. To begin with, we can easily support a distributed object-oriented formal specification style, which is ideal for modeling a network system; this is due to rewriting logic’s natural support for modeling distributed object systems [15]. Furthermore, object-oriented features such as inheritance were key in meeting the composability challenge (5). The modeling of resource-sensitive behavior (challenge (2)) is also a consequence of our object-oriented specification style: the key idea is to model the different resources explicitly as additional *objects* in the distributed object configuration of the system. In particular, we explicitly modeled network *links*, their capacity, their transmission delays, and the dropping of packets when links are full. Of course, meeting challenge (1) is most natural, since Real-Time Maude is by design a real-time formal specification and analysis language based on real-time rewrite theories [22]. In such theories, both 0-time, instantaneous transitions, and time-advancing “tick” transitions can be naturally specified by rewrite rules. The modeling of probabilistic behavior (challenge (3)) was also key to our application and deserves some discussion. The appropriate extension of rewriting logic to model a wide range of probabilistic systems, including probabilistic distributed algorithms, is the concept of a *probabilistic rewrite theory* [11]. This is a proper extension

of rewriting logic; however, for purposes of *simulating* the behavior of probabilistic systems using sampling techniques, one can associate to a probabilistic rewrite theory an ordinary rewrite theory that does the sampling, and can use Maude to execute its behaviors. This is exactly the approach taken in our modeling of the probabilistic algorithms in the AER/NCA suite, which we were able to do without any need for additional extensions to either Real-Time Maude or its underlying real-time rewriting logic formalism.

This leaves us with challenge (4), and the associated topic of passing from an informal to a formal specification and then formally analyzing the relevant properties. That is: (i) how was the formal specification of the AER/NCA system arrived at?; (ii) how well were we able to formally express the relevant system *properties*, involving in this case both correctness and performance aspects in a closely-related way?; (iii) how did we *formally analyze* such properties in Real-Time Maude; and (iv) what problems did we uncover as a result of those analyses and what were the *practical advantages* of making those discoveries?

Regarding (i), we started with a well-documented informal use-case-based specification of AER/NCA (available at [18, App. B]), provided to us by the designers and implementers of AER/NCA. We also benefited greatly from extended discussions with Mark Keaton and Steve Zabele, which were invaluable in making sure that we were correctly capturing the intended meaning of the informal specifications. In a sense this is the most important phase of any formal specification and analysis effort, and, as mentioned above, a phase not itself amenable to full formalization: there is nothing like an *algorithm* to pass from the informal to the formal specification.<sup>1</sup> As was to be expected, our interactions with the designers became a fruitful two-way street, in which our initial formalization attempts, particularly due to the fact that Real-Time Maude specifications are *executable*, required making explicit many implicit assumptions and uncovered a number of errors in the informal specification, even before any serious analysis was performed. An important, unexpected lesson learned from our extended discussions with some of the AER/NCA designers was the seemingly paradoxical fact that rewrite rules were much more intuitive and helpful to network engineers than the informal use-case descriptions. They explained to us that they had found use-cases ill-suited to gain an understanding of the protocols, and had translated them into state-transition diagrams to gain a better intuition about protocol behavior. Since rewrite rules are just parametric descriptions of local state transitions in a distributed system, this provided the level of description that network engineers were looking for, with the added bonus of executability and formal analysis.

Regarding (ii), the fact that time and resources were explicitly modeled in our specification made the formalization of relevant system properties —where we also benefited much from informal descriptions of such properties provided to us in discussions with Mark Keaton and Steve Zabele—relatively easy in several ways. First, time and resource utilization dynamics were directly inspectable in *executions* simulating the behavior of the different protocol components and their compositions. Second, even though we used standard linear time temporal logic (LTL) —which has no built-in notion of real-time and is in this sense less expressive than the various real-time temporal logics— to express the more sophisticated system properties, the fact that time and resource utilization was explicitly represented in the *states* of our model made it possible in practice to express in LTL all the desired properties —often involving real-time aspects and resource utilization—

---

<sup>1</sup>It is however possible to use formal methods and tools to support the passage from informal to formal specifications: a fairly large body of work on the formal underpinnings of UML, as well as work on passing from scenarios to system specifications and code [7] are good examples of work in this direction. But it does not follow from this that a *full* formalization of the entire process is possible: just the simple fact that *ethical* decisions are often involved in the choice of the relevant properties, particularly for safety-critical systems, seems to us a clear indication that it isn't.

by an adequate choice of *state predicates*, which queried the states for the relevant basic properties. Regarding (iii), three kinds of formal analysis were performed:

- *symbolic simulation*, by executing the specification starting from different initial states
- *breadth-first search* analysis, to find violations of safety properties, and
- *LTL model checking*, for *time-bounded* properties.

Provided the model of time used is discrete, as it was in our AER/NCA specification, breadth-first search analysis becomes a complete semi-decision procedure (if a safety violation exists, it will be found, although in practice this requires that sufficient memory is available). Similarly, under the discrete time assumption and reasonable requirements about our specification, LTL model checking of time-bounded properties is in fact a decision procedure. However, in the case of AER/NCA this kind of completeness cannot be claimed about our analyses. This is due to the *probabilistic* nature of several of the protocol components, and the corresponding sampling performed in the probabilistic transitions using a pseudo-random number generator object. As a consequence, the states we visited were determined by our choice of the pseudo-random number generator function: we would have visited different states if we had chosen a different such function. In summary, this just means that all errors we found in our analyses were always genuine errors; but there may be analyses not showing any errors for which, with a different way of sampling the probabilistic transitions, the same analysis could have uncovered a genuine error.

Regarding (iv), Real-Time Maude analysis uncovered a series of subtle and significant errors, which were easily traced to errors in the design of the original protocol suite. In particular, such formal analysis helped us to discover *all* design errors which were found independently by the protocol designers. None of these errors were disclosed to us as known by the designers until after we had found them. In addition, Real-Time Maude analysis found design errors which were *not* found during extensive traditional simulation and testing by the protocol designers. While some of these additional errors uncovered during Real-Time Maude analysis could be —and *were*— easily corrected, others indicated the need for a more thorough redesign of the original protocol. In our experience, Real-Time Maude analysis, apart from actually discovering *more* errors, required much less effort than traditional testing, because the executable formal specification can be subjected to exhaustive mechanical analysis without further work, and because it is easy to define different network topologies from which the specification can be analyzed in a variety of ways.

This paper is organized as follows. Sections 2 and 3 give a brief overview of, respectively, the AER/NCA protocol and Real-Time Maude. Section 4 describes how we met the modeling challenges described above, and how the AER/NCA protocol was specified and analyzed in Real-Time Maude. It includes parts of the informal specification to compare the two specification styles and to show how to get from a use-case based specification to a formal Real-Time Maude specification. Section 5 gives some concluding remarks. Finally, the Real-Time Maude tool —together with a user manual and related papers— and both the original informal use-case specification and the executable Real-Time Maude specification of the AER/NCA protocol suite are available at <http://www.ifi.uio.no/RealTimeMaude>.

## 2 The AER/NCA Protocol Suite

The AER/NCA protocol suite [9] combines several state-of-the-art techniques to achieve adaptive reliable multicast in *active networks*<sup>2</sup>. The protocol suite consists of a collection of composable protocol components supporting *active error recovery* (AER) and *nominee-based congestion avoidance* (NCA) features, and makes use of the possibility of having some processing capabilities at “active nodes” between the sender and the receivers to achieve scalability and efficiency. A high-level overview of the protocol suite, together with architectural requirements and simulation results, is given in [9]. The protocol itself was originally specified informally as a set of use cases. The Real-Time Maude formalization and analysis described in this paper led to a new version of the detailed informal protocol specification.

The goal of reliable multicast is to send a sequence of data packets from a sender to a group of receivers. Packets may be lost due to congestion in the network, and it must be ensured that each receiver eventually receives each data packet. Most multicast protocols are either not scalable or do not guarantee delivery for reasons which include the following [9]:

- To ensure reliability, the sender must be given feedback from the receivers, either by acknowledging the reception of data packets (ACK), or by signaling the lack of an expected packet (NAK). When there are many receivers, and each one frequently sends (positive or negative) acknowledgments to the single sender, then the sender—and the links closest to it in the network—easily become overwhelmed by this traffic.
- If there are many receivers, then some packet will be lost *somewhere* most of the time, keeping the sender busy with retransmissions. Furthermore, the sender has to multicast the repair packet to all the receivers—even though only a small group may have lost the packet—thereby increasing congestion in the network, or the sender must unicast the repair packet to the receivers, which is not desirable either for efficiency purposes when the losses are high.

The main design goal of the protocol is to minimize as much as possible the number of packet transmissions to achieve efficient, reliable, and scalable multicast. In addition, the protocol should find the appropriate sending rate to ensure that there is some bandwidth left for competing unicast TCP sessions.

### 2.1 Repair Servers

To overcome the above problems, Kasera et al. [9] suggested the use of *active services* at strategic locations inside the network. These active services can execute application-level programs inside routers, or, equivalently, on servers co-located with routers along the physical multicast distribution tree. By caching packets, these active services can subcast lost packets directly to “their” receivers, thereby localizing loss recovery, making loss recovery more efficient while solving the problem of retransmission scoping. They call such an active service, which may have fairly limited buffering capacity, a *repair server*. If a repair server does not have the missing packet in its cache, it aggregates all the negative acknowledgments (NAKs) it receives, and sends only one request for the lost packet toward the sender, solving the problem of feedback implosion at the sender.

---

<sup>2</sup>Active networks allow users to inject programs into the nodes of the network.

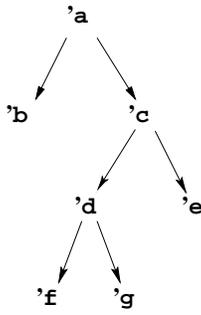


Figure 1: A multicast distribution tree.

*Terminology:* In this work, we abstract from routers which do not support active services, so that we regard the *multicast distribution tree* as having the *sender* at its root, the *receivers* in the multicast group as its leaf nodes, and the *repair servers* as its internal nodes. In this tree, the first node on the path from a node  $n$  to the root is called the *parent* of  $n$ . The *siblings* and the *children* of a node can be defined analogously. The parent is sometimes also denoted the *upstream* node of a given node, and children are denoted as the *downstream* nodes. We use the expression *the (upstream) repair server of a node  $n$*  to denote the parent of node  $n$ , which is therefore a repair server or the sender. For example, the multicast tree in Fig. 1 has sender 'a, repair servers 'c and 'd, and receivers 'b, 'e, 'f, and 'g. The sibling of node 'c is 'b, and the repair server of node 'c is 'a.

## 2.2 Overview of the Protocol

The AER/NCA protocol suite consists of the following four interconnected components:

- The *repair service* (RS) component deals with packet losses and tries to ensure that each packet is eventually received by all receivers in the multicast group. To enhance efficiency, loss recovery should happen as close as possible to the nodes where the losses were detected.
- The *rate control* (RC) component of the protocol aims at dynamically adjusting the rate by which the sender sends (original) data packets, so that the frequency decreases when many packets are lost (as the loss of a substantial number of packets indicates congestion due to a too high frequency in the sending of packets), and increases in time intervals when no, or few, packet losses are detected.
- The sender needs feedback about discovered packet losses to adjust its sending rate. The *nomination* (NOM) component aims at finding the “worst” receiver, based on the loss rates and the distance to the sender. The sender takes only the losses reported from this *nominee* receiver into account when determining the sending rate, instead of letting *all* receivers report their loss rates (which would result in too many messages being sent around just to determine the loss rate).
- The RTT component computes various *round trip time* values (the time it takes for a packet to travel from a given node to another given node, and back) in the network. These values

are needed for determining the sending rate and the nominee, and to decide how frequently to check for missing packets.

We give a high-level description below of the behavior of the protocol.

### 2.2.1 Active Error Recovery

The *sender* sends a sequence of data packets along the multicast distribution tree to the multicast group. It handles messages indicating the loss of a packet by resending the packet (the packets are identified by their sequence numbers), if the desired retransmission has not been done before.

When a data packet reaches a *repair server*, the repair server *caches* the packet before subcasting it downstream in the multicast distribution tree.

When a *repair server* or a *receiver* discovers that it has not received a data packet, it signals this by sending a *NAK* packet with the sequence number of the missing packet to its upstream repair server. This *NAK*-sending may be somewhat delayed, and is repeated at certain retransmission intervals until the missing data packet is received. The repair server or receiver needs an estimate of its round trip time to the sender to know how frequently it should resend this *NAK* packet. When a repair server discovers a packet loss, it also subcasts a *NAK* packet downstream to indicate to its children that it has started the loss recovery process for the missing data packet. Finally, if a repair server receives a *NAK* from downstream, it subcasts the missing data packet if it has it.

### 2.2.2 Nominee-based Congestion Avoidance

The frequency at which the sender sends new data packets was not specified above. The task of the NCA part of the protocol suite is to find a sending rate such that the multicast session does not overly congest any path from the sender to the receiver. That is, the most congested path should be identified, and the sending frequency should be adjusted so that this path is not overly congested, in order to ensure that there is enough bandwidth for competing TCP-sessions, and that the worst receiver can handle the onslaught of packets. The congestion control part consists of two subtasks, each of which is treated by a separate protocol component:

- Finding the worst receiver.
- Adjusting the sending rate according to the loss rate at the worst receiver.

**Rate control.** Whenever the worst receiver, the *nominee* receiver, receives an original data packet, it sends a *congestion control message* (CCM) packet to the sender. The frequency increases when CCM packets are received by the sender, and decreases when the sender hasn't seen a CCM packet for a while. The actual computations are sophisticated.

**Nominee receiver selection.** The *nominee* receiver is the receiver with the highest value of  $rtt \times \sqrt{lpe}$ , where *rtt* is the estimated *round trip time* to the source, and *lpe* is the *loss probability estimate*.

Each *receiver* uses a sliding window to record the sequence of original packets it has received, based upon which it can compute its estimated loss rate, its *lpe*. Each receiver sends a *congestion status message* (CSM) packet with its current *lpe* and *rtt* estimates to its upstream repair server at regular intervals.

Each *repair server* receives CSM packets from each of its children, containing the nominee in the subtree where the child is the root, together with the associated *rtt* and *lpe* values. The repair server computes the nominee in its own subtree based on these received values. Whenever there is a change in the nominee or its data, the repair server sends a CSM packet to its upstream repair server with these new data. It also resends this data at regular intervals, and blanks out its data if its nominee data has not been updated for a while.

The *sender* chooses the nominee receiver based on the CSM packets from its children. When a new nominee receiver is found, the sender sends a *nominee activation message* (NAM) packet to the old and the new receiver to notify them of their change of status (since (only) the new receiver should report back acknowledgments). Finally, the sender resends a NAM packet to the nominee at regular intervals.

### 3 Real-Time Maude

Real-Time Maude [24, 23] is a language and tool extending Maude [3, 4] to support the formal specification and analysis of *real-time* and *hybrid* systems. The specification formalism is based on *real-time rewrite theories* [22] —an extension of *rewriting logic* [2, 14]— and emphasizes *ease* and *generality* of specification. It is particularly suitable to specify distributed real-time systems in an object-oriented style.

Real-Time Maude specifications are *executable* under reasonable assumptions, so that a first form of formal analysis consists in simulating the system’s progress in time by *timed rewriting*. This can be very useful for debugging the specification; but of course, any such execution gives us only *one* behavior among the many possible concurrent behaviors of the systems. To gain further assurance about a system design one can use *model checking* techniques that explore many different behaviors from a given initial state of the system. *Timed search* and *time-bounded linear temporal logic model checking* can analyze *all* behaviors (possibly relative to a chosen treatment of time, in case we have a dense time domain) from a given initial state up to a certain duration. By restricting search and model checking to behaviors up to a given duration, the set of reachable states can often be restricted to a finite set, which can then be subjected to model checking.

Real-Time Maude offers an alternative to informal specifications and their testing on simulation tools and testbeds by:

- providing a precise formal specification of the system which, being executable, can be simulated and tested directly;
- allowing the specification to be analyzed in many different ways, not just by simulating a few behaviors of the system, but by exhaustively exploring a wide range of different scenarios; and
- allowing the user to define the appropriate forms of communication at a high level of abstraction, instead of having to use a fixed set of communication primitives.

On the other side of the spectrum, Real-Time Maude complements *formal* tools such as the timed/hybrid automaton-based tools Kronos [27], UPPAAL [12], and HyTech [8] by providing a more general specification formalism which supports well the specification and analysis of “infinite-state” systems with different communication and interaction models and with advanced object-oriented and modularity features. Such systems usually fall outside the decidable fragments supported by the aforementioned tools. Finally, some tools geared toward modeling and analyzing larger real-time systems, such as, e.g., IF [1], extend timed automaton techniques with explicit UML-inspired constructions for modeling objects, communication, and some notion of data types. Real-Time Maude complements such tools not only by the full generality of the specification language, but, most importantly, by its simplicity and clarity: A simple and intuitive formalism is used to specify both the data types (by *equations*) and dynamic and real-time behavior of the system (by *rewrite rules*). Furthermore, the operational semantics of a Real-Time Maude specification is clear and easy to understand.

Real-Time Maude is implemented in Maude as an extension of Full Maude [4, Part II]. The tool achieves high performance by exploiting as much as possible the underlying Maude engine.

### 3.1 Preliminaries: Object-Oriented Specification in Maude

Since Real-Time Maude specifications extend Maude specifications, we first recall object-oriented specification in Maude. A Maude module specifies a *rewrite theory* of the form  $(\Sigma, E \cup A, \phi, R)$ , where  $(\Sigma, E \cup A)$  is a *membership equational logic* [16] theory with  $\Sigma$  a signature,  $E$  a set of conditional equations and memberships, and  $A$  a set of equational axioms such as associativity, commutativity, and identity, so that equational deduction is performed *modulo* the axioms  $A$ . The theory  $(\Sigma, E \cup A)$  specifies the system’s state space as an algebraic data type.  $\phi$  is a function which associates to each function symbol  $f \in \Sigma$  its *frozen*<sup>3</sup> argument positions [4], and  $R$  is a collection of *labeled conditional rewrite rules* specifying the system’s local transitions, each of which has the form<sup>4</sup>

$$[l] : t \longrightarrow t' \text{ if } \bigwedge_{i=1}^n u_i \longrightarrow v_i \wedge \bigwedge_{j=1}^m w_j = w'_j,$$

where  $l$  is a *label*. Intuitively, such a rule specifies a *one-step transition* from a substitution instance of  $t$  to the corresponding substitution instance of  $t'$ , *provided* the condition holds; that is, corresponding substitution instances of the  $u_i$  can be rewritten (possibly in several steps) to those of the  $v_i$ , and the substitution instances of the equalities  $w_j = w'_j$  follow from  $E \cup A$ . The rules are implicitly universally quantified by the variables appearing in the  $\Sigma$ -terms  $t$ ,  $t'$ ,  $u_i$ ,  $v_i$ ,  $w_j$ , and  $w'_j$ . The rewrite rules are applied *modulo* the equations  $E \cup A$ .<sup>5</sup>

We briefly summarize the syntax of Maude. *Functional* modules and *system* modules are, respectively, equational theories and rewrite theories, and are declared with respective syntax `fmod ... endfm` and `mod ... endm`. *Object-oriented* modules provide special syntax to specify concurrent object-oriented systems, but are entirely reducible to system modules; they are declared with the

---

<sup>3</sup>Rewrites cannot take place in a frozen argument position of a function symbol, so that a term  $f(t_1, \dots, t_i, \dots, t_n)$  will *not* rewrite to  $f(t_1, \dots, u_i, \dots, t_n)$  when  $t_i$  rewrites to  $u_i$  if  $i \in \phi(f)$ .

<sup>4</sup>In general, the condition of such rules may not only contain rewrites  $u_i \longrightarrow v_i$  and equations  $w_j = w'_j$ , but also *memberships*  $t_k : s_k$ ; however, the specifications in this paper do not use this extra generality.

<sup>5</sup>Operationally, a term is reduced to its  $E$ -normal form modulo  $A$  before any rewrite rule is applied in Maude. Under the coherence assumption [26] this is a complete strategy to achieve the effect of rewriting in  $E \cup A$ -equivalence classes.

syntax (`omod ... endom`).<sup>6</sup> Immediately after the module's keyword, the *name* of the module is given. After this, a list of imported submodules can be added. One can also declare *sorts* and *subsorts* and *operators*. Operators are introduced with the `op` keyword. They can have user-definable syntax, with underbars ‘`_`’ marking the argument positions, and are declared with the sorts of their arguments and the sort of their result. Some operators can have equational *attributes*, such as `assoc`, `comm`, and `id`, stating, for example, that the operator is associative and commutative and has a certain identity element. Such attributes are then used by the Maude engine to match terms *modulo* the declared axioms. There are three kinds of logical statements, namely, *equations*—introduced with the keywords `eq`, or, for conditional equations, `ceq`—*memberships*—declaring that a term has a certain sort and introduced with the keywords `mb` and `cmb`—and *rewrite rules*—introduced with the keywords `r1` and `cr1`. The mathematical variables in such statements are either explicitly declared with the keywords `var` and `vars`, or can be introduced on the fly in a statement without being declared previously, in which case they must have the form `var : sort`. Finally, a comment is preceded by ‘`***`’ or ‘`---`’ and lasts till the end of the line.

In object-oriented Maude modules one can declare *classes* and *subclasses*. A class declaration

```
class C | att1 : s1, ... , attn : sn .
```

declares an object class  $C$  with attributes  $att_1$  to  $att_n$  of sorts  $s_1$  to  $s_n$ . An *object* of class  $C$  in a given state is represented as a term

$$\langle O : C \mid att_1 : val_1, \dots, att_n : val_n \rangle,$$

where  $O$  is the object's name or identifier, and where  $val_1$  to  $val_n$  are the current values of the attributes  $att_1$  to  $att_n$  and have sorts  $s_1$  to  $s_n$ . Objects can interact with each other in a variety of ways, including the sending of messages. A message declaration

```
msg m : p1 ... pn → Msg
```

defines the name of the message and the sorts of its parameters. In a concurrent object-oriented system, the state, which is usually called a *configuration*, has typically the structure of a *multiset* made up of objects and messages. Multiset union for configurations is denoted by a juxtaposition operator (empty syntax) that is declared associative and commutative and having the `none` multiset as its identity element, so that order and parentheses do not matter, and so that rewriting is *multiset rewriting* supported directly in Maude. The dynamic behavior of concurrent object systems is axiomatized by specifying each of its concurrent transition patterns by a rewrite rule. For example, the rule

```
r1 [1] : m(0,w) < 0 : C | a1 : x, a2 : y, a3 : z > =>
          < 0 : C | a1 : x + w, a2 : y, a3 : z > m'(y,x)
```

defines a (parameterized family of) transition(s) in which a message  $m$  having arguments  $0$  and  $w$  is consumed by an object  $0$  of class  $C$ , with the effect of altering the attribute  $a1$  of the object and of generating a new message  $m'(y, x)$ . By convention, attributes, such as  $a3$  in our example, whose values do not change and do not affect the next state of other attributes need not be mentioned in

---

<sup>6</sup>In Real-Time Maude, being an extension of Full Maude, module declarations and execution commands must be enclosed by a pair of parentheses.

a rule. Attributes like `a2` whose values influence the next state of other attributes or the values in messages, but are themselves unchanged, may be omitted from right-hand sides of rules. Thus the above rule could also be written

```
r1 [1] : m(0,w) < 0 : C | a1 : x, a2 : y > => < 0 : C | a1 : x + w > m'(y,x) .
```

A *subclass* inherits all the attributes and rules of its superclasses.

### 3.2 Object-Oriented Specification in Real-Time Maude

A Real-Time Maude *timed module* (syntax `(tmod ... endtm)`) specifies a *real-time rewrite theory* [22, 23], that is, a rewrite theory  $\mathcal{R} = (\Sigma, E \cup A, \phi, R)$ , such that:

1.  $(\Sigma, E \cup A)$  contains an equational subtheory  $(\Sigma_{TIME}, E_{TIME}) \subseteq (\Sigma, E \cup A)$ , satisfying the *TIME* axioms in [22], which specifies a sort `Time` as the time domain (which may be discrete or dense). Although a timed module is parametric on the time domain, Real-Time Maude provides some predefined modules specifying useful time domains. For example, the modules `NAT-TIME-DOMAIN-WITH-INF` and `POSRAT-TIME-DOMAIN-WITH-INF` define the time domain to be, respectively, the natural numbers and the nonnegative rational numbers, and contain the subsort declarations `Nat < Time` and `PosRat < Time`. These modules also add a supersort `TimeInf`, which extends the sort `Time` with an “infinity” value `INF`.
2. The sort of the “states” of the system has the designated sort `System`.
3. The rules in  $R$  are decomposed into:
  - “ordinary” rewrite rules model instantaneous change which is assumed to take zero time, and
  - *tick (rewrite) rules* that model the elapse of time in a system. Such tick rules must be of the form  $l : \{t\} \rightarrow \{t'\}$  **if cond**, where  $t$  and  $t'$  are of sort `System`, `{_}` is a built-in constructor of a new sort `GlobalSystem` which takes a term of sort `System` as argument, and where we have associated to such a rule a term  $u$  of sort `Time` intuitively denoting the *duration* of the rewrite. In Real-Time Maude, tick rules, together with their durations, are specified with the syntax

```
cr1 [l] : {t} => {t'} in time u if cond .
```

The initial state of a real-time system so specified must always have the form  $\{t_0\}$  (for  $t_0$  a ground term of sort `System`). The form of the tick rules ensures uniform time elapse in all parts of a system. We can then describe any finite *computation* from the initial state  $\{t_0\}$  as a sequence of one-step  $\mathcal{R}$ -rewrites  $\{t_0\} \rightarrow \{t_1\} \rightarrow \dots \rightarrow \{t_n\}$  with the rules in  $R$ , some of which may be instantaneous, and some tick rules. Furthermore, we assume that all the  $t_i$  are ground terms. The *duration* of such a computation is by definition the sum  $\sum_i^k \sigma_{i_j}(u_{i_j})$  corresponding to all the substitution instances of the terms  $u_{i_j}$  in all rewrite steps  $\{t_{i_j}\} \rightarrow \{t_{i_{j+1}}\}$  involving a tick rule with duration term  $u_{i_j}$  and a substitution  $\sigma_{i_j}$  [22].

*Timed object-oriented modules* (syntax `(tomod ... endtom)`) extend both object-oriented and timed modules to provide support for object-oriented specification of real-time systems. Timed

object-oriented modules include subsorts such as `MsgConfiguration`, `ObjectConfiguration`, and `NEConfiguration`, denoting, respectively, multisets of messages, multisets of objects, and non-empty configurations. The sort `Configuration` is declared to be a subsort of the sort `System`.

### 3.3 Rapid Prototyping and Formal Analysis in Real-Time Maude

We summarize below the Real-Time Maude analysis commands used in our case study. All Real-Time Maude analysis commands are described in [20], and their mathematical semantics is given in [23]. Note that, when the tick rule is *time-nondeterministic* (as in our case), then all analysis is performed with respect to the chosen *time sampling strategy* treatment of the tick rule(s) [24, 23].

#### 3.3.1 Rapid Prototyping: Timed Rewriting

Real-Time Maude’s *timed rewrite* command simulates *one* behavior of the system *up to a certain duration*. It is written with syntax

```
(trew t in time <= limit .)
```

where *t* is the term to be rewritten (“the initial state”), and *limit* is a ground term of sort `Time`. Our tool also provides facilities for *tracing* the rewrite steps performed in a simulation (see [20]).

#### 3.3.2 Search and Model Checking

Real-Time Maude provides a variety of search and model checking commands for further analyzing timed modules by exploring *all* possible behaviors—up to a given number of rewrite steps, duration, or satisfaction of other conditions—that can be nondeterministically reached from the initial state.

First of all, Real-Time Maude extends Maude’s *search* command—which uses a breadth-first strategy to search for states that are reachable from the initial state which match the *search pattern* and satisfy the *search condition*—to search for states and deadlocks which can be reached within a given time interval from the initial state. The search command has syntax

```
(tsearch t arrow pattern such that cond timeInterval .)
```

where *t* is the initial state (of sort `GlobalSystem`), *arrow* is either `=>*` (search for states reachable in zero or more steps) or `=>!` (search for “deadlocked” states which cannot be further rewritten), *pattern* is the search pattern, *cond* is a semantic condition on the variables in the search pattern, and *timeInterval* has either of the forms `with no time limit`, `in time op r`, or `in time-interval between op r and op' r'`, where each *op* and *op'* is one of `<`, `<=`, `>`, or `>=`, and *r* and *r'* are terms of sort `Time`. The command then returns all the states that are solutions of the search, but can be restricted to search only for at most *n* solutions by writing `(tsearch [n] ...)`. The `such that`-condition may be omitted.

Real-Time Maude provides commands for analyzing all behaviors from the initial state by searching for the *earliest* and the *latest* time when a certain state is reached for the first time. The command

```
(find earliest t =>* pattern such that cond .)
```

finds the earliest state reachable from  $t$  which is matched by  $pattern$  and satisfies  $cond$ . The command

```
(find latest t =>* pattern such that cond timeLimit .)
```

searches through all behaviors in a breadth-first way, and finds the *first* occurrence of a  $pattern$ -state satisfying  $cond$  in each behavior. Among these states, the state which took the longest time to reach is returned. The execution of this command will loop or return “not found in all computations” if there is a behavior in which the desired state cannot be reached within the time limit.  $timeLimit$  has either of the forms `with no time limit`, `in time < r`, or `in time <= r`.

Real-Time Maude has also commands for checking some simple temporal properties using *breadth-first* search techniques.<sup>7</sup> An example is the `check/untilStable` command which has the syntax

```
(check t |= pattern1 such that cond1 untilStable
      pattern2 such that cond2 timeLimit .)
```

It checks whether, for each behavior, a state matched by  $pattern_2$  and satisfying  $cond_2$  is found within the time limit, and each state following a  $pattern_2$ -state (within the time limit) is itself a  $pattern_2$ -state satisfying  $cond_2$ . In addition, each state in a behavior must be a  $pattern_1$ -state satisfying  $cond_1$  before a  $pattern_2$ -state is reached.

Finally, Real-Time Maude extends Maude’s *linear temporal logic model checker* [5, 4] to check whether each behavior “up to a certain time,” as explained in [23], satisfies a temporal logic formula. Restricting the computations to their time-bounded prefixes means that properties can be model checked in specifications that do not allow *Zeno behavior*, since (assuming a certain criterion for advancing time) only a finite set of states can be reached from an initial state. Because of the time-boundedness, liveness properties which hold in a specification may not hold for the time-bounded computations, and safety properties that do not hold for all computations may hold for all computations within the given time bound. Temporal logic model checking must be done in a module which includes both the module `TIMED-MODEL-CHECKER` and the module to be analyzed. *State propositions*, possibly parameterized, should be declared as operators of sort `Prop`, and their semantics should be given by (possibly conditional) equations of the form

$$\{statePattern\} \mid= prop = b$$

for  $b$  a term of sort `Bool`, which defines the state proposition  $prop$  to hold in all states  $\{t\}$  where  $\{t\} \mid= prop$  evaluates to `true`. It is not necessary to define explicitly the states in which  $prop$  does not hold. We may also define *clocked propositions*, which take the elapsed time into account, and which are defined

$$\{statePattern\} \text{ in time } r \mid= prop = b$$

A temporal logic *formula* is constructed by state propositions and temporal logic operators such as `True`, `False`,  $\sim$  (negation),  $\wedge$ ,  $\vee$ ,  $\rightarrow$  (implication),  $\square$  (“always”),  $\diamond$  (“eventually”),  $U$  (“until”), and  $W$  (“weak until”). The command

---

<sup>7</sup>Since the temporal logic model checker uses *depth-first* search techniques, there are cases in which the `check` command terminates even without a time limit, and where the temporal logic model checker would loop. One such example is shown in Section 4.13.

(mc *t* |=*t* *formula* *timeLimit* .)

is the timed model checking command which checks whether the temporal logic formula *formula* holds in all behaviors up to duration *timeLimit* starting from the initial state *t*. As explained in [23], timed model checking problems are equivalent to ordinary linear temporal logic (LTL) problems on a *transformed* (rewriting logic) theory. The Real-Time Maude system performs internally this transformation and sends the transformed model checking problem to the underlying Maude LTL model checker.

## 4 Formal Specification and Analysis of the AER/NCA Protocol Suite in Real-Time Maude

We summarize in this section the Real-Time Maude specification of the AER/NCA protocol suite. The given formal specification is based on version 1.0 of the informal specification of AER/NCA, as well as on consultations with the protocol developers. The paper [21] briefly outlined the specification and analysis of the AER/NCA protocol suite in version 1.0 of Real-Time Maude, and the thesis [18] presented that specification in its entirety. This paper describes in more detail the specification and analysis of the protocol suite in version 2.1 of our tool. This new version offers a simpler way of modeling object-oriented systems, as well as an entirely new set of efficiently implemented analysis commands [23]. In particular, all the analyses in [21, 18], for which user-defined strategies were needed, can now be performed directly using Real-Time Maude commands.

The formalization phase helped clarifying ambiguities and forced us to make implicit, but essential, knowledge explicit. After specifying the protocol, we executed it using Real-Time Maude’s default interpreter, and we found a first set of errors this way. The flaws detected during the prototyping of earlier versions of the specification were corrected, leading to a better specification, and more features, for example in the form of more accurate modeling of communication, were also added.

The specification is given in an object-oriented style, following the specification techniques suggested in [22, 23]. Although the four protocol components are closely inter-related, it is nevertheless important to analyze each component separately, as well as in combination. We manage this task by using object-oriented features such as inheritance.

We start this section by presenting our treatment of time. The time it takes for a packet to travel through a *link* between two nodes plays a crucial role in the protocol, since it determines the round trip between nodes, the retransmission intervals, the nominee receiver, and so on. We therefore need a fairly detailed model of communication —suggested to us by the protocol developers— which is presented in Section 4.6. In Section 4.9 we outline the class hierarchy which allows the analysis of the protocol components in isolation and in combination. Section 4.10 presents, for comparison purposes, both the informal specification and the Real-Time Maude specification of the RTT component. Sections 4.12, 4.14, and 4.16 present key aspects of the Real-Time Maude specification of, respectively, the NOM, RC, and RS components. Section 4.18 outlines how these components lead to a specification of the combined protocol. Sections 4.11, 4.13, 4.15, 4.17, and 4.18 give some examples of Real-Time Maude analysis of the protocol components and of the composite protocol.

## 4.1 Modeling the Time Domain

The protocol is parametric on the choice of a concrete time domain. We use the natural numbers as the time domain<sup>8</sup> by just importing into our specification the built-in Real-Time Maude module `NAT-TIME-DOMAIN-WITH-INF`, which defines the time domain `Time` to be the natural numbers, defines the subsort `NzTime` to be set of non-zero natural numbers, and defines a supersort `TimeInf` of `Time` with an additional infinity value `INF` and extends the standard functions on `Time` to `TimeInf`.

## 4.2 Modeling Time Elapse

In [22, 23] we suggested some specification techniques which have proved useful for specifying object-oriented real-time systems. In such systems it is often convenient to use functions `delta` and `mte` to define, respectively, the effect of time advance on a configuration of objects and messages, and the *maximum time elapse* allowed in a configuration before some action must be taken, and to let these functions distribute over the elements in a configuration. The operators  $\delta$  and *mte* should be declared to be *frozen* operators to avoid ill-timed rewrites or rewrites of the time domain [22, 23].

```
vars NECF NECF' : NEConfiguration .    var R : Time .

op delta : Configuration Time -> Configuration [frozen (1)] .
eq delta(none, R) = none .
eq delta(NECF NECF', R) = delta(NECF, R) delta(NECF', R) .

op mte : Configuration -> TimeInf [frozen (1)] .
eq mte(none) = INF .
eq mte(NECF NECF') = min(mte(NECF), mte(NECF')) .
```

The equations define how the functions distribute over the objects and messages in a configuration. To completely specify these functions, they must be defined for single objects as illustrated later in this paper.

Time elapse is modeled by the single tick rule

```
var OC : NEObjectConfiguration .    var R : Time .

crl [tick] : {OC} => {delta(OC, R)} in time R if R <= mte(OC) [nonexec] .
```

This tick rule is *time-nondeterministic*, as time may advance by *any* amount `R` less than or equal to `mte(OC)`, and is *nonexecutable* until we define a *time sampling strategy*. Since the time domain is the natural numbers, we could also cover the entire time domain by having a tick rule which advances time by one time unit. However, we notice that each instantaneous action in the protocol is triggered by an event such as the expiration of a timer or the reception of a message. That is, nothing “interesting” can happen in the time intervals *between* the “current” time and time `mte(OC)` thereafter. For efficiency purposes, we will therefore choose a time sampling strategy which always advances time as much as possible, and still “cover” all relevant behaviors.

---

<sup>8</sup>We could equally well have used the rational numbers as time domain, but the naturals numbers are sufficient.

The use of the variable `OC` of sort `NEObjectConfiguration` requires that the configuration only consists of objects when the `tick` rule is applied, and therefore forces messages which are not being transmitted over links to be treated without delay, because the above rule will not match, and therefore time will not advance, when there are messages present in the configuration.

This tick rule is the *only* tick rule in our specification; all other rules are instantaneous rules.

### 4.3 Timers

Many communication protocols, including AER/NCA, use *timers* to force an action to take place at a certain time. We model a timer belonging to an object as a class attribute of sort `TimeInf`, whose value is a time value  $r$  when the timer should expire in time  $r$ , and whose value is `INF` when the timer is turned off. For example, a timer `xTimer` in a class  $C$  may be declared

```
class C | ..., xTimer : TimeInf, ... .
```

We can force an action to take place when a timer expires by not letting time elapse beyond the expiration time of the timer, and by turning off or resetting the timer when the action associated with the expiration of the timer is performed. The functions `mte` and `delta` are typically defined on such classes  $C$  as follows:

```
var O : Oid .          var TI : TimeInf .          var R : Time .
eq delta(< O : C | ..., xTimer : TI, ... >, R) =
    < O : C | ..., xTimer : TI monus R, ... > .
eq mte(< O : C | ..., xTimer : TI, ... >) = TI .
```

The last equation should be replaced by an equation of the form

```
eq mte(< O : C | ..., xTimer : TI, ... >) = min(TI, ...) .
```

when the maximum time elapse depends also on other attributes than `xTimer`. An object may have more than one timer attribute, and one attribute could hold many timers.

### 4.4 Object Identifiers

We abstract from object addresses and define the set of object identifiers to be the set of *quoted identifiers*: subsort `Qid < Oid` . We also define *sets of object identifiers* and a supersort `DefOid` of object identifiers with a “default” value `noOid` (corresponding to a `null` pointer) as follows:

```
sorts DefOid OidSet .      subsorts Oid < DefOid OidSet .
op noOid : -> DefOid .
op none : -> OidSet .
op __ : OidSet OidSet -> OidSet [assoc comm id: none prec 15] .
```

## 4.5 A Set of Frequently Used Variables

To avoid repeating the declarations of variables in this exposition, we list below some variables used frequently in the formal specification. Other variables will be introduced together with the declarations of their sorts, and are considered to be declared throughout the paper.

```
var Q : Qid .
vars DO DO' DO'' : DefOid .
vars M M' : Msg .
vars R R' R'' R''' R'''' R''''' : Time .
var NZR : NzTime .
vars NZN NZN' NZN'' NZN''' NZN'''' : NzNat .
vars O O' O'' O''' O'''' : Oid .
var OS : OidSet .
vars MC MC' : MsgConfiguration .
vars TI TI' TI'' TI''' : TimeInf .
vars N N' N'' N''' : Nat .
vars X Y Z X' : Bool .
```

## 4.6 Modeling Communication and the Communication Topology

We abstract from the passive nodes in the network, and model the multicast communication topology by the multicast distribution tree which has the sender as its root, the receivers in the multicast group as its leaf nodes, and the repair servers as its internal nodes. The appropriate classes for these objects are defined as follows:

```
class Sendable | children : OidSet .
class Receivable | repairserver : DefOid .
class Sender . subclass Sender < Sendable .
class Receiver . subclass Receiver < Receivable .
class Repairserver . subclass Repairserver < Sendable Receivable .
```

For example, in a state representing the topology in Fig. 1 on page 6, the object (with identifier) 'c belongs to a subclass of the class `Repairserver`, and its `children` attribute has the value 'd 'e. (Its `repairserver` attribute should be set (to 'a) by the protocol; it has the value `noOid` initially.) The routing table is given implicitly by the multicast tree, which, in turn, is given by the objects together with the values of their `children` attributes. A multicast follows this multicast tree.

### 4.6.1 Links

Network performance and congestion (and the resulting loss of packets) are critical metrics in the AER/NCA protocol and should be explicitly modeled to faithfully analyze the protocol.

Packets are sent through bidirectional links, which model edges in a multicast distribution tree. The time it takes for a packet to arrive at a link's target node depends on the size of the packet, the number of packets in the link, and the link speed and propagation delay of the link, while the capacity of the link and the number of packets in it determines whether packets are lost. All these factors affect the network performance and the degree of congestion and are modeled by the following `Link` class, which is declared as follows:

```
class Link | up : Oid, down : Oid, bound : NzNat, propDelay : NzTime, linkSpeed : NzNat,
downMsgs : MsgList, downSize : Nat, upMsgs : MsgList, upSize : Nat .
```

where `up` and `down` denote the end nodes of the link, `bound` its capacity, `propDelay` its propagation delay, `linkSpeed` its link speed in megabits per second, and `downMsgs` is its buffer of length `downSize` of messages being sent downstream along the link. The link buffers the packets in lists which are declared as follows:

```
sort MsgList .      subsort Msg < MsgList .
op nil : -> MsgList .
op _+_ : MsgList MsgList -> MsgList [assoc id: nil] .
```

A packet is stored in the link as  $dly(p, r)$ , where  $r$  is the time until the packet can be delivered. The attempt to send a packet  $p$  through the link from  $a$  to  $b$  takes place by the method call/message `send(p, a, b)`. This `send`-request is treated by the link by discarding the packet if the link is full, and by computing the transmission delay and adding the packet to its buffer otherwise (rule `intoLinkDown`). When the packet has been in the link for the time it takes for the packet to travel through the link (i.e., its delay has reached 0), the link “delivers” the packet  $p$  by sending the message  $p$  from  $a$  to  $b$  to the global configuration (rule `outOfDownLink`), where it should be treated by object  $b$ .

```
vars ML ML' : MsgList .

crl [intoLinkDown] :
  send(M, 0, 0')
  < 0'' : Link | up : 0, down : 0', bound : N, propDelay : NZR,
      linkSpeed : NZN, downMsgs : ML, downSize : N' >
=>
  if N' < N then
    < 0'' : Link | downMsgs :
      ML + dly(M, max(NZR, greatestDly(ML)) + transDelay(M, NZN)),
      downSize : N' + 1 >
  else < 0'' : Link | > fi
  if leastDly(ML) /= 0 .

r1 [outOfDownLink] :
  < 0 : Link | down : 0', up : 0'', downMsgs : dly(M, 0) + ML, downSize : s N >
=>
  < 0 : Link | downMsgs : ML, downSize : N >
  (M from 0'' to 0') .
```

The treatment of packets from node `down` to node `up` is symmetric. The packet wrappers are declared as follows:

```
msg send : Msg Oid Oid -> Msg .
msg _from_to_ : Msg Oid Oid -> Msg .
msg dly : Msg Time -> Msg .
```

The *transmission delay* of a packet in a link is the packet size divided by the link speed, and the total delay of a packet entering a link is

$$\max(\text{propagation delay, } \mathit{maxDelayInLink}) + \text{transmission delay,}$$

where *maxDelayInLink* is the current delay of the last packet entered in the link, and is 0 if there are no packets in the link. Data packets are usually around 1500 bytes large, and all other kinds of packets are 64 bytes large. We declare the sorts `Packet` (for 64 bytes packets) and `LargePacket` (for 1500 bytes packets) as subsorts of the sort `Msg`, and define a function which computes the transmission delay of a given packet and link speed (in megabits per second) as follows:

```

sorts Packet LargePacket .      subsorts Packet LargePacket < Msg .

op transDelay : Msg NzNat -> Time .
var SMALLPACKET : Packet .      var LARGEPACKET : LargePacket .
eq transDelay(SMALLPACKET, NZN) = (64 * 8 + ((NZN * 1000) monus 1)) quo (NZN * 1000) .
eq transDelay(LARGEPACKET, NZN) = (1500 * 8 + ((NZN * 1000) monus 1)) quo (NZN * 1000) .

```

The functions `leastDly` and `greatestDly` compute the least and greatest delay of a message in a message list:

```

op leastDly : MsgList -> TimeInf .      op greatestDly : MsgList -> Time .
eq leastDly(nil) = INF .                eq greatestDly(nil) = 0 .
eq leastDly(dly(M, R) + ML) = R .      eq greatestDly(ML + dly(M, R)) = R .

```

Packets can be sent to a *group* of objects using the `multiSend` operator:

```

msg multiSend : Msg Oid OidSet -> Configuration .
eq multiSend(M, 0, (0' OS)) = send(M, 0, 0') multiSend(M, 0, OS - 0') .
eq multiSend(M, 0, none) = none .

```

The “timed” behavior of a link object is defined by the `mte` and `delta` functions. The delay associated with each message in the link’s buffer can be seen as a timer (intended to force the release of the message at the appropriate time), so that time acts on a link by decreasing the delay of each packet in the link’s buffer according to the time elapsed, and so that time does not advance beyond the time the first packet in the link is ripe for delivery:

```

eq delta(< Q : Link | downMsgs : ML, upMsgs : ML' >, R) =
  < Q : Link | downMsgs : ML minus R, upMsgs : ML' minus R > .

op _minus_ : MsgList Time -> MsgList . --- decrease delay of messages
eq nil minus R = nil .
ceq (ML + ML') minus R = (ML minus R) + (ML' minus R) if ML /= nil and ML' /= nil .
eq dly(M, R) minus R' = dly(M, R monus R') .

eq mte(< Q : Link | downMsgs : ML, upMsgs : ML' >) = min(leastDly(ML), leastDly(ML')) .

```

It may be worth noticing that nothing was said about communication aspects in the informal specification of the AER/NCA protocol suite. Giving a formal executable specification has the advantage of making explicit the communication assumptions.

## 4.6.2 The State of the System

The global state of the system will have the form  $\{t\}$ , where  $t$  is a configuration that consists of: (i) “node” objects, which are instances of subclasses of the classes `Sender`, `Repairserver`, or `Receiver`, (ii) links, which are objects of class `Link`, (iii) messages sent to and from the links, and (iv) unicast packets. The state may also contain some additional objects such as a random number generator, and/or objects representing a simplified view of the “environment.”

## 4.7 Random Numbers for Probabilistic Features

The AER/NCA protocol suite is a probabilistic protocol suite in that there are many places where a “randomly varying” value, “uniformly distributed” within a certain interval, is needed. To model such probabilistic features, we have one object of class `RandomNGen` with one attribute `seed` which carries the global “seed” for the function

```
op random : Nat -> Nat .
eq random(N) = ((104 * N) + 7921) rem 10609 .
```

which generates a pseudo-random sequence of natural numbers and which is an instance of a class of “good” pseudo-random number generators given in [10].

## 4.8 A Clock Class

Most classes have a clock attribute; they can be defined as subclasses of the following class:

```
class Clock | clock : Time .
```

## 4.9 The Class Hierarchy

The protocol components do not operate independently of each other. Some transitions are *composite transitions* which involve actions from different components. One such example is the reception of a data packet by the nominee receiver which involves detecting lost packets (RS component), updating the receiver’s *lpe* (NOM component), and acknowledging the data packet (RC component). Most transitions, however, are *independent transitions*, which only involve actions in one protocol component. Although the informal specification describes the behavior of the components (only) when all the components are executed together, it is important to be able to execute and analyze each component in isolation, as well as the protocol with all the components combined together.

The Real-Time Maude specification is designed using multiple class inheritance, so that each of the four protocol components RTT, NOM, RC, and RS can be executed separately as well as in combination. Figure 2 shows the class hierarchy for sender objects (with some classes omitted). Objects of the class `RTTsenderAlone` should be used in the initial state when the RTT part of the protocol is analyzed separately, while the sender object in the composite protocol should be an instance of the class `SenderCombined`. Since `RTTsenderAlone` and `SenderCombined` are subclasses of the class `RTTsender`, rules which model independent transitions should involve objects of class `RTTsender` to allow for maximal reuse of these rules. For composite transitions, we have defined their behavior when executed in a *single* component in rules involving objects of class

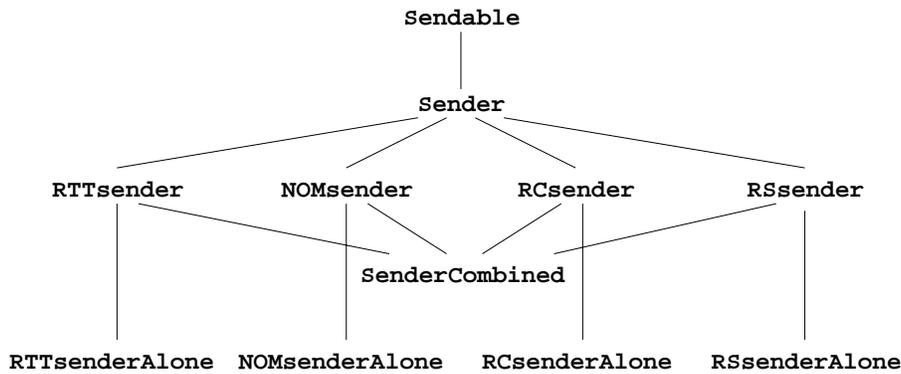


Figure 2: The sender class hierarchy.

`RTTsenderAlone`, and their behavior when executed in the composite protocol in rules involving objects of class `SenderCombined` as illustrated in Section 4.18. These techniques could also be used to specify the composition of two or three of the protocol components. The class hierarchies for repair servers and receivers are entirely similar.

## 4.10 The Formal and Informal Specifications of the RTT Component

In this section we present in detail both the informal specification and the Real-Time Maude specification of the RTT component.

The task of the RTT part of the protocol is to find, for each repair server and receiver object, the following values:

- `sourceRTT`: The round trip time (RTT) from the sender to the object.
- `maxUpRTT`: The maximal RTT from the object’s upstream repair server to any of that repair server’s children.

The round trip time values should be recently estimated values. The informal specification does not say anything about whether we are interested in the round trip times of large or small packets, or the time it takes for a small packet to go upstream plus the time it takes for a large packet to go downstream. The protocol presented in the formal specification below computes the round trip times of small packets, and can easily be changed by declaring the `getRTTRequest` and/or `getRTTResponse` packets to have sort `LargePacket`.

### 4.10.1 Class Declarations

The “state variables” of the nodes are declared as follows in the informal specification:

`Sender`:

*maxDownRTTSetTime*: Time that `maxDownRTT` was last updated.  
 Initialized to the `currentTime` in milliseconds.

*maxDownRTT*: Value currently being used as the largest RTT received from a directly supplied downstream receiver or repair server. Initialized to -1.  
*maxRecentDownRTT*: Largest RTT received from a directly supplied downstream receiver or repair server since *maxDownRTT* was last set. Initialized to 0.  
*sourceRTT*: RTT to the sender. Always 0.

Repair servers:

*resendInterval*: Time between successive Get-RTT requests.  
 Initialized to the value of *initialGetRTTCycleTime*.  
*maxDownRTTSetTime*: Time that *maxDownRTT* was last updated.  
 Initialized to the *currentTime* in milliseconds.  
*maxDownRTT*: Value currently being used as the largest RTT received from a directly supplied downstream receiver or repair server. Initialized to -1.  
*maxRecentDownRTT*: Largest RTT received from a directly supplied downstream receiver or repair server since *maxDownRTT* was last set. Initialized to 0.  
*maxUpRTT*: Largest RTT observed by the nearest upstream repair server (or sender).  
 Value is used to derive the suppression timer value. Initialized to -1.  
*myUpRTT*: RTT observed to the nearest upstream repair server (or sender).  
 Initialized to -1.  
*sourceRTT*: RTT to the sender. Value is used to derive the retransmit timer value.  
 Initialized to -1.

Receivers:

*resendInterval*: Time between successive Get-RTT requests.  
 Initialized to the value of *initialGetRTTCycleTime*.  
*maxUpRTT*: Largest RTT observed by the nearest upstream repair server (or sender). Value is used to derive the suppression timer value. Initialized to -1.  
*myUpRTT*: RTT observed to the nearest upstream repair server (or sender).  
 Initialized to -1.  
*sourceRTT*: RTT to the sender. Value is used to derive the retransmit timer value.  
 Initialized to -1.

We use the constant *INF* as “default” value instead of -1. In the Real-Time Maude specification, each state variable corresponds to an attribute in the class *RTTsender*, *RTTrepairserver*, or *RTTreceiver*. In addition, receivers and repair servers have an attribute *getRTTResendTimer* corresponding to the timer mentioned in the use cases. Since the repair servers perform many of the same transitions as the sender and the receivers, we find it convenient to define the superclasses *RTTsendable* and *RTTreceivable*, for, respectively, sender and repair servers, and repair servers and receivers. The class hierarchy and the packets involved in the protocol are given as follows:

```

*** All RTT objects are subclasses of RTT:
class RTT | sourceRTT : TimeInf .

*** Sender and repairserver:
class RTTsendable | maxDownRTT : TimeInf, maxDownRTTSetTime : Time,
                  maxRecentDownRTT : Time .
subclass RTTsendable < RTT Sendable Clock .

*** Repair server and receiver:
class RTTreceivable | resendInterval : Time maxUpRTT : TimeInf,

```

```

        myUpRTT : TimeInf, getRTTResendTimer : TimeInf .
subclass RTTreceivable < RTT Receivable Clock .

*** Classes for both stand-alone and combined protocols:
class RTTsender .           subclass RTTsender < RTTsendable .
class RTTreceiver .       subclass RTTreceiver < RTTreceivable .
class RTTrepairserver .   subclass RTTrepairserver < RTTreceivable RTTsendable .

*** Classes for stand-alone protocol only:
class RTTsenderAlone .     subclass RTTsenderAlone < RTTsender .
class RTTreceivableAlone . subclass RTTreceivableAlone < RTTreceivable .
class RTTreceiverAlone .  subclass RTTreceiverAlone < RTTreceiver RTTreceivableAlone .
class RTTrepairserverAlone .
subclass RTTrepairserverAlone < RTTrepairserver RTTreceivableAlone .

```

#### 4.10.2 Packet Declarations

Packets used in the RTT component are given as follows in the informal specification:

The algorithm functions via request-response messages exchanged between subscribers (receivers or repair servers) and their nearest upstream providers (repair servers or the sender). The Get-RTT request message and the Get-RTT response message have, respectively, the formats:

```

-----
| xmitTime | upRTT |           and           | xmitTime | peerGroupRTT | globalRTT |
-----

```

These packets are declared as follows in Real-Time Maude:

```

msg getRTTRequest : Time TimeInf -> Packet .
*** Usage: getRTTRequest(xmitTime, upRTT).

msg getRTTResponse : Time TimeInf TimeInf -> Packet .
*** Usage: getRTTResponse(xmitTime, peerGroupRTT, globalRTT)

```

#### 4.10.3 Specification of the Use Cases

We describe in this section the dynamics of the RTT component by presenting each use case in the informal specification followed by the corresponding rewrite rule(s) in the formal specification. The use cases can be summarized as follows:

**Use Case R.1.** in the informal specification defines the initial values of the state variables. The formal specification handles initialization by analyzing the protocol from initial states where the attributes have the given initial values.

**Use Case R.2.** A `startRTT` message initiates the sender, which sends *source path message* (SPM) packets downstream. Upon reception of such a packet, a node initiates its timer to a “random” value between 0 and 30 milliseconds, and subcasts the SPM packet downstream.

**Use Case R.3.** When a timer expires, its node sends a `getRTTRequest` packet to its (upstream) repair server with the current time, and the node’s current estimate of the RTT to its repair server.

**Use Case R.4.** When the upstream repair server receives this `getRTTRequest` packet, it may need to update its own `maxDownRTT`, and sends a `getRTTResponse` back to the downstream node, with the timestamp unchanged, and its current `sourceRTT` and current `maxDownRTT` values.

**Use Case R.5.** When the originator of this packet exchange receives the response, it can compute the “latest” RTT to its upstream repair server by just taking the current time minus the timestamp. Having this 1-step RTT, it adds this to the received `sourceRTT` value of its upstream repair server and gets its new `sourceRTT` estimate. It also compares the received `maxDownRTT` value with its own `maxUpRTT` estimate. Finally, the timer interval may be updated.

## Use Case R.2. Processing the First Received SPM Packet

This use case begins when the first SPM packet is received at a receiver or repair server. Each receiver or repair server starts a `Get-RTT` resend timer with a duration of a random variate, uniformly distributed between 0 and 1.0, times `implosionSuppressionInterval`

This use case ends when the `Get-RTT` resend timer has been set.

The execution of the *composite* protocol starts with the sender sending *source path message* (SPM) packets to its multicast group. To execute the RTT component in isolation, we use a message `startRTT` to start the RTT component, which the sender does by multicasting a SPM packet with sequence number 0 to the multicast group:

```
msg startRTT : Oid -> Msg .
rl [startRTT] :
  startRTT(0)
  < 0 : RTTsenderAlone | children : OS >
=>
  < 0 : RTTsenderAlone | >
  multiSend(SPMPacket(0), 0, OS) .
```

In the RTT component, such SPM packets are used to set the `repairserver` attributes and to start the protocol by initializing the timer to a random value between 0 and 30, which is the nominal value of the constant `implosionSuppressionInterval`. Upon the reception of the first SPM packet (`SPMPacket(0)`), a repair server must set its timer and subcast the SPM packet downstream (rule `R2rs`), while a receiver just sets its timer (rule `R2rcv`). The `RandomNGen` object provides the seed for computing the new “random” initial value of the timer:

```
rl [R2rs] :
  (SPMPacket(0) from 0' to 0)
  < 0'' : RandomNGen | seed : N >
```

```

    < 0 : RTTrepairserverAlone | children : OS >
=>
    < 0'' : RandomNGen | seed : random(N) >
    < 0 : RTTrepairserverAlone | repairserver : 0', getRTTResendTimer : random(N) rem 31 >
    multiSend(SPMPacket(0), 0, OS) .

r1 [R2rcv] :
    (SPMPacket(0) from 0' to 0)
    < 0'' : RandomNGen | seed : N >
    < 0 : RTTreceiverAlone | >
=>
    < 0'' : RandomNGen | seed : random(N) >
    < 0 : RTTreceiverAlone | repairserver : 0', getRTTResendTimer : random(N) rem 31 > .

```

### Use Case R.3. Get-RTT Resend Timer Service Routine

This use case begins when the Get-RTT resend timer expires at a receiver or repair server. Each receiver or repair server resets the Get-RTT resend timer using the current value of `resendInterval`, and subsequently sends a Get-RTT request packet to the nearest upstream repair server (or sender). The Get-RTT request packet fields are set as follows:

```

    xmitTime = currentTime
    upRTT    = myUpRTT

```

This use case ends when the Get-RTT request packet has been sent.

This use case, which describes how a node initiates a request/response round when its timer expires, is modeled formally as follows. (Remember that the `clock` attribute shows the current time.)

```

r1 [R3] :
    < 0 : RTTreceivable | clock : R, repairserver : 0', resendInterval : R',
                          myUpRTT : TI, getRTTResendTimer : 0 >
=>
    < 0 : RTTreceivable | getRTTResendTimer : R' >
    send(getRTTRequest(R, TI), 0, 0') .

```

### Use Case R.4. Processing a Received Get-RTT Request Packet

This use case begins when a Get-RTT request packet is received at a repair server or sender. The following processing is performed (`xmitTime` and `upRTT` are Get-RTT request packet fields):

```

    if (upRTT > maxDownRTT) {
        maxDownRTT          = upRTT
        maxDownRTTSetTime = currentTime in milliseconds
        maxRecentDownRTT   = 0
    }

```

```
else { if (upRTT > maxRecentDownRTT) { maxRecentDownRTT = upRTT } }
```

A check is then made to determine if `maxDownRTT` should be updated to the value of `maxRecentRTT` by comparing the update time against the current time in milliseconds:

```
if (currentTime > (maxDownRTTSetTime + updateWindowLength)) {
    maxDownRTT      = maxRecentDownRTT
    maxDownRTTSetTime = currentTime in milliseconds
    maxRecentDownRTT = 0
}
```

The repair server or sender then sends a Get-RTT response packet. The Get-RTT response packet fields are set as follows:

```
xmitTime      = xmitTime,
peerGroupRTT  = maxDownRTT
globalRTT     = sourceRTT
```

This use case ends when the repair server or sender sends a Get-RTT response packet.

This use case describes how Get-RTT request packets are treated. In the formal specification we split the treatment of `getRTTRequest` packets according to whether the `upRTT` value, that is, the second parameter of the packet, and the value of the attribute `maxDownRTT` are INF or time values (we exploit that the variables  $R, R', \dots$  range over time values, and  $TI, TI', \dots$  range over the sort `TimeInf` comprising time values and the value INF):

\*\*\* `upRTT` and `maxDownRTT` are both time values, and `upRTT > maxDownRTT`.

```
cr1 [R4a] :
  (getRTTRequest(R, R') from 0 to 0')
  < 0' : RTTsendable | clock : R'', sourceRTT : TI, maxDownRTT : R'' >
=>
  < 0' : RTTsendable | maxDownRTT : R', maxDownRTTSetTime : R'', maxRecentDownRTT : 0 >
  send(getRTTResponse(R, R', TI), 0', 0)      if R'' < R' .
```

\*\*\* `upRTT` is still a time value, but `maxDownRTT` is now INF.

```
rl [R4b] :
  (getRTTRequest(R, R') from 0 to 0')
  < 0' : RTTsendable | clock : R'', sourceRTT : TI, maxDownRTT : INF >
=>
  < 0' : RTTsendable | maxDownRTT : R', maxDownRTTSetTime : R'', maxRecentDownRTT : 0 >
  send(getRTTResponse(R, R', TI), 0', 0) .
```

The case where `upRTT` is INF, and the case where both `upRTT` and `maxDownRTT` are time values and `upRTT <= maxDownRTT`, are modeled by two additional rules in the same style.

## Use Case R.5. Processing a Received Get-RTT Response Packet

This use case begins when a Get-RTT response packet is received at a receiver or repair server. The receiver or repair server computes the local round trip time as:

```
myUpRTT = (currentTime - xmitTime)
```

Next, if either the value of `maxUpRTT` is equal to -1 or the value of `peerGroupRTT` is equal to -1, it immediately issues another Get-RTT request as follows:

```
if ((maxUpRTT == -1) || (peerGroupRTT == -1))
{ resendInterval = myUpRTT
  Cancel the current Get-RTT resend timer
  Start the Get-RTT timer with a duration of resendInterval
  Send a Get-RTT request packet (xmitTime = currentTime, upRTT = myUpRTT) }
```

The resend interval is doubled and limited to the value of `maxGetRTTCycleTime`:

```
resendInterval = (resendInterval * 2)
if (resendInterval > maxGetRTTCycleTime) { resendInterval = maxGetRTTCycleTime }
```

`maxUpRTT` is set to either the value of `peerGroupRTT` from the packet or the new value of `myUpRTT`, whichever is greater:

```
if (myUpRTT > peerGroupRTT) { maxUpRTT = myUpRTT }
else { maxUpRTT = peerGroupRTT }
```

Finally, if the value of `globalRTT` in the response packet is nonnegative, the receiver or repair server then computes and sets the value of `sourceRTT` as the sum of the value of `globalRTT` and the value of `myUpRTT`:

```
if (globalRTT >= 0) { sourceRTT = (globalRTT + myUpRTT) }
```

The treatment of `getRTTResponse` packets in the formal specification is divided into two cases, depending on whether both the `maxUpRTT` attribute value and the *peerGroupRTT* value (the second parameter) in the received packet are time values. In the combined protocol, other actions must also be taken when new RTT values are found. Therefore, the following rules apply to objects of class `RTTreceivableAlone`:

```
cr1 [R5a] :
  (getRTTResponse(R, TI, TI') from 0 to 0')
  < 0' : RTTreceivableAlone | clock : R', sourceRTT : TI'',
    resendInterval : R'', maxUpRTT : TI''' >
=>
  < 0' : RTTreceivableAlone | sourceRTT : (if TI' /= INF then
    TI' + (R' minus R) else TI'' fi),
    resendInterval : min(2 * (R' minus R), 3000),
    maxUpRTT : (if TI' /= INF then
    max(R' minus R, TI)
    else (R' minus R) fi),
    myUpRTT : R' minus R,
    getRTTResendTimer : R' minus R >
  send(getRTTRequest(R', R' minus R), 0', 0)
  if TI''' == INF or TI == INF .
```

\*\*\* Neither `peerGroupRTT` nor `maxUpRTT` has INF value:

```

rl [R5b] :
  (getRTTResponse(R, R', TI) from 0 to 0')
  < 0' : RTTreceivableAlone | clock : R'', sourceRTT : TI',
                                resendInterval : R''', maxUpRTT : R'''' >
=>
  < 0' : RTTreceivableAlone | sourceRTT : (if TI /= INF
                                           then TI + (R'' monus R) else TI' fi),
                                resendInterval : min(2 * R''', 3000),
                                maxUpRTT : max(R'' monus R, R'),
                                myUpRTT : R'' monus R > .

```

#### 4.10.4 Real-Time Behavior

Finally, we need to specify how time acts on RTT objects in the stand-alone protocol. A repair server or receiver has a timer attribute on which `mte` and `delta` work as described in Section 4.2. The objects also have a `clock` attribute which must be updated as time elapses:

```

eq delta(< 0 : RTTsenderAlone | clock : R >, R') =
  < 0 : RTTsenderAlone | clock : R + R' > .
eq delta(< 0 : RTTreceivableAlone | clock : R, getRTTResendTimer : TI >, R')
  = < 0 : RTTreceivableAlone | clock : R + R', getRTTResendTimer : TI monus R' > .

eq mte(< 0 : RTTsenderAlone | >) = INF .
eq mte(< 0 : RTTreceivableAlone | getRTTResendTimer : TI >) = TI .

```

### 4.11 Formal Analysis of the RTT Component in Real-Time Maude

This section illustrates how the RTT component has been analyzed using the Real-Time Maude tool.

#### 4.11.1 Defining a Time Sampling Strategy

Before any analysis can be undertaken, we must select a time sampling strategy to guide the application of the time-nondeterministic tick rule given in Section 4.2. As mentioned there, even though a strategy which advances time by one time unit in each tick would cover the time domain, we use for efficiency purposes a strategy which increases time by the maximum amount possible, since no instantaneous rule can be applied before time has advanced “as much as possible.” We declare this time sampling strategy by the command

```
Maude> (set tick max .)
```

and note that this strategy will apply to the analysis of all the protocol components.

### 4.11.2 Prototyping the RTT Component

In an object-oriented timed module `AER-RTT1` which imports the module `AER-RTT` specifying the RTT component, we define the following initial state `RTTstate2`. This state has the topology given in Fig. 1 and is parameterized by the initial value of the seed used by the random number generator:

```

op RTTstate2 : Nat -> GlobalSystem .
eq RTTstate2(N) =
({ startRTT('a)
  < 'a : RTTsenderAlone | clock : 0, sourceRTT : 0, children : 'b 'c, maxDownRTT : INF,
    maxDownRTTSetTime : 0, maxRecentDownRTT : 0 >
  < 'b : RTTreceiverAlone | ATTS-RCVR >
  < 'c : RTTrepairserverAlone | children : 'd 'e, ATTS-RS >
  < 'd : RTTrepairserverAlone | children : 'f 'g, ATTS-RS >
  < 'e : RTTreceiverAlone | ATTS-RCVR >
  < 'f : RTTreceiverAlone | ATTS-RCVR >
  < 'g : RTTreceiverAlone | ATTS-RCVR >
  < 'random : RandomNGen | seed : N >
  < 'ab : Link | up : 'a, down : 'b, bound : 5, propDelay : 21, linkSpeed : 1, ATTS-LINK >
  < 'ac : Link | up : 'a, down : 'c, bound : 21, propDelay : 28, linkSpeed : 3, ATTS-LINK >
  < 'cd : Link | up : 'c, down : 'd, bound : 9, propDelay : 23, linkSpeed : 1, ATTS-LINK >
  < 'ce : Link | up : 'c, down : 'e, bound : 4, propDelay : 17, linkSpeed : 1, ATTS-LINK >
  < 'df : Link | up : 'd, down : 'f, bound : 12, propDelay : 5, linkSpeed : 10, ATTS-LINK >
  < 'dg : Link | up : 'd, down : 'g, bound : 12, propDelay : 5, linkSpeed : 10, ATTS-LINK >}) .

```

where *ATTS-RCVR* stands for

```

clock : 0, sourceRTT : INF, repairserver : noOid, resendInterval : 200,
maxUpRTT : INF, myUpRTT : INF, getRTTResendTimer : INF ,

```

*ATTS-RS* stands for the multiset union (juxtaposition) of *ATTS-RCVR* and

```

maxDownRTT : INF, maxDownRTTSetTime : 0, maxRecentDownRTT : 0 ,

```

and *ATTS-LINK* stands for `downMsgs : nil, downSize : 0, upMsgs : nil, upSize : 0`.

Fig. 3 shows the multicast distribution tree of `RTTstate2`, where the number associated with each link indicates how much time it takes (namely, the propagation delay plus the transmission delay) for a small packet to travel through the link when the link is otherwise empty. For example, in otherwise empty links, the round trip time to the source from the nodes 'c', 'd', and 'e' is, respectively, 58, 106, and 94, and the `maxUpRTT` of these nodes is, respectively, 58, 48, and 48.

Real-Time Maude's timed rewrite command can be used to simulate *one* behavior of the RTT protocol up to time 1000:<sup>9</sup>

---

<sup>9</sup>The output of Real-Time Maude executions will be manually tabulated for readability purposes, and parts of the output omitted in the exposition will be replaced by '...'

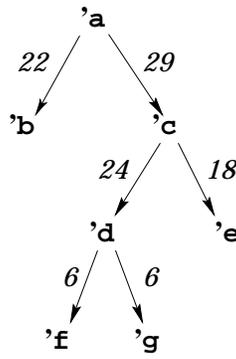


Figure 3: The multicast distribution tree corresponding to `RTTstate2`.

```
Maude> (trew RTTstate2(1) in time <= 1000 .)
```

```
Result ClockedSystem :
{< 'a : RTTsenderAlone | children : 'b 'c, clock : 907, maxDownRTTSetTime : 124,
    maxDownRTT : 58, maxRecentDownRTT : 58, sourceRTT : 0 >
  < 'b : RTTreceiverAlone | maxUpRTT : 58, sourceRTT : 44, ... >
  < 'c : RTTrepairserverAlone | maxUpRTT : 58, sourceRTT : 58, ... >
  < 'd : RTTrepairserverAlone | maxUpRTT : 48, sourceRTT : 106, ... >
  < 'e : RTTreceiverAlone | maxUpRTT : 48, sourceRTT : 94, ... >
  < 'f : RTTreceiverAlone | maxUpRTT : 12, sourceRTT : 118, ... >
  < 'g : RTTreceiverAlone | maxUpRTT : 12, sourceRTT : 118, ... >
  ... } in time 907
```

These `sourceRTT` and `maxUpRTT` values are as expected.

#### 4.11.3 Further Formal Analysis of the RTT Component

Although prototyping the RTT component gave the desired result, such execution explores just *one* of many possibly behaviors of a system, arbitrarily chosen by Real-Time Maude’s default interpreter. To gain further understanding—and to uncover errors—the specification can be subjected to further formal analysis using the search and model checking commands of Real-Time Maude.

The main property the stand-alone RTT protocol should satisfy is that, as long as no more than one packet travels simultaneously in the same direction in the same link, then:

- each computation will reach a state with the desired `sourceRTT` and `maxUpRTT` values within given time and depth limits (reachability); and
- once the correct values are found, they will not change within the given time limit (stability).

In addition, since a nominee receiver must be found before the whole protocol can start the transmission of data packets, and the RTT values are needed to find a nominee receiver, it is useful to know

- how long time it takes (in the worst case) to find the RTT values.

The first and last of these issues can be checked by Real-Time Maude's `find latest` command. The following command checks whether correct `sourceRTT` and `maxUpRTT` values of the objects in the topology given in Fig. 3 will be reached in *all* behaviors from initial state `RTTstate2(1)`, and the longest time needed to do so:

```
Maude> (find latest RTTstate2(1) =>*)
  {< 'b : RTTreceiverAlone | sourceRTT : 44, maxUpRTT : 58, ATTS1:AttributeSet >
  < 'c : RTTrepairserverAlone | sourceRTT : 58, maxUpRTT : 58, ATTS2:AttributeSet >
  < 'd : RTTrepairserverAlone | sourceRTT : 106, maxUpRTT : 48, ATTS3:AttributeSet >
  < 'e : RTTreceiverAlone | sourceRTT : 94, maxUpRTT : 48, ATTS4:AttributeSet >
  < 'f : RTTreceiverAlone | sourceRTT : 118, maxUpRTT : 12, ATTS5:AttributeSet >
  < 'g : RTTreceiverAlone | sourceRTT : 118, maxUpRTT : 12, ATTS6:AttributeSet >
  C:Configuration} in time < 5000 .)
```

to which Real-Time Maude answers

```
Result: {< 'b : RTTreceiverAlone | maxUpRTT : 58, sourceRTT : 44, ... >
  < 'g : RTTreceiverAlone | maxUpRTT : 12, sourceRTT : 118, ... >
  ... } in time 255
```

That is, it takes at most time 255 to reach a state with the desired RTT values. (A `find earliest` check showed that *earliest* possible time, in which the desired values can be found, is 181.)

The remaining task is therefore to check whether the correct RTT values can be altered once they are found. Unbounded model checking cannot check this property, since an infinite number of states can be reached from the initial state (the `clock` attribute, and therefore also other values such as the time stamps, can assume an infinite number of values). We can check the property for each computation “up to a certain time  $r$ ” in either of the two following ways. The first option is to use Real-Time Maude's built-in checker for `untilStable` properties:

```
Maude> (check RTTstate2(1) |= {C:Configuration} untilStable
  {< 'b : RTTreceiverAlone | sourceRTT : 44, maxUpRTT : 58, ATTS1:AttributeSet >
  < 'c : RTTrepairserverAlone | sourceRTT : 58, maxUpRTT : 58, ATTS2:AttributeSet >
  ...
  < 'g : RTTreceiverAlone | sourceRTT : 118, maxUpRTT : 12, ATTS6:AttributeSet >
  C:Configuration} in time < 1000 .)
```

Result: the property holds.

The other way of checking the stability property is to use Real-Time Maude's linear temporal logic model checker. The following module imports the timed model checker and the module `AER-RTT1`, and defines a proposition `CORRECT-RTT-VALUES`:

```
(tmod MODEL-CHECK-RTT is protecting AER-RTT1 .
  including TIMED-MODEL-CHECKER .

  op CORRECT-RTT-VALUES : -> Prop .
```

```

eq {< 'b : RTTreceiverAlone | sourceRTT : 44, maxUpRTT : 58 >
...
< 'g : RTTreceiverAlone | sourceRTT : 118, maxUpRTT : 12 >
C:Configuration}
|=
CORRECT-RTT-VALUES = true .
endtm)

```

We can then model check the temporal property that, within time 1000, the correct values are found in each behavior ( $\langle \rangle$  CORRECT-RTT-VALUES), and a correct value will not change once it is found (CORRECT-RTT-VALUES  $\Rightarrow$   $[\ ]$  CORRECT-RTT-VALUES)<sup>10</sup> :

```

Maude> (mc RTTstate2(1) |=t (<> CORRECT-RTT-VALUES) /\
(CORRECT-RTT-VALUES => [\ ] CORRECT-RTT-VALUES)
in time < 1000 .)

result Bool :
true

```

## 4.12 Formal Specification of the NOM Component

This section introduces the nominee selection component of the protocol and presents the most crucial parts of its formal specification.

An important goal for the AER/NCA protocol suite is “TCP-friendliness,” which mandates that a multicast session must not receive more bandwidth than competing TCP sessions on any of the source-to-destination paths in the multicast tree [9]. In order to achieve TCP-friendliness, the worst path in a multicast tree is determined as the path on which a TCP session will receive the least bandwidth, namely, the path with the highest value of  $r_{tt} * \sqrt{lpe}$ , where  $r_{tt}$  is the round trip time from the receiver to the sender, and  $lpe$  is the loss probability estimate of the receiver. The protocol behaves as follows, and determines the worst path and nominates the multicast receiver at the end of this path to send acknowledgments to the sender:

- Each receiver estimates its end-to-end packet loss probability (its  $lpe$ ) using a fixed size *sliding window*. Each receiver periodically unicasts its  $lpe$  value and its current round trip time estimate  $r_{tt}$  in a *congestion status message* (CSM) to its upstream repair server.
- Based on CSMs from its children, a repair server identifies the “worst” receiver in its subtree and unicasts the CSM of this worst receiver to its nearest upstream repair server.
- The sender receives periodic CSMs from its downstream repair servers and receivers, and uses the same method to select the worst receiver in the entire multicast group.
- Once the sender has identified the worst receiver, it unicasts (with no repair server intervention) a *nominee activation message* (NAM) to this receiver soliciting acknowledgments from it, and unicasts a NAM to the previous, if any, nominee receiver, to let the previous nominee know that it is no longer the nominee receiver. The sender resends a NAM every seven seconds to the nominee receiver until a different nominee is identified.

---

<sup>10</sup>Remember that  $P \Rightarrow Q$  is an abbreviation of  $[\ ] (P \rightarrow Q)$ .

### 4.12.1 Sender Protocol

The sender class is declared as follows. `NAMTimer` is used to periodically send NAM packets to the current nominee, and `csmNominee`, `csmLPE`, `csmRTT`, and `csmSetTime` denote, respectively, the current nominee receiver, its *lpe* and *rtt* values, and the last time these values were updated. (The sort `DefRat` is a sort which adds an element `noRat` to the built-in sort `Rat` of the rational numbers.)

```
class NOMsender | NAMTimer : TimeInf, csmNominee : DefOid, csmLPE : DefRat,
                csmRTT : TimeInf, csmSetTime : Time .
subclass NOMsender < Sender Clock .

class NOMsenderAlone . subclass NOMsenderAlone < NOMsender .
```

The crucial rule is the following rule, specifying the handling of a CSM packet from one of the children of the sender:

```
msg csmPacket : DefNat TimeInf DefOid -> Packet . *** Usage: csmPacket(lpe, rtt, rcvr)
msg NAMPacket : Bool -> Packet . *** Usage: NAMPacket(isNominee)

vars DR DR' : DefRat .

rl [D2D3] :
  (csmPacket(DR, TI, DO') from 0 to 0')
  < 0' : NOMsenderAlone | clock : R, csmNominee : DO, csmLPE : DR',
    csmRTT : TI', csmSetTime : R', NAMTimer : TI'' >
=>
  if updateNomValues(DR, TI, DO', DO, DR', TI', R, R') then
    < 0' : NOMsenderAlone | csmNominee : DO', csmLPE : DR,
      csmRTT : TI, csmSetTime : R,
      NAMTimer : (if DO /= DO' then 7000 else TI'' fi) >
    (if DO /= DO' and DO' /= noOid *** Notify new nominee DO'
     then (NAMPacket(true) from 0' to DO') else none fi)
    (if DO /= DO' and DO /= noOid *** Notify previous nominee DO
     then (NAMPacket(false) from 0' to DO) else none fi)
  else < 0' : NOMsenderAlone | > fi .
```

The function `updateNomValues` takes the received and the stored nominee values, as well as the current time and the last time the nominee-values were updated, and returns `true` iff the nominee values should be updated. (See the Real-Time Maude specification for the definition of the function `updateNomValues`.) The `NAMPackets` to the new and old nominee receivers should be unicast without going through the links, and could in a first abstraction be seen as having no delay. Therefore, the `NAMPackets` already have the “ready-to-read” form.

### 4.12.2 Repair Server Protocol

A repair server stores the values of the receiver with the most congested path in its subtree in its attributes `csmLPE`, `csmRTT`, and `csmAddress`. The `csmTimer` attribute is used to send the current nominee values to the upstream repair server at regular intervals:

```

class NOMrepairserver | csmLPE : DefRat, csmRTT : TimeInf, csmAddress : DefOid,
                      csmSetTime : Time, csmTimer : TimeInf .
subclass NOMrepairserver < Repairserver Clock .

class NOMrepairserverAlone . subclass NOMrepairserverAlone < NOMrepairserver .

```

The crucial rule is the one handling a CSM packet from a child. If the received values indicate that the subtree has a new nominee, or that the current nominee's *rtt* and *lpe* values are changed, then the new values are sent upstream to the node's repair server:

```

rl [F2F3] :
  (csmPacket(DR, TI, DO') from 0 to 0')
  < 0' : NOMrepairserver | clock : R, repairserver : 0'', csmAddress : DO,
                          csmLPE : DR', csmRTT : TI', csmSetTime : R' >
=>
  if updateNomValues(DR, TI, DO', DO, DR', TI', R, R')
  then (< 0' : NOMrepairserver | csmAddress : DO', csmLPE : DR, csmRTT : TI,
                                csmSetTime : R, csmTimer : 7000 >
       send(csmPacket(DR, TI, DO'), 0', 0''))
  else < 0' : NOMrepairserver | > fi .

```

### 4.12.3 Receiver Protocol

The receiver updates a *sliding window* with the sequence number of the data packets it receives to estimate its loss probability. The following declarations define the interface of the sliding window module WINDOW given in our specification. (Note that `windowLPE` returns `noRat` if no elements have been added to an `initWindow`.)

```

sort Window .
op initWindow : NzNat -> Window .    *** Empty window with given max size
op size : Window -> Nat .           *** No of elements currently in window
op add : NzNat Window -> Window .   *** Adds a sequence number to a window
op windowLPE : Window -> DefRat .   *** lpe of a window

```

The `msgWindow` attribute in the following class is the sliding window for storing message numbers, and `isNominee` is a flag which is set (to `true`) when the receiver is the nominee receiver:

```

class NOMreceiver | isNominee : Bool, sourceRTT : TimeInf, msgWindow : Window,
                  csmTimer : TimeInf .
subclass NOMreceiver < Receiver .

class NOMreceiverAlone . subclass NOMreceiverAlone < NOMreceiver .

```

We use a simplified form `dataPacket(seqNo, timeStamp)` of data packets in this protocol, and treat the reception of a data packet by inserting its sequence number into the receiver's sliding window:

```

var W : Window .
rl [E2] :
  (dataPacket(NZN, R) from 0 to 0')
  < 0' : NOMreceiverAlone | msgWindow : W >
=>
  < 0' : NOMreceiverAlone | msgWindow : add(NZN, W) > .

```

A receiver sends a `csmPacket` with its current `sourceRTT` and `lpe` values to its repair server when the `csmTimer` expires. According to the informal specification, the `lpe` estimate is considered unreliable if the size of the window is less than 150, so the default value `noRat` is sent instead. However, for more convenient prototyping, we changed the size bound from 150 to 2 below<sup>11</sup>.

```

rl [E3] :
  < 0 : NOMreceiver | csmTimer : 0, msgWindow : W, repairserver : 0', sourceRTT : TI >
=>
  < 0 : NOMreceiver | csmTimer : 5000 >
  send(csmPacket(if (size(W) < 2) then noRat else windowLPE(W) fi, TI, 0), 0, 0') .

```

The `isNominee` attribute is updated according to received status in the `NAMPacket`:

```

rl [E4] :
  (NAMPacket(X) from 0' to 0)
  < 0 : NOMreceiverAlone | >
=>
  < 0 : NOMreceiverAlone | isNominee : X > .

```

### 4.13 Analyzing the NOM Component

The NOM component is supposed to find the nominee receiver, which is crucial since only the nominee receiver acknowledges the reception of data packets, and without such acknowledgments the rate control part may slow or block the sending of new data packets. Execution of earlier versions of the *combined* protocol failed to send more than one data packet, because the first packet was sent before a nominee was chosen. We therefore changed the protocol to wait some time before sending the first packet, so that a nominee would have been found. An important property the NOM protocol should satisfy is that some receiver must have its `isNominee` flag set to `true` within a reasonable amount of time. A second important property to ensure acknowledgment of each data packet is that, at any time after a nominee has been found for the first time, there should be *some* receiver with its `isNominee` flag set to `true`. A third important property is that the *correct* nominee is chosen.

To be able to analyze the specification of the NOM component in isolation, we add an *environment* object which defines *what* original data packets are received by the receivers as well as *when* these packets are received. In addition, we fix the RTT values. In the following initial state, the receiver 'f' will receive three data packets, with sequence numbers 2, 3, and 4, arriving at times 14996, 14999, and 15031 respectively.

---

<sup>11</sup>Although the `lpe` estimates are then considered less reliable, this avoids having long initial computation segments that could cause combinatorial explosion when performing formal analysis. Furthermore, the design errors we found did not depend on the specific value chosen for the size bound.

```

op NOMstate2 : Nat -> GlobalSystem .
eq NOMstate2(N) =
  ({ startNOM('a)
    < 'a : NOMsenderAlone | clock : 0, children : 'b 'c, NAMTimer : INF, csmRTT : 0,
      csmLPE : noRat, csmSetTime : 0, csmNominee : no0id >
    < 'b : NOMreceiverAlone | sourceRTT : 44, repairserver : no0id, isNominee : false,
      msgWindow : initWindow(4), csmTimer : INF >
    < 'c : NOMrepairserverAlone | clock : 0, children : 'd 'e, repairserver : no0id,
      csmLPE : noRat, csmRTT : 0, csmAddress : no0id,
      csmSetTime : 0, csmTimer : INF >
    < 'd : NOMrepairserverAlone | ... >
    < 'e : NOMreceiverAlone | sourceRTT : 94, msgWindow : initWindow(4), ... >
    < 'f : NOMreceiverAlone | sourceRTT : 118, msgWindow : initWindow(4), ... >
    < 'g : NOMreceiverAlone | sourceRTT : 118, msgWindow : initWindow(4), ... >
    < 'random : RandomNGen | seed : N >
  LINKS
  < 'env : Env | msgsFromEnv :
    dly(dataPacket(1, 1) from 'a to 'b, 5001)   dly(dataPacket(4, 1) from 'a to 'b, 5004)
    dly(dataPacket(2, 1) from 'd to 'f, 14996)  dly(dataPacket(3, 1) from 'd to 'f, 14999)
    dly(dataPacket(4, 1) from 'd to 'f, 15031)  dly(dataPacket(4, 1) from 'd to 'g, 5002)
    dly(dataPacket(5, 1) from 'd to 'g, 15000)  dly(dataPacket(6, 1) from 'd to 'g, 15004)
    dly(dataPacket(1, 1) from 'c to 'e, 5003)   dly(dataPacket(2, 1) from 'c to 'e, 5018)
    dly(dataPacket(16, 1) from 'c to 'e, 15001) > }) .

```

where *LINKS* stands for the same set of link objects given in the state *RTTstate2* above. There are no *lpe* values before time 5000 when starting from this initial state, since no data packets have been received. In that case, the receiver with the largest *rtt* value, i.e., 'f or 'g, should be the chosen nominee. The nominee is not changed by the packets arriving around time 5000, since 'f's *lpe* is still undefined and 'g's *lpe* is 3/4, while the *lpe* of 'b and 'e is 1/2. Finally, after the messages received around time 15000, receiver 'e, with its loss rate of 75% (since only packet 16 fits in the window which can store elements within an interval of length 4), should be the nominee.

We can first use Real-Time Maude's rewriting command to check the nominees at time 4500, 14500, and 20000:

```
Maude> (trew NOMstate2(1) in time <= 5000 .)
```

```
Result ClockedSystem :
```

```

{< 'a : NOMsenderAlone | csmLPE : noRat, csmNominee : 'f, csmRTT : 118, ... >
 < 'c : NOMrepairserverAlone | csmAddress : 'f, ... >
 < 'f : NOMreceiverAlone | isNominee : true, msgWindow : window(nil,4,0), ... >
 ... } in time 4764

```

```
Maude> (trew NOMstate2(1) in time <= 14500 .)
```

```
Result ClockedSystem :
```

```

{< 'a : NOMsenderAlone | csmLPE : noRat, csmNominee : 'f, csmRTT : 118, ... >
 < 'c : NOMrepairserverAlone | csmAddress : 'f, ... >
 < 'f : NOMreceiverAlone | isNominee : true, msgWindow : window(nil,4,0), ... >
 ... } in time 14490

```

```
Maude> (trew NOMstate2(1) in time <= 20000 .)
```

```

result ClockedSystem :
  {< 'a : NOMsenderAlone | csmNominee : 'e, ... >
   < 'e : NOMreceiverAlone | isNominee : true, ... >
   ... } in time 19764

```

It should also be possible to reach a state where 'g is the nominee instead of 'f within time 5000:

```

Maude> (tsearch [1] NOMstate2(1) =>*
      {< 'g : NOMreceiverAlone | isNominee : true, ATTS:AttributeSet >
       C:Configuration}
      in time <= 5000 .)

```

Solution 1

```

C:Configuration <- < 'a : NOMsenderAlone | csmNominee : 'g, ... > ... ;
TIME_ELAPSED:Time <- 490

```

There can be no other nominee than 'f and 'g before time 15000:

```

Maude> (tsearch [1] NOMstate2(1) =>*
      {< 'a : NOMsenderAlone | csmNominee : 0:Oid, ATTS:AttributeSet >
       C:Configuration} such that 0:Oid /= 'f /\ 0:Oid /= 'g
      in time <= 15000 .)

```

No solution

The receiver 'e should eventually be the nominee, but *not* before time 15000:

```

Maude> (find earliest NOMstate2(1) =>*
      {< 'e : NOMreceiverAlone | isNominee : true, ATTS:AttributeSet >
       C:Configuration} .)

```

Result:

```

{< 'e : NOMreceiverAlone | isNominee : true, ... > ... } in time 19504

```

Sooner or later 'e must be the nominee receiver in all possible behaviors:

```

Maude> (find latest NOMstate2(1) =>*
      {< 'e : NOMreceiverAlone | isNominee : true, ATTS:AttributeSet >
       C:Configuration} with no time limit .)

```

Result:

```

{< 'e : NOMreceiverAlone | isNominee : true, ... > ... } in time 19504

```

Furthermore, once 'e is the receiver it should remain so:

```

Maude> (check NOMstate2(1) |=
      {C:Configuration} untilStable
      {< 'e : NOMreceiverAlone | isNominee : true, ATTS:AttributeSet >
       C:Configuration} in time <= 30000 .)

```

Result: the property holds.

The protocol seems to find the *correct* nominees. It remains to check how much time the protocol needs to find a nominee, and that there will always be a nominee once a nominee is found. The first of these properties can be checked as follows:

```
Maude> (find latest NOMstate2(1) =>* {< 0:Oid : NOMreceiverAlone | isNominee : true,
                                     ATTS:AttributeSet >
                                     C:Configuration} with no time limit .)
```

Result:

```
{< 'f : NOMreceiverAlone | isNominee : true, ... > ... } in time 490
```

Finally we check whether there is a behavior —after some receiver has been nominated and is aware of it— in which no receiver has its `isNominee` flag set to `true` (and that some packet therefore may not be acknowledged).

```
Maude> (check NOMstate2(1) |=
        {C:Configuration} untilStable
        {< 0:Oid : NOMreceiverAlone | isNominee : true, ATTS:AttributeSet >
         C:Configuration} with no time limit .)
```

Result: the property does not hold. Counterexample:

```
{< 'a : NOMsenderAlone | csmNominee : 'e, ... >
 < 'b : NOMreceiverAlone | isNominee : false, ... >
 < 'e : NOMreceiverAlone | isNominee : false, ... >
 < 'f : NOMreceiverAlone | isNominee : false, ... >
 < 'g : NOMreceiverAlone | isNominee : false, ... >
 (NAMPacket(true) from 'a to 'e) .... } in time 19504
```

This shows one troubling scenario in which no receiver is aware of it being the nominee, and so no receiver will acknowledge a data packet received at this moment.<sup>12</sup>

(Some of the above properties can also be expressed in temporal logic. However, the temporal logic model checker failed, as expected, to terminate when checking the unlimited “until-stable” property above.) An important aspect of the specification is that once either `'f` or `'g` is found to be the nominee receiver, the nominee should not change until `'e` becomes the nominee receiver. This property cannot be expressed using Real-Time Maude’s search and until-commands, so we must use temporal logic to express it. The module

```
(tomod MODEL-CHECK-NCA is including TIMED-MODEL-CHECKER .
  protecting NCA-NOM1 .

  op nomineeExists : -> Prop .
  eq {< 0:Oid : NOMreceiverAlone | isNominee : true > C:Configuration}
    |=
    nomineeExists = true .

  op nomineeIs : Oid -> Prop .
```

---

<sup>12</sup>Because of this problem the informal protocol has since been changed to a version in which each data packet is equipped with a field which says which receiver should be nominee receiver for that packet.

```

eq {< 0:Oid : NOMreceiverAlone | isNominee : true > C:Configuration}
  |=
  nomineeIs(0:Oid) = true .

op becomingNominee : Oid -> Prop .
eq {(NAMPacket(true) from 0:Oid to 0':Oid) C:Configuration}
  |=
  becomingNominee(0':Oid) = true .
endtom)

```

defines the properties `nomineeExists`, which holds if some receiver has its `isNominee` flag set, `nomineeIs(o)`, which holds when `o`'s `isNominee` flag is set, and `becomingNominee(o)`, which holds when there is a message to receiver `o` stating that `o` is the nominee.

The following model checking command checks the whole expected behavior from the given initial state: First there is no nominee, then either `'f` or `'g` is the nominee and it stays that way until `'a` sends the NAM packet to `'e`, then it stays so until `'e` reads the NAM packet and is the nominee, and then `'e` becomes the nominee and remains the nominee:

```

Maude> (mc NOMstate2(1) |=t
  ~ nomineeExists U
    ((nomineeIs('f) \ / nomineeIs('g))
     /\
     ((nomineeIs('f) -> (nomineeIs('f) U
                          (becomingNominee('e) U ([] nomineeIs('e))))))
     /\
     (nomineeIs('g) -> (nomineeIs('g) U
                          (becomingNominee('e) U ([] nomineeIs('e))))))
  in time <= 50000 .)

Result Bool :
  true

```

The above temporal property sums up the desired “untimed” behavior of the system. We should also check the same behavior in its “timed” version: There is no nominee until either `'f` or `'g` becomes the nominee within time 500; then this nominee stays the nominee, but not *past* time 20000. In the meantime, `'e` must become the nominee, but not before time 15300, and `'e` remains the nominee ever after. Real-Time Maude also allows model checking *clocked* properties, where the properties may depend on the time elapsed to reach a certain state. In the following module, the *clocked* proposition `nomineeIsBefore(o, r)` holds for all states where `o` is the nominee receiver and where the total time elapse is less than or equal to `r`. This property, and the symmetric `nomineeIsAfter`, can be defined as follows:

```

(tomod MODEL-CHECK-CLOCKED-NOM is including MODEL-CHECK-NCA .

ops nomineeIsBefore nomineeIsAfter : Oid Time -> Prop .

eq {< 0:Oid : NOMreceiverAlone | isNominee : true > C:Configuration} in time R:Time
  |=
  nomineeIsBefore(0:Oid, R':Time) = R:Time <= R':Time .

```

```

eq {< 0:0id : NOMreceiverAlone | isNominee : true > C:Configuration} in time R:Time
|=
  nomineeIsAfter(0:0id, R':Time) = R':Time <= R:Time .
endtom)

```

We can now check whether all behaviors satisfy the expected timed behavior:

```

Maude> (mc NOMstate2(1) |=t
  ~ nomineeExists U
    ((nomineeIsBefore('f, 500) \/\ nomineeIsBefore('g, 500))
     /\
      ((nomineeIsBefore('f, 500) ->
        (nomineeIsBefore('f, 20000) U
          (becomingNominee('e) U ([] nomineeIsAfter('e, 15300))))))
      /\
        (nomineeIsBefore('g, 500) ->
          (nomineeIsBefore('g, 20000) U
            (becomingNominee('e) U ([] nomineeIsAfter('e, 15300)))))))
  in time <= 50000 .)

Result Bool :
  true

```

#### 4.14 Specification of the Rate Control Component

The rate control component aims at dynamically adjusting the sending rate of data packets based on acknowledgments from the *nominee* receiver. The protocol does not decide *when* the next packet can be sent, since the reception of an acknowledgment could alter the earliest time the next packet could be sent. Instead, the protocol is used to decide *whether* a packet can be sent at the *current* time.

The nominee receiver acknowledges each data packet that it has not seen before by sending a *congestion control message* (`ccmPacket`) with the sequence number of its *lowest outstanding* data packet (together with the timestamp and the retransmission flag of the received data packet) to the sender.

The sender part is fairly complex and we just outline the main ideas. The sender has a “window” of sequence numbers, defining the set of data packets which can be sent, starting with the sequence number of the lowest outstanding data packet at the nominee receiver (the attribute `winLowerSeq` in the class declared below), and containing `[winSize]` elements. The other limitation on sending is that a new packet should not be sent if the `maxBurstCount` value, denoting the number of original data packets sent since the last `ccmPacket` was received, is greater than a certain limit.

Much of the protocol is concerned with increasing and decreasing the `winSize` value based on acknowledgments, or lack thereof, from the nominee. Each time a `ccmPacket` is received (in normal recovery mode), the value of the `winSize` attribute increases by one until a threshold is reached, after which `winSize` increases by the much smaller amount  $\frac{1}{\lfloor \text{winSize} \rfloor}$ . The sender goes into *fast repair* mode when it has received three `ccmPackets` with the same lowest outstanding sequence number. In fast recovery mode, the number of messages that can be sent is lowered by dividing `winSize` by 2.

The main rule, which updates the various attributes upon the reception of a CCM packet from the nominee, is not given here but can be found in [18, p. 183]. (A CCM packet from a receiver which is not the current nominee is ignored.)

The packets communicated, namely `ncarcDataPacket` which is used for a data packet in the stand-alone RC component, and `ccmPacket` are declared as follows:

```
msg ncarcDataPacket : NzNat Time Bool -> LargePacket .
*** Usage: ncarcDataPacket(seqNo, timestamp, retransmissionFlag)

msg ccmPacket : Oid Nat Time Bool -> Packet .
*** Usage : ccmPacket(originator, sequenceNo, timestamp, retransmissionFlag)
```

**The Sender Protocol.** The sender object has 17 attributes used to determine the sending rate. We just list a few in the following declaration of the sender:

```
class RCsender | csmNominee : DefOid, winSize : PosRat, winLowerSeq : Nat,
                maxSeqSent : Nat, ccmTimeout : Time, ccmTimer : TimeInf, ... .
subclass RCsender < Clock Sender .
```

where `csmNominee` denotes the current nominee, `maxSeqSent` denotes the highest sequence number of packets sent, `winLowerSeq` denotes the lowest outstanding sequence number at the nominee, and `winSize` denotes the window size of sequence numbers that can be sent. The `ccmTimer` is used to detect the lack of reception of CCM packets for a period of time.

To make an execution of the stand-alone RC protocol interesting, we add a sending mechanism that attempts to send a new packet every millisecond. We also add to the sender a list which contains the times the sender has been allowed to send original data packets. These “extra” attributes are declared in the following subclass `RCsenderAlone`, and are therefore not part of the formal specification of the (combined) protocol.

```
class RCsenderAlone |
    sendDataTimer : TimeInf, *** Attempt to send new data packet upon its expiration.
    sentTimes : TimeList . *** Times of sendings of data packets.
subclass RCsenderAlone < RCsender .
```

The function

```
op sendAllowed : Nat Nat PosRat ... -> Bool .
*** Usage: sendAllowed(maxSeqSent, winLowerSeq, winSize, ...)
```

decides whether a new packet can be sent. It checks whether the next data packet, having sequence number `maxSeqSent + 1`, is in the window of sequence numbers that can be sent, and whether the `maxBurstCount` value is sufficiently low.

A new data packet is sent, if possible, by the following rule when the `sendDataTimer` expires:

```

rl [G12] :
  < 0 : RCsenderAlone | children : OS, clock : R, sendDataTimer : 0, maxSeqSent : N,
    winLowerSeq : N', winSize : WS, sentTimes : TL, ... >
=>
  if sendAllowed(N, N', WS, ...) then
    < 0 : RCsenderAlone | sendDataTimer : 1, maxSeqSent : N + 1, sentTimes : TL ++ R,
      ccmTimer : (if N' <= N then ... else INF fi), ... >
    multiSend(ncarcDataPacket(N + 1, R, false), 0, OS)
  else < 0 : RCsenderAlone | sendDataTimer : 1 > fi .

```

The `ccmTimer` expires when no CCM packet has been received for a while. The window size is reset to 1, and the lowest sequence number is set to the next sequence number to be sent, essentially starting the rate control all over, albeit with some other values (which we don't present here) changed to reflect that no packet has been acknowledged for a while.

```

rl [G13] :
  < 0 : RCsender | ccmTimer : 0, winSize : WS, winLowerSeq : N, maxSeqSent : N', ... >
=>
  < 0 : RCsender | ccmTimer : INF, winSize : 1, winLowerSeq : N' + 1, ... > .

```

**The Repair Server Protocol.** A repair server simply subcasts packets from its repair server downstream and propagates packets from its children upstream. The details can be found in the Real-Time Maude specification [18, p. 186].

**The Receiver Protocol.** The receiver stores the sequence numbers of the data packets it has seen in the attribute `msgRcvd` of sort `OrderedNzNatList` of lists of nonzero natural numbers, to be able to find its lowest outstanding sequence number.

```

class RCreceiver | isNominee : Bool .
class RCreceiverAlone | msgRcvd : OrderedNzNatList .
subclass RCreceiverAlone < RCreceiver .

```

The nominee receiver acknowledges a data packet it has not seen before by sending a `ccmPacket` upstream to the sender by way of the intermediate repair servers:

```

var NNL : OrderedNzNatList .

rl [H1] :
  (ncarcDataPacket(NZN, R, X) from 0 to 0')
  < 0' : RCreceiverAlone | isNominee : Y, msgRcvd : NNL >
=>
  < 0' : RCreceiverAlone | msgRcvd : add(NZN, NNL) >
  (if (Y and (not (NZN inList NNL)))
    then send(ccmPacket(0', lowestOutstanding(add(NZN, NNL)), R, X), 0', 0)
    else none fi) .

```

## 4.15 Analyzing the Rate Control Component

The rate control protocol is tested by attempting to send a new data packet every millisecond, and by recording in the state the time stamp of each new data packet sent, and by recording the messages lost. The list of sending times and packet losses could then be inspected to get a feeling for the sending rate.

What results do we expect from executing the specification? In the beginning, the sending window is (1), and packet 1 is sent. When the acknowledgment of packet 1 is received, both the lowest window number, and the size of the window, are increased. That is, the window is now (2, 3), and both these data packets are sent. When the acknowledgment of packet 2 is received, the window is (3, 4, 5), and packets 4 and 5 can be sent. The sequence of sending and responses can be given as

message sent	1	2, 3	4, 5	6, 7	8, 9	10, 11	...
response		1	2	3	4	5	...

The sending rate would grow exponentially were it not for congestion delays, since acknowledgments of packets 2 and 3 would come at the same time, allowing packets 4 to 7 to be sent, and upon reception of the acknowledgment for these latter packets, the sender could send packets 8 to 15, and so on. The exponential increase should flatten out when the window size has reached a certain threshold, after which the window size should only increase by its inverse. However, increasing the sending frequency could result in packets getting lost. In the full protocol, the lost packets would be repaired, but we did not add any repair service to the stand-alone RC protocol. The system should therefore get stuck when the first data packet is lost, as the receiver will require it but cannot get it. Eventually, the `ccmTimer` should expire, which will reset `winLowerSeq` to next packet to be sent, but will also reset the window size to 1, starting the slow sending once again.

Despite these expectations, a simulation of one behavior of the protocol from a state with one sender 'a, one repair server 'b, and one receiver 'c gives the following result:

```
Maude> (trew RCstate1 in time <= 6800 .)

Result ClockedSystem :
{< 'a : RCsenderAlone | clock : 6800, children : 'b, csmNominee : 'c, ccmTimer : 55,
  backOffFactor : 1, nextNoToSend : 31, sendDataTimer : 1,
  maxSeqSent : 30, winLowerSeq : 28, winSize : 3, sentTimes :
    (6000 ++ 6067 ++ 6068 ++ 6134 ++ 6201 ++ 6202 ++ 6268 ++ 6269 ++ 6280 ++
    6335 ++ 6347 ++ 6359 ++ 6402 ++ 6414 ++ 6426 ++ 6469 ++ 6481 ++ 6493 ++
    6536 ++ 6548 ++ 6560 ++ 6603 ++ 6615 ++ 6627 ++ 6670 ++ 6682 ++ 6694 ++
    6737 ++ 6749 ++ 6761), ... >
< 'b : RCrepairserverAlone | children : 'c, repairserver : 'a >
< 'c : RCreceiverAlone | isNominee : true, msgRcvd :
  (1 ++ 2 ++ 3 ++ 4 ++ 5 ++ 6 ++ 7 ++ 8 ++ 9 ++ 10 ++ 11 ++ 12 ++ 13 ++ 14
  ++ 15 ++ 16 ++ 17 ++ 18 ++ 19 ++ 20 ++ 21 ++ 22 ++ 23 ++ 24 ++ 25 ++ 26
  ++ 27 ++ 28 ++ 29), ... >
... } in time 6800
```

The sending rate does not increase exponentially—it is rather the case that four data packets can be sent every 67 milliseconds. After receiving `ccmPacket` requesting packet three, both data packet number four and five should be allowed to be sent according to our analysis above. However, only packet number four is sent, at time 6134. Why cannot packet five be sent at time 6135?

After using Real-Time Maude’s tracing facilities to analyze the rewrite sequence leading to this state, it turns out that the problem is that the `ccmTimer` is set to the RTT value when a new packet is sent, and it expires just when the `ccmPacket` for number three arrives. (This is not unnatural, since the timer is set to exactly the RTT.) The system may choose to send either packet 4 (and then 5) upon the reception of the `ccmPacket`, and then handle the timer, or it could handle the timer first. In the above execution, the timer was treated first, so that the window was “reinitialized” (with `winSize` set to one) and then packet four was sent, and then this process essentially repeats itself over and over. One could also notice that no packets are lost since the sending rate never increases. If this is not the desired behavior of the protocol, one may want to change the initialization of the `ccmTimer` in Use Case G12, especially since `backOffFactor` is reset to 1 in Use Case G2 each time a `ccmPacket` is received.

To further analyze the protocol, we could search for all states reachable in exactly time 6700:

```
Maude> (tsearch RCstate1 =>* {C:Configuration}
      in time-interval between >= 6700 and <= 6700 .)
```

This search fails to terminate within reasonable time, due to the high degree of nondeterminism in the specification. Instead, we “split” the search above into searches for states reachable in time 6200, and from there to states reachable in time 6500, and from there to states reachable in time 6700. Out of the many thousand solutions to that chain of searches, one result is the state

```
{< 'a : RCsenderAlone | sentTimes :
  (6000 ++ 6066 ++ 6067 ++ 6132 ++ 6133 ++ 6144 ++ 6145 ++ 6198 ++ 6199 ++
  6210 ++ 6211 ++ 6222 ++ 6223 ++ 6234 ++ 6235 ++ 6264 ++ 6265 ++ 6276 ++
  6277 ++ 6288 ++ 6289 ++ 6300 ++ 6301 ++ 6312 ++ 6313 ++ 6324 ++ 6325 ++
  6336 ++ 6337 ++ 6348 ++ 6349 ++ 6360 ++ 6361 ++ 6372 ++ 6373 ++ 6384 ++
  6385 ++ 6548 ++ 6614 ++ 6615 ++ 6680 ++ 6681 ++ 6692 ++ 6693), ... >
< 'c : RCreceiverAlone | isNominee : true, msgRcvd :
  (1 ++ 2 ++ 3 ++ 4 ++ 5 ++ 6 ++ 7 ++ 8 ++ 9 ++ 10 ++ 11 ++ 12 ++ 13 ++ 14
  ++ 15 ++ 16 ++ 17 ++ 18 ++ 20 ++ 22 ++ 24 ++ 26 ++ 28 ++ 30 ++ 32 ++ 34
  ++ 36 ++ 38 ++ 39 ++ 40)>
... } in time 6700
```

This result is more in line with the expected behavior of the protocol. The first packets are sent pairwise with (exponentially?) increasing frequency. The first packet lost is number 19, so packet 38 cannot be sent. The system therefore gets stuck after packet 37 is sent (time 6385). After quite a long time (163 milliseconds) the `ccmTimer` expires and reinitializes the window and the sending can continue, with slower sending rate. However, packet 19 is still not repaired, and will cause further problems.

To summarize the analysis effort of this admittedly difficult protocol component, we see that its behavior is highly nondeterministic – probably more so than intended. Some of the behaviors seem undesirable and some are more in line with our expectations. Real-Time Maude’s tracing capabilities allowed us to trace the spurious behavior and to suggest changes in the original protocol to remedy the problem.

## 4.16 Specification of the Repair Service Component

The repair service component of the AER/NCA protocol suite specifies a system which receives variable-sized data blocks from a sender application, places the data in a number of data packets, and is responsible for transmitting *all* data packets to a multicast group of receiver applications, so that the original data blocks can be recovered. The overall goal is to ensure reliability while transmitting as few packets as possible.

The component exhibits the following behavior: The data packets are multicast by the sender and are intercepted by the repair servers, which cache the packets before subcasting them downstream. Each data packet has a sequence number making it possible to detect that packets have been lost. Upon detection of the loss of data packets, a receiver or repair server waits for some time before requesting a repair (a retransmission of the lost packet) for each missing packet from its upstream repair server. If a missing packet is not received within a certain time, the receiver or repair server retransmits its repair request. A sender or a repair server which has the missing packet in its cache treats a repair request by trying to estimate whether the packet is currently being repaired, and, if that is not the case, then it subcasts (a repair for) that packet. If the repair server does not have the missing packet, then the repair request should be propagated upstream. A subcast repair packet is also intercepted by the next downstream repair servers and cached but only forwarded further down the tree if there is a pending repair request for the repair.

In addition, the sender regularly multicasts *source path message* (SPM) packets with the sequence number of the last (original) data packet transmitted. These (SPM) packets are used to discover packet losses and to set up the reverse paths from the each receiver to the sender through repair servers (that is, to set the correct value of the `repairserver` attribute).

Again, we present only a small part of the Real-Time Maude specification.

### 4.16.1 Sender Protocol

To model the interaction with the application layer of the protocol, the sender receives data blocks of the form `dataBlock(content, sender, size,  $\Delta$ )` from the application layer, where *content* is an identifier of the content of the data, *sender* is the sender object (useful in case there are more than one multicast group), *size* is the number of data packets needed to represent the data block, and  $\Delta$  is the interval between each attempt at sending a new (original) data packet<sup>13</sup>.

The data packets are represented as terms of sort `LargePacket` of the form

$$\text{dataPacket}(\text{content}, \text{seqNo}, \text{timeStamp}, \text{firstSeqFlag}, \text{segFlag}, \text{retransFlag}),$$

where *content* is the content, or *payload*, of the data packet<sup>14</sup>, *seqNo* is the sequence number of the data packet, and is one plus the number of original data packets already sent in the current execution of the protocol, *timeStamp* is the time the packet is sent, *firstSeqFlag* is true iff the packet is the first data packet of a data block, *segFlag* is true iff the data packet is neither the first nor the last packet of a data block, and *retransFlag* is true iff the packet is a repair packet.

---

<sup>13</sup> $\Delta$  is therefore the actual sending rate when the repair service component is executed independently of the rate control component. If  $\Delta$  is 1, then sending is attempted every time unit.

<sup>14</sup>In our abstraction, the payload is the identifier of the data block being transmitted

The sender should not necessarily repair a packet when it receives a repair request (in the form of a `NAKPacket`). The reason is that a repair for the packet may already have been subcast to its children as a response to another repair request (and repairing the packet again would be unnecessary and inefficient). However, if a repair is lost, the packet needs to be repaired again. To estimate whether a repair should be performed, the sender stores for each data packet a *NAK count*, which estimates how many “rounds” the packet has been repaired. (The NAK count is also influenced by a repair performed by a repair server further down the tree.)

The sender class is declared as follows. The sender stores the parts of the data blocks it has yet to send in an attribute `unsentDataPackets`. This attribute is a *list* (where list concatenation is denoted by juxtaposition) of elements of the form `data(content, m, size, Δ)`, where  $m - 1$  is the number of data packets already multicast from the data block, and *size* and  $\Delta$  correspond to the same parameters of the data block. The `nextSeq` attribute denotes the sequence number of the next packet to be sent, `SPMTimer` a timer used to force the sending of `SPMPackets` every 4000 milliseconds., `reTransBuf` is the retransmission buffer containing already transmitted data packets together with their corresponding NAK count in messages of the form `MsgAndNak(data packet, NAKcount)`. Finally, `sendDataTimer` is a timer upon whose expiration the sender sends a new data packet:

```
class RSsender | nextSeq : NzNat, SPMTimer : TimeInf, reTransBuf : MsgConfiguration,
                unsentDataPackets : DataBuffer, sendDataTimer : TimeInf .
subclass RSsender < Sender Clock .

class RSsenderAlone .   subclass RSsenderAlone < RSsender .
```

We omit here the rules modeling the treatment of data blocks from the application layer and the initialization phase. The following rule models the multicast of `SPMPackets` with the sequence number of the last data packet transmitted when the `SPMTimer` expires:

```
rl [A2] :
  < 0 : RSsender | nextSeq : s N, children : OS, SPMTimer : 0 >
=>
  < 0 : RSsender | SPMTimer : 4000 >
  multiSend(SPMPacket(N), 0, OS) .
```

The data packet with sequence number `nextSeq` is multicast to the sender’s children when `sendDataTimer` expires:

```
var DB : DataBuffer .

rl [A3send] :
  < 0 : RSsenderAlone | children : OS, clock : R, nextSeq : NZN'',
                      unsentDataPackets : data(0', NZN, NZN', R') DB,
                      reTransBuf : RTB, sendDataTimer : 0 >
=>
  < 0 : RSsenderAlone | unsentDataPackets : removeFirst(data(0', NZN, NZN', R') DB),
                      reTransBuf :
                        RTB MsgAndNAK(dataPacket(0', NZN'', R, NZN == 1,
                                                  NZN /= 1 and NZN /= NZN', false), 0),
                      nextSeq : NZN'' + 1,
                      sendDataTimer : sendRate(removeFirst(data(0', NZN, NZN', R') DB)) >
  multiSend(dataPacket(0', NZN'', R, NZN == 1, NZN /= 1 and NZN /= NZN', false), 0, OS) .
```

(The function `removeFirst` removes the first “data packet” from the data buffer, and `sendRate` finds the time between consecutive sending attempts.)

When a child has not received a data packet that it should have received, or needs to see the data packet again (to repair a packet that is no longer in its cache), it requests a retransmission of the data packet by sending a `NAKPacket` message containing the sequence number of the missing packet as well as its NAK count. If the NAK count of the request (denoted by  $N$  in the rule below) is less than or equal to the sender’s NAK count for the packet (denoted by  $N'$ ), then this repair is considered underway and is not performed:

```

msg NAKPacket : NzNat Nat Bool -> Packet .
*** Usage: NAKPacket(seqNo, NAKcount, fastRepairFlag)

rl [A4] :
  (NAKPacket(NZN, N, X') from 0'' to 0)
  < 0 : RSsender | children : OS, reTransBuf :
    (MC MsgAndNAK(dataPacket(0', NZN, R, X, Y, Z), N')) >
=>
  < 0 : RSsender | reTransBuf :
    (MC MsgAndNAK(dataPacket(0', NZN, R, X, Y,
      if N' < N then true else Z fi),
      max(N, N')) >
  (if N' < N then multiSend(dataPacket(0', NZN, R, X, Y, true), 0, OS) else none fi) .

```

#### 4.16.2 Receiver Protocol

A receiver receives data packets and forwards them to a corresponding receiver *application* in increasing order of their sequence numbers. Received data packets that cannot be forwarded to the application, because some data packets with lower sequence numbers are missing, are stored in the `dataBuffer` attribute.

When the receiver detects the loss of a data packet, it suppresses its repair request for a short time (in case some of its “siblings” or its repair server have also detected the loss) before sending a NAK-request for the lost packet to its repair server. The receiver retransmits its request for the missing data packet if it does not receive a response to the repair request within a reasonable amount of time, which depends on the round trip time to the source (since a repair request may be propagated all the way up to the source in the worst case).

The receiver class in the RS component is declared as follows. `isNominee` and `fastRepairFlag` are flags which are set to `true` when the receiver is the nominee receiver; `readNextSeq` is the sequence number of the first missing data packet; `smoothedRTT` is the smoothed sender to receiver RTT estimate; `suppressT0` is the time before sending a NAK packet upon detection of a loss; `retransT0` is the time between consecutive repair requests for the same packet; `dataBuffer` contains the packets received but not forwarded to the receiver application; and `dataInfo` is the *repair information*, where the receiver stores, for each *missing* data packet, the information about the recovery attempts for the missing data packets in a term

$$\text{info}(\text{seqNo}, \text{supprTimer}, \text{retransTimer}, \text{NAKcount}),$$

where `seqNo` is the sequence number of the data packet, `supprTimer` is the value of the suppression timer for the data packet, `retransTimer` is the value of the retransmission timer of the data packet,

and  $NAKcount$  is the NAK count of the data packet, denoting how many times a repair for the data packet has been attempted. Elements of a sort `DataInfo` are multisets of `info` terms, where multiset union is denoted by an associative and commutative juxtaposition operator.

```

class RSreceiver | isNominee : Bool, fastRepairFlag : Bool, readNextSeq : NzNat,
                  smoothedRTT : Int, rttVariance : Int, suppressTO : Time,
                  retransTO : Time, dataBuffer : MsgConfiguration,
                  dataInfo : DataInfo, SPMread : Bool .
subclass RSreceiver < Receiver .

class RSreceiverAlone . subclass RSreceiverAlone < RSreceiver .

```

We start our selection of crucial rules with the rule that handles a received data packet which has not been seen before. That is, the packet has not been delivered to the receiver application (since the sequence number `NZN` is greater than or equal to the `readNextSeq` value `NZN'`), and it is not in the `dataBuffer` of the receiver:

```

vars DI DI' : DataInfo .
var LAST-SEQNO-DELIVERED-TO-APP : Nat .

crl [B3b] :
  (dataPacket(Q, NZN, R, X, Y, Z) from 0 to 0')
  < 0'' : RandomNGen | seed : N >
  < 0' : RSreceiverAlone | fastRepairFlag : X' readNextSeq : NZN', suppressTO : R',
                        dataBuffer : MC, dataInfo : DI >
=>
  < 0'' : RandomNGen | seed : (if X' then N else ... fi) >
  < 0' : RSreceiverAlone | readNextSeq : max(NZN', LAST-SEQNO-DELIVERED-TO-APP + 1),
                        dataBuffer : rmDelivered(MC dataPacket(Q, NZN, R, X, Y, Z),
                                                NZN'),
                        dataInfo : remove(updateDI(cancelTimers(DI, NZN),
                                                    X', max(NZN', highestSeqNo(MC) + 1,
                                                    highestNzNAKSeqNo(DI) + 1),
                                                    NZN - 1, R', N),
                                          LAST-SEQNO-DELIVERED-TO-APP) >
  toApp((dataPacket(Q, NZN, R, X, Y, Z) MC), NZN', 0')
  if NZN' <= NZN and not (NZN seqNoIn MC)
  /\ LAST-SEQNO-DELIVERED-TO-APP :=
     lastSeqNo(toApp(MC dataPacket(Q, NZN, R, X, Y, Z), NZN')) .

```

If the received data packet has sequence number `NZN`, then all data packets with sequence numbers less than `NZN` which have not been received must be considered missing. Upon detection of missing packets for which loss recovery has not yet started, the receiver initiates loss recovery by adding to the `dataInfo` attribute an element `info( $n'$ ,  $r$ , INF, 1)` for each missing packet  $n'$ . The suppression time,  $r$ , is 0 when the receiver is in fast repair mode, and a random value in the interval 0 to 1.5 times the value of `suppressTO` otherwise (all this is done by the function `updateDI`; the function `remove` just removes information about those packets which are delivered to the receiver *application*).

Furthermore, the message buffer must be updated when a new data packet is received, and if the received data packet is the first missing packet, then this, and the following packets until the next missing data packet, must be sent to the receiver application. The message buffer and the `readNextSeq`

attribute must be updated accordingly. The test `NZN seqNoIn MC` holds true if in the message configuration `MC` there is a data packet with sequence number `NZN`. The function `toApp` yields the packets to transmit to the receiver application, and `lastSeqNo` computes the highest sequence number in such a transmission. `updateDI(repair_info, fastRepairFlag, lowSeqNo, highSeqNo, suppressTO, randomSeed)` creates a new repair state for each data packet between `lowSeqNo` and `highSeqNo` and is defined as follows:

```

eq updateDI(DI, X, N, N', R, N'') =
  if N' < N then DI
  else (info(N, suppr(X, R, N''), INF, 1)
        updateDI(removeSeqNo(DI, N), X, N + 1, N', R,
                  if X then N'' else random(N'') fi))
  fi .

```

where `suppr(fastRepairFlag, suppressTO, randomSeed)` finds the initial value for the suppression timer. `cancelTimers` turns off the suppression and retransmission timers for the received data packet, in case they were set.

The treatment of repairs detected by receiving a `SPMPacket` is similar.

In the following rule, a repair request for a missing data packet with sequence number `NZN'` is sent upstream when the suppression timer for the packet expires, and the retransmission timer is started. The receiver gives up if it has unsuccessfully requested a repair 48 times:

```

rl [B5] :
  < 0 : RSreceiver | readNextSeq : NZN, fastRepairFlag : X,
                    dataBuffer : MC, repairserver : 0',
                    dataInfo : (info(NZN', 0, TI, N) DI), retransTO : R >
=>
  if (NZN' seqNoIn MC) or (NZN' < NZN)
  then < 0 : RSreceiver | dataInfo : (info(NZN', INF, TI, N) DI) >
  else (if 48 < N then
        ERROR("Use case B5, too high NAK count for RSreceiver",
              (dataInfo : (info(NZN', 0, TI, N) DI)))
        else < 0 : RSreceiver | dataInfo : (info(NZN', INF, R, N) DI) >
        send(NAKPacket(NZN', N, X), 0, 0')
        fi)
  fi .

op ERROR : String AttributeSet -> Configuration .

```

When the retransmission timer for a (lost) data packet expires, a new repair session is initiated by starting the suppression timer and by increasing the NAK count for the packet.

The final rule we show here is the treatment of a NAK packet for a packet for which there is no NAK state. The NAK packet indicates that a repair is underway upstream, and that the receiver should not request a repair on its own. Therefore, a new NAK state is created for the packet, with the retransmission timer set (so that a repair is requested later if the “promised” repair is not successful):

```

crl [B7b] :
  (NAKPacket(NZN, N, X) from 0 to 0')
  < 0' : RReceiver | readNextSeq : NZN', retransTO : R, dataInfo : DI >
=>
  < 0' : RReceiver | dataInfo :
    (if NZN' <= NZN then (info(NZN, INF, R, N) DI) else DI fi) >
  if not (NZN in DI) .

```

### 4.16.3 Repair Server Protocol

A repair server has a (bounded) cache in which it stores some of the data packets it has seen, as well as information about the repairs it is involved in. When a repair server discovers that some data packets are missing, it tries to repair lost data packets in the same way that a receiver does. In addition, a repair server sends `NAK` packets to its children notifying them that a repair process has been initiated.

Original data packets are stored in the cache and are subcast downstream. Repair data packets are stored and are subcast downstream only if a `NAK` is pending for the packet. A repair request from a child is treated by subcasting the requested data packet if it is in the cache (and the packet is not considered to be under reparation), and by starting its own repair procedure for the packet otherwise.

A cache is represented by a value `bBuf(msgList, size, bound)` of sort `BoundedBuffer`, where the data packets are stored in the message list `msgList`. Our specification is based on the (unstated but reasonable) assumption that data packets are stored in the buffer in their order of reception, and not in order of increasing sequence numbers.

The repair server stores information about repairs in terms of the form

```

nakState(seqNo, NAKcount, downNAKcount, NAKpending, supprTimer, retransTimer),

```

where `seqNo` is the sequence number of the data packet, `NAKcount` is the upstream `NAK` count (which roughly corresponds to the number of upstream repair requests), `downNAKcount` is the downstream `NAK` count used in communication with the children, `NAKpending` is true when a child is waiting for the repair of the packet, and `supprTimer` and `retransTimer` are timer values which are used for upstream repair requests in the same way that they were used in the receiver protocol. The union of such `NAK` states is given by juxtaposition.

The repair server classes for the repair service part of the protocol are declared as follows, where `maxSeqRcvd` this time denotes the highest sequence number of any received data packet:

```

class RSrepairserver | SPMread : Bool, maxSeqRcvd : Nat, smoothedRTT : Int,
  rttVariance : Int, retransTO : Time, suppressTO : Time,
  SPMWaitTimer : TimeInf, NAKStates : NAKstates,
  dataBuffer : BoundedBuffer .
subclass RSrepairserver < Repairserver .

class RSrepairserverAlone . subclass RSrepairserverAlone < RSrepairserver .

```

The following rule treats the reception of a data packet which is not stored in the cache (the condition `not (NZN seqNoIn BB)` below). The packet is added to the `dataBuffer` cache of the

repair server, and is subcast to its children if it is an original packet (the condition `not Z` below), or if a NAK is pending for the sequence number. In addition, error recovery must be initiated for each packet that is detected to be lost with the reception of this data packet (namely, those packets with sequence number lower than the just received packet which have not been received). For each of these, a new NAK state is added by the function `updateNAKstates`, and a `NAKPacket` is sent to all its children by the function `newNAKpackets`:

```

var NS : NAKstates .   var BB : BoundedBuffer .

cr1 [C4b] :
  (dataPacket(Q, NZN, R, X, Y, Z) from 0' to 0'')
  < 0''' : RandomNGen | seed : N >
  < 0'' : RSrepairserver | children : OS, maxSeqRcvd : N', dataBuffer : BB,
    NAKStates : NS, suppressTO : R' >
=>
  < 0'' : RSrepairserver | maxSeqRcvd : max(NZN, N'),
    dataBuffer : add(BB, dataPacket(Q, NZN, R, X, Y, Z)),
    NAKStates :
      (if (N' == 0) then cancelTimersAndSetFalse(NZN, NS)
       else updateNAKstates(
         cancelTimersAndSetFalse(NZN, NS),
         N', NZN - 1, N, R)
       fi) >
  < 0''' : RandomNGen | seed : repeatRandom(N, noOfNewNAKs(NS, N', NZN - 1)) >
  (if ((not Z) or NAKpending(NS, NZN)) then
    multiSend(dataPacket(Q, NZN, R, X, Y, Z), 0'', OS) else none fi)
  newNAKpackets(NS, N', NZN - 1, 0'', OS)
  if not (NZN seqNoIn BB) .

```

In the above rule, `newNAKpackets(ns, n, n', q, os)` sends `NAKpackets` to the children *os* from *q* for the sequence numbers in the interval  $n + 1$  to  $n'$  for which there is no NAK state in *ns*.

`updateNAKstates(ns, n, n', n'', r)` updates the NAK state *ns* with new NAK states for sequence numbers between  $n + 1$  and  $n'$  which are not already in *ns*.  $n = 0$  leaves *ns* unchanged. *n''* is the seed to be used by the random number generator, and *r* is the value of `suppressTO`. Finally, `NAKpending(nakStates, seqNo)` holds when a NAK state exists and a NAK is pending for packet *seqNo*.

A `NAKPacket` with NAK count *N* for sequence number *NZN* received from the *upstream* repair server indicates that repair number *N* (for packet *NZN*) is attempted upstream, and should be forwarded to the children if their NAK count is estimated to be smaller than *N* or the NAK state does not exist. Furthermore, the retransmission timer should be reset for round  $N + 1$ . The `NAKPacket` comes from the upstream repair server of this object if the sender of the packet has the same value as the object's `repairserver` attribute:

```

*** Packet number NZN not in cache; update and forward if higher NAK count received:
cr1 [C7b] :
  (NAKPacket(NZN, N, X) from 0 to 0')
  < 0' : RSrepairserver | repairserver : 0, dataBuffer : BB, retransTO : R,
    NAKStates : (nakState(NZN, N', N'', Y, TI, TI') NS),
    children : OS >

```

```

=>
< 0' : RSrepairserver | NAKStates : (nakState(NZN, max(N, N'),
                                     if N' < N then N else N'' fi,
                                     Y, INF, R) NS) >
(if N' < N then multiSend(NAKPacket(NZN, N, false), 0', OS) else none fi)
if not (NZN seqNoIn BB) .

*** Packet number NZN not in cache, and no NAK state for packet NZN:
cr1 [C7c] :
(NAKPacket(NZN, N, X) from 0 to 0')
< 0' : RSrepairserver | repairserver : 0, dataBuffer : BB, retransTO : R,
                       children : OS, NAKStates : NS >

=>
< 0' : RSrepairserver | NAKStates : (nakState(NZN, N, N, true, INF, R) NS) >
multiSend(NAKPacket(NZN, N, false), 0', OS)
if not (NZN seqNoIn BB) /\ not (NZN in NS) .

```

The following rules C8c and C8e are among the five rules treating a `NAKPacket` from a child (that the packet comes from downstream is seen by the sender `0` being an element of the object's `children` attribute). Essentially, if the data packet exists in the cache, it is subcast to the children. Otherwise, if the received NAK count is greater than the object's NAK count, it must start a new repair, and if the received NAK count is greater than the downstream NAK count, the other children should be aware of this current repair request. Note that if the `NAKPacket` has the fast repair flag (`X`) set, the possible upstream repair should be undertaken with no delay. We show first the rule which treats the case when the data packet is not in the cache, but where a NAK state exists for the packet:

```

*** Data packet NZN not in buffer, NAK state exists for NZN:
cr1 [C8c] :
(NAKPacket(NZN, N, X) from 0 to 0')
< 0' : RSrepairserver | repairserver : 0'', children : 0 OS,
                       NAKStates : (nakState(NZN, N', N'', Y, TI, TI') NS),
                       retransTO : R, dataBuffer : BB >

=>
(if N' < N then
  < 0' : RSrepairserver | NAKStates : (nakState(NZN, N, N, true,
                                               if X then INF else TI fi,
                                               if X then R else TI' fi) NS) >
  multiSend(NAKPacket(NZN, N, false), 0', 0 OS)
  (if X then send(NAKPacket(NZN, N, true), 0', 0'') else none fi)
else
  (if N'' < N then
    < 0' : RSrepairserver | NAKStates : (nakState(NZN, N', N, Y, TI, TI') NS) >
    multiSend(NAKPacket(NZN, N, false), 0', 0 OS)
    else < 0' : RSrepairserver | > fi)
  fi)
if not (NZN seqNoIn BB) .

```

The following rule shows the treatment of a `NAKPacket` from downstream when the requested packet `NZN` is in the repair server's cache, and where the repair server has a NAK state for the packet. If the NAK count `N` of the request is greater than the NAK count `N'` in the NAK state, then the packet is repaired; otherwise the repair request is ignored:

```

rl [C8e] :
  (NAKPacket(NZN, N, X) from 0 to 0')
  < 0' : RSrepairserver | children : 0 OS,
    dataBuffer :
      bBuf(ML + dataPacket(Q, NZN, R, Y, Z, XX) + ML', N'', NZN'),
      NAKStates : (nakState(NZN, N', N'', X', TI, TI') NS) >
=>
  if N' < N then
    < 0' : RSrepairserver | NAKStates : (nakState(NZN, N, N, false, TI, TI') NS) >
    multiSend(dataPacket(Q, NZN, R, Y, Z, true), 0', 0 OS)
  else < 0' : RSrepairserver | > fi .

```

The repair of missing packets is similar to that for receivers. A repair is attempted when a suppression timer expires for a packet with a NAK pending (since a repair is not needed if none of the repair server's children misses the packet), and is re-attempted when the retransmission timer expires for such a packet, in which case the NAK packet is also subcast downstream.

#### 4.16.4 Behavior in Time

The function `delta` increases the clocks and decreases the timers with the elapse of time, and the function `mte` makes sure that time progress stops whenever a timer expires. These functions must be extended to accommodate multiple timers in the `dataInfo` attribute of the receivers and the `NAKStates` attribute of the repair servers. We show their definitions for receivers:

```

eq mte(< 0 : RSreceiverAlone | dataInfo : DI >) = mte(DI) .
op mte : DataInfo -> TimeInf .
eq mte((none).DataInfo) = INF .
ceq mte(DI DI') = min(mte(DI), mte(DI')) if DI /= none /\ DI' /= none .
eq mte(info(NZN, TI, TI', N)) = min(TI, TI') .

eq delta(< 0 : RSreceiverAlone | dataInfo : DI >, R) =
  < 0 : RSreceiverAlone | dataInfo : delta(DI, R) > .
op delta : DataInfo Time -> DataInfo .
eq delta((none).DataInfo, R) = (none).DataInfo .
ceq delta(DI DI', R) = delta(DI, R) delta(DI', R) if DI /= none /\ DI' /= none .
eq delta(info(NZN, TI, TI', N), R) = info(NZN, TI monus R, TI' monus R, N) .

```

#### 4.16.5 The Application Layer

The repair service component and the combined protocol receive data blocks from a sender in the *application layer* of the protocol, and should relay packets (in order) to the receivers in the application layer. For simulation purposes, we have defined a simplistic model of the sender and the receivers at the application level. In that model, the sender application is an object of class `SenderApplication`, which stores a list of data blocks, with delays indicating the times each data block should be transmitted to the protocol layer, and sends the data blocks to the sender object. Each application-level receiver object of class `RcvrApplication` stores the concatenation of the data packets it has received from its associated receiver in the protocol.

## 4.17 Prototyping the Repair Service Component

To execute the repair service protocol we define an initial state `RSstate`, in which the sender application object wants to use the protocol to multicast data blocks comprising 21 data packets to the receiver applications. Rewriting this initial state should have led to a state where all receiver applications had received all packets. Instead, the execution gave the following result:

```
Maude> (trew RSstate2(113) in time <= 20000 .)

result ClockedSystem :
  {ERROR("Use case B5, too high NAK count for RSreceiver",
        dataInfo : (info(18,0,INF,49) info(19,INF,INF,1) info(20,INF,INF,1)))
  ... } in time 17559
```

By using Real-Time Maude's tracing facilities to trace the execution leading to the `ERROR`-state, we could easily find the errors in the formal and informal specifications. The problem is that when a repair server has repaired a lost packet, and the repair is lost as well, then the repair server will not try to repair the packet again if the packet is no longer in its cache, thinking that it has already repaired the packet. (In particular, the fault is in rule `C8c` above, where a `NAKPacket` is treated. If the fast repair flag (`X` in the rule) is not set, and both timers for the `nakState` are turned off, then they will not be turned on in this rule (and therefore no repair will be initiated by these timers), and neither will a repair request to the upstream repair server be sent. The execution shows that we can indeed arrive at a situation where both timers are turned off.)

In another test configuration, 35 data packets should be sent from 'a' to 'c' via 'b'. The sending rate is one new packet every 5 milliseconds, and the propagation delay of 100 milliseconds in the link from 'a' to 'b' ensures that many packets will be lost along this link with bound 10. The other link is much faster and should not lose many packets. An execution gives the following result:

```
Maude> (trew RSstate3(113) in time < 400000 .)

Result ClockedSystem :
  {< 'a : RSsenderAlone | unsentDataPackets : nil, ... >
   < 'b : RSrepairserverAlone | children : 'c, repairserver : 'a,
     dataBuffer : bBuf(dataPacket('first, 10, 5145, false, true, true) +
                        dataPacket('first, 29, 5240, false, true, true) +
                        dataPacket('first, 20, 5195, false, true, true) +
                        dataPacket('first, 35, 5270, false, false, true) +
                        dataPacket('first, 34, 5265, false, true, true) +
                        dataPacket('first, 33, 5260, false, true, true) +
                        dataPacket('first, 31, 5250, false, true, true) +
                        dataPacket('first, 22, 5205, false, true, true) +
                        dataPacket('first, 26, 5225, false, true, true) +
                        dataPacket('first, 21, 5200, false, true, true),10,10),
     NAKStates : (nakState(10, 1, 1, false, INF, INF)
                  ...
                  nakState(22, 1, 1, false, INF, INF)
                  nakState(24, 2, 2, false, INF, INF)
                  nakState(25, 1, 1, false, INF, INF)
                  ... ), ... >
```

```

< 'c : RSreceiverAlone | repairserver : 'b, isNominee : true,
  dataInfo : (info(24, INF, INF, 2)
              info(25, INF, INF, 1)
              ...
              info(35, INF, INF, 1)),
  dataBuffer : (dataPacket('first, 24, 5215, false, true, true)
                dataPacket('first, 25, 5220, false, true, true)
                ...
                dataPacket('first, 35, 5270, false, false, true)),
  retransT0 : 420, suppressT0 : 4, readNextSeq : 23,
  fastRepairFlag : true, ... >
... } in time 400000

```

After time 400000 it would be expected that the receiver application would have received the 35 data packets the sender application wanted to send to the receiver application using the protocol. Instead, only 22 packets have been relayed to the receiver application.

This execution reveals another serious problem: Packet 23 is missing at both the receiver and the repair server. Furthermore, neither of them has any repair state for this packet, which is alarming since all “holes” should be repaired when a received SPM packet or data packet indicates that some packets are lost. The tracing analysis led to the following understanding of the undesired behavior:

'a sends data packets to 'b. The first 9 data packets enter the link, which then becomes full. However, 'a continues to send packets, so the packets 10–22 are lost, and 23 is the next packet in the link which is not lost. When 'b reads packet 23 in rule C4b, it discovers the holes 10–22, and sends NAK packets for these to 'c with NAK count 1. At the same time, it also subcasts packet 23 to 'c. This makes 14 messages sent from 'b to 'c at the same time. Since the link can only take 10 packets, four packets are lost, among them data packet 23. Back in the link from 'a to 'b packet 24 is lost, and 25 arrives safely at 'b, which then discovers the hole for 24 and sends a NAK for 24 to 'c, as well as the data packet 25. 'c then reads and stores the NAK packet for 24 (with NAK count 1) in rule B7b, and then reads the packet 25 in rule B3b. This was the golden opportunity to discover the hole for packet 23, since it will never be discovered by the repair server, which had the packet. Instead of discovering the hole at 23, there was a NAK state for packet 24, and therefore 'c only initiates repairs from 24 and up in rule B3b, missing packet 23. Therefore, no repair will be attempted for packet 23. This same scenario can also be shown to exist in the informal specification.

#### 4.18 Specification and Analysis of the Combined Protocol

This section briefly sketches the specification and execution of the composition of the four protocol components that make up the AER/NCA suite of protocols. As mentioned in Section 4.9, we use object-oriented inheritance techniques to define the combined protocol. A sender in the combined protocol is an object of the following class `SenderCombined`:

```

class SenderCombined .
subclass SenderCombined < RTTsender NOMsender RCsender RScsender .

```

The `SenderCombined` class inherits all the attributes and rules of its superclasses. The definition of the receivers and the repair servers in the combined protocol is analogous:

```

class RepairserverCombined .
subclass RepairserverCombined < RTTrepairserver NOMrepairserver .
subclass RepairserverCombined < RCrepairserver RSrepairserver .

class ReceiverCombined .
subclass ReceiverCombined < RTTreceiver NOMreceiver RCreceiver RSreceiver .

```

While a “combined object” can perform all the rules defined on its superclasses, there are some composite transitions in which the different components must *synchronize* their actions when the components are combined. For example, the multicast of a *new* data packet is mainly a concern of the *repair service* component, but the *rate control* component must be consulted to check whether a new packet can be sent at the current time. In the combined protocol, we have therefore combined the parts dealing with sending new data packets from these two components into a single rule, presented below, involving objects of class `SenderCombined`. Similarly, when a new data packet is received by a receiver, the *repair service* component must buffer the received packet and check for missing packets, the *nominee selection* component must update its sliding window used to compute its loss probability estimate, and the *rate control* component must acknowledge the packet in case the receiver is the nominee receiver. There are only five such “combined” rules in our specification, out of a total of 76 rules.

The following rule models the attempt to send a new data packet in the combined protocol when the `sendDataTimer` expires. If the rate control part (the `sendAllowed` test) does not permit sending, then another attempt will be made one millisecond later. Otherwise, the `sendDataTimer` is reset to the sending interval of the data block:

```

r1 [A3sendG11G12] :
  < 0 : SenderCombined | children : OS, clock : R, sendDataTimer : 0,
    unsendDataPackets : data(Q, NZN, NZN', R') DB,
    reTransBuf : RTB, nextSeq : NZN'', maxSeqSent : N,
    winLowerSeq : WLS, csmNominee : DO, winSize : WS, ... >
=>
  if sendAllowed(N, WLS, WS, FRF, MBC)
  then
    < 0 : SenderCombined | unsendDataPackets : removeFirst(data(Q, NZN, NZN', R') DB),
      reTransBuf :
        (RTB MsgAndNAK(dataPacket(Q, NZN'', R, NZN == 1,
          NZN /= 1 and NZN /= NZN', false), 0)),
      nextSeq : NZN'' + 1,
      sendDataTimer :
        sendRate(removeFirst(data(Q, NZN, NZN', R') DB)),
      ... >
    multiSend(dataPacket(Q, NZN'', R, NZN == 1, NZN /= 1 and NZN /= NZN', false),
      0, OS)
  else < 0 : SenderCombined | sendDataTimer : 1 > fi .

```

We have executed one behavior of the combined protocol with two initial states, corresponding to the two initial states which invalidated the repair service component. In contrast to the execution of that component, all packets were delivered (in order) to each receiver in the single executions of the combined protocol provided by Real-Time Maude’s `trew` command. This was probably due to the presence of the rate control component, that adjusted the sending rate according to the

packet losses, thereby avoiding the extensive loss of packets which led to the above-described faulty behavior of the repair service component. The resulting states show that some data packets were indeed lost, but that they were successfully repaired.

Although the current combined protocol executes as expected, we found a significant flaw/omission in an earlier version of the protocol during execution: Only one data packet could be sent because the data packet was sent before a nominee was found. No receiver would then acknowledge the first data packet, and the second packet could not be sent before the first one was acknowledged. We solved this problem by changing the rewrite rules, so that the first data packet is not sent until a nominee is found.

#### 4.19 Summary of the Analysis Efforts

We have analyzed the four protocol components and the combined protocols by defining *some* initial states, and by analyzing, for each such initial state,

- *one* possible behavior from the initial state using timed rewriting.

For the RTT, NOM, and RC components we could also analyze

- *all* possible behaviors – up to a certain duration, and w.r.t. the choices of “random” values for the probabilistic parts of the protocol – from the initial state, using time-bounded search and temporal logic model checking.

Such analysis of the *RTT* component showed that the correct RTT values are found reasonably quickly, and that they are unchanged thereafter.

Timed rewriting analysis of the *NOM* component indicated that the correct nominee receivers are found. Using timed model checking we showed that the correct nominees are found at the appropriate times in all behaviors from the chosen initial states. However, using model checking we discovered the troubling scenario where, at some stages, no node is aware of it being the nominee.

The situation was somewhat “reversed” for the *RC* component, where timed rewriting yielded an unwanted behavior, which could be traced using the tool’s tracing capabilities; whereas the use of timed search showed that there exist behaviors, from the same initial state, which have the desired properties.

Timed rewriting was sufficient to find flaws in the *RS* component, which were then traced. Timed rewriting of different initial states gave the desired result where all packets were delivered to the receiver applications.

Finally, timed rewriting in the *combined protocol* yielded states where all packets were delivered to all receivers, even for those topologies for which the stand-alone RS component failed. This positive result was probably due to the addition of the rate control mechanism, which reduced the packet losses. Nevertheless, the flaws in the RS and NOM components carry over to the combined protocol; they are just more difficult to find. The difficulties have to do with the combinatorial explosion of states, given the size and degree of nondeterminism of the combined protocol; for these reasons we could not find the errors within reasonable time using timed search and model checking although we knew they were there. This is in fact one advantage of having modularly decomposed the protocol and having analyzed each of its components.

For all the analyses reported in this paper Real-Time Maude returned an answer within reasonable times (a few seconds to a few minutes). Therefore, except for the intrinsic combinatorial explosions alluded to above, we found that in practice the tool was quite usable for analyses of the kind performed.

#### 4.20 The Updated Informal Specification

Apart from minor changes such as correcting small errors and typos, making the state and communication assumptions explicit, and modifying the values of some constants used, the version 1.1 of the informal specification, that incorporated the results of our formal analysis, updates the above described version 1.0 in the following ways:

1. To solve the problem mentioned in Section 4.13, where there could be states with no nominee receiver, so that some data packet were not acknowledged even if no packets were lost, the data packets now have an additional field for the current nominee. Each receiver updates its `isNominee` flag according to this field upon the reception of *original* data packets. (A slight modification is that only the reception of original data packets is acknowledged.) Since data packets may get lost, also SPM packets are equipped with a field for nominee receivers. It is still the case that the sender only treats CCM packets from the *current* nominee receiver, and not necessarily from the receiver which was the nominee when the data packet was sent. Furthermore, in the initial phase there is no nominee, so the previously described deadlock scenario, where the first data packets are not acknowledged, could still happen if the sender sends data packets before a nominee receiver is found. As mentioned in Section 4.18, we have tried to avoid this deadlock by waiting for a certain amount of time before sending data packets. A better solution would be to modify the rule `A3sendG11G12` to disallow sending of data packets until a nominee is found, by adding the condition `D0 /= no0id` to the `if`-test.
2. In an attempt to solve the first troubling scenario of the repair service component mentioned in Section 4.17, where a requested packet is not repaired by a repair server because a NAK state exists for a packet which is no longer in the repair server's cache, the repair server now deletes the NAK state of a data packet which is removed from the repair server's cache.

## 5 Conclusions

We have discussed in detail our formalization and analysis in Real-Time Maude of the AER/NCA active network protocol suite. Being a quite complex distributed system with essential real-time and probabilistic features, and with performance requirements essential to its design and correct functioning, the modeling of AER/NCA presented a number of interesting challenges. We have explained how those challenges were successfully met by Real-Time Maude. As a fruit of this modeling and analysis work, important errors were found, and valuable insights were gained. First, all the errors in the use-case informal specification that the designers were familiar with, but did not tell us about, were independently uncovered by our analysis. Furthermore, several more subtle design errors not known to the designers, which impaired the intended correct behavior of AER/NCA and which were not discovered by traditional simulation and testing of an actual implementation, were found.

An important encouraging lesson learned was the intuitive appeal of Real-Time Maude specifications to network engineers, comparing in fact favorably with informal use-case specifications, and the associated low-threshold adoption barrier for rewriting logic based specification languages like Maude and Real-Time Maude. This agrees with our experience in teaching rewriting logic based formal methods to undergraduate students at the University of Oslo [19]. Both the simple direct representation of state transitions by rewrite rules, and the executable nature of the specifications—that allow a user to view them as programs in a programming language, with minimal or no acquaintance with the formal foundations—seem to be crucial aspects of this low adoption barrier.

More generally, there is by now ample experience on the usefulness and adequacy of rewriting logic for specifying and analyzing distributed systems in general and network systems in particular (see the survey [17], and recent advanced case studies such as [25, 6]). The present case study is a further substantial confirmation of this general experience for network applications in which real-time and resource sensitive behavior are crucial aspects to model. A more recent Real-Time Maude analysis of a new multicast protocol proposed by the IETF [13] further confirms this experience. A promising area with several ongoing Real-Time Maude specification efforts is wireless communication protocols. A final point—indeed quite relevant for wireless communication and for networked embedded systems—is the natural convergence of real-time and probabilistic specifications, something already exemplified by our AER/NCA case study. This convergence offers an exciting research opportunity to combine the best methods and tools developed so far for real-time rewrite theories and for probabilistic rewrite theories, and to develop new methods to fruitfully analyze probabilistic real-time specifications.

## Acknowledgments

We are grateful to Mark Keaton and Steve Zabele for their invaluable cooperation during the specification and analysis of a previous version, in Real-Time Maude 1.0, of the AER/NCA protocol suite. Their explanation of AER/NCA and related issues, their feedback to our specification efforts, and their suggestions of suitable initial states for the analysis parts were essential for the modeling and analysis described in this paper. Partial support of this research by ONR Grant N00014-02-1-0715, by NSF Grant CCR-0234524, by DARPA through Rome Labs. Contract F30602-97-C-0312, and by The Norwegian Research Council is gratefully acknowledged.

## References

- [1] M. Bozga, J.-C. Fernandez, L. Ghirvu, S. Graf, J.-P. Krimm, and L. Mounier. IF: A validation environment for timed asynchronous systems. In E. A. Emerson and A. P. Sistla, editors, *Proc. Computer Aided Verification (CAV'2000)*, volume 1855 of *Lecture Notes in Computer Science*, pages 543–547. Springer, 2000.
- [2] R. Bruni and J. Meseguer. Generalized rewrite theories. In J. C. M. Baeten, J. K. Lenstra, J. Parrow, and G. J. Woeginger, editors, *Proc. 30th International Colloquium on Automata, Languages and Programming (ICALP 2003)*, volume 2719 of *Lecture Notes in Computer Science*, pages 252–266. Springer, 2003.
- [3] M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and J. F. Quesada. Maude: Specification and programming in rewriting logic. *Theoretical Computer Science*, 285:187–243, 2002.
- [4] M. Clavel, F. Dúran, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and C. Talcott. *Maude Manual (Version 2.1)*, April 2004. <http://maude.cs.uiuc.edu>.

- [5] S. Eker, J. Meseguer, and A. Sridharanarayanan. The Maude LTL model checker. In F. Gadducci and U. Montanari, editors, *Fourth International Workshop on Rewriting Logic and its Applications*, volume 71 of *Electronic Notes in Theoretical Computer Science*. Elsevier, 2002.
- [6] A. Goodloe, C. A. Gunter, M. McDougall, A. Sridharanarayanan, and M.-O. Stehr. Formal specification of Sectrace: A protocol to set up security associations and policies in IPSec networks. [http://formal.cs.uiuc.edu/stehr/sectrace\\_eng.html](http://formal.cs.uiuc.edu/stehr/sectrace_eng.html).
- [7] D. Harel. From play-in scenarios to code: an achievable dream. In *Proc. FASE'00, 3rd Intl. Conf. on Fundamental Approaches to Software Engineering*, volume 1783 of *Lecture Notes in Computer Science*, pages 22–34. Springer, 2000.
- [8] T. A. Henzinger, P.-H. Ho, and H. Wong-Toi. HyTech: A model checker for hybrid systems. *Software Tools for Technology Transfer*, 1:110–122, 1997. See also HyTech home-page at <http://www-cad.eecs.berkeley.edu/~tah/HyTech/>.
- [9] S. Kasera, S. Bhattacharyya, M. Keaton, D. Kiwior, J. Kurose, D. Towsley, and S. Zabele. Scalable fair reliable multicast using active services. *IEEE Network Magazine (Special Issue on Multicast)*, 14(1):48–57, 2000.
- [10] D. E. Knuth. *The Art of Computer Programming: Seminumerical Algorithms*, volume 2. Addison-Wesley, second edition edition, 1981.
- [11] Nirman Kumar, Koushik Sen, José Meseguer, and Gul Agha. A rewriting based model of probabilistic distributed object systems. In *Proc. Formal Methods for Open Object-Based Distributed Systems (FMOODS 2003)*, volume 2884 of *Lecture Notes in Computer Science*, pages 32–46. Springer, 2003.
- [12] K. G. Larsen, P. Pettersson, and W. Yi. UPPAAL in a nutshell. *Software Tools for Technology Transfer*, 1(1–2):134–152, October 1997. See also UPPAAL home-page at <http://www.uppaal.com/>.
- [13] E. Lien. Formal modelling and analysis of the NORM multicast protocol using Real-Time Maude. Master’s thesis, Department of Linguistics, University of Oslo, 2004.
- [14] J. Meseguer. Conditional rewriting logic as a unified model of concurrency. *Theoretical Computer Science*, 96:73–155, 1992.
- [15] J. Meseguer. A logical theory of concurrent objects and its realization in the Maude language. In G. Agha, P. Wegner, and A. Yonezawa, editors, *Research Directions in Concurrent Object-Oriented Programming*, pages 314–390. MIT Press, 1993.
- [16] J. Meseguer. Membership algebra as a logical framework for equational specification. In F. Parisi-Presicce, editor, *Proc. WADT'97*, volume 1376 of *Lecture Notes in Computer Science*, pages 18–61. Springer, 1998.
- [17] J. Meseguer. Rewriting logic and Maude: a wide-spectrum semantic framework for object-based distributed systems. In S. Smith and C.L. Talcott, editors, *Formal Methods for Open Object-based Distributed Systems, FMOODS 2000*, pages 89–117. Kluwer, 2000.
- [18] P. C. Ölveczky. *Specification and Analysis of Real-Time and Hybrid Systems in Rewriting Logic*. PhD thesis, University of Bergen, 2000. Available at <http://maude.cs.uiuc.edu/papers>.
- [19] P. C. Ölveczky. Formal modeling and analysis of distributed systems in Maude. Course book for INF3230, Dept. of Informatics, University of Oslo, 2004.
- [20] P. C. Ölveczky. *Real-Time Maude 2.1 Manual*, 2004. <http://www.ifi.uio.no/RealTimeMaude/>.
- [21] P. C. Ölveczky, M. Keaton, J. Meseguer, C. Talcott, and S. Zabele. Specification and analysis of the AER/NCA active network protocol suite in Real-Time Maude. In H. Hussmann, editor, *Fundamental Approaches to Software Engineering (FASE 2001)*, volume 2029 of *Lecture Notes in Computer Science*, pages 333–347. Springer, 2001.
- [22] P. C. Ölveczky and J. Meseguer. Specification of real-time and hybrid systems in rewriting logic. *Theoretical Computer Science*, 285:359–405, 2002.

- [23] P. C. Ölveczky and J. Meseguer. Real-Time Maude 2.1. In N. Martí-Oliet, editor, *Fifth International Workshop on Rewriting Logic and its Applications*, Electronic Notes in Theoretical Computer Science. Elsevier, 2004. To appear.
- [24] P. C. Ölveczky and J. Meseguer. Specification and analysis of real-time systems using Real-Time Maude. In T. Margaria and M. Wermelinger, editors, *Fundamental Approaches to Software Engineering (FASE 2004)*, volume 2984 of *Lecture Notes in Computer Science*, pages 354–358. Springer, 2004.
- [25] M.-O. Stehr, C. Talcott, and G. Denker. Towards a formal specification of the Spread group communication system. [http://formal.cs.uiuc.edu/stehr/spread\\_eng.html](http://formal.cs.uiuc.edu/stehr/spread_eng.html).
- [26] P. Viry. Equational rules for rewriting logic. *Theoretical Computer Science*, 285:487–517, 2002.
- [27] S. Yovine. Kronos: A verification tool for real-time systems. *Software Tools for Technology Transfer*, 1(1/2), 1997. See also Kronos home-page at <http://www-verimag.imag.fr/TEMPORISE/kronos/>.