# A Client-Server Approach to Virtually Synchronous Group Multicast: Specifications and Algorithms

Idit Keidar                     Roger Khazan

Massachusetts Institute of Technology Lab for Computer Science
545 Technology Square, Cambridge, MA 02143, USA
Email: {idish, roger}@theory.lcs.mit.edu     URL: http://theory.lcs.mit.edu/{~idish, ~roger}

## Abstract

*This paper presents a formal design for a novel group multicast service that provides virtually synchronous semantics in asynchronous fault-prone environments. The design employs a client-server architecture in which group membership is maintained not by every process but only by dedicated membership servers, while virtually synchronous group multicast is implemented by service end-points running at the clients. Specifically, the paper defines service semantics for the client-server interface, that is, for the group membership service. The paper then specifies virtually synchronous semantics for the new group multicast service, as a collection of commonly used safety and liveness properties. Finally, the paper presents new algorithms that use the defined group membership service to implement the specified properties. The algorithm that provides the complete virtually synchronous semantics executes in a single message round in parallel with the membership service's agreement on views, and is therefore more efficient than previously suggested algorithms providing such semantics.*

**Keywords**: Group Communication, Virtual Synchrony, Reliable Multicast, Formal Modeling, I/O Automata.

## 1. Introduction

Group communication systems are powerful building blocks that facilitate the development of fault-tolerant distributed applications (see [1, 19, 5] for discussion of the utility of group communication systems). Group communication provides the notion of *group abstraction*, which allows processes to be easily organized in multicast groups. Group communication systems typically integrate two types of services: group membership and reliable group multicast. The membership service maintains a listing of the currently active and connected group members and delivers this information to its clients whenever it changes. The output of the membership service is called a *view*. Reliable multicast services that deliver messages to the current view members complement the membership service. In this paper, we present a novel group multicast service.

Group communication systems usually run in asynchronous fault-prone environments. In such environments, group communication systems generally provide some variant of *virtual synchrony* semantics which synchronize membership notifications with regular messages and thus simulate a "benign" world in which message delivery is reliable within the set of connected processes. Such semantics are especially useful for constructing fault-tolerant applications that maintain consistent replicated state of some sort (e.g., [7, 14]). The key aspect of virtual synchrony is the semantics of interleaving of message send and delivery events with view delivery events. In order to reason about this interleaving, we associate message send and delivery events with views: we say that an event e occurs at a process p *in view* v if v was the last view delivered to p before e, or a default initial view $v_p$ if no such view was delivered.

Many variants of virtual synchrony semantics have been suggested (e.g., [16, 8, 19, 17, 7]). A key property specified by nearly all of these (e.g., [16, 8, 19, 17]) is that processes moving together from a view v to another view v' deliver the same messages in v. This property allows applications to avoid costly re-synchronization following certain view changes. Our service specification includes this property, as well as several additional safety and liveness properties. We present our specifications in Section 4.

During the period in which the group communication service is attempting to reach agreement on a view, processes may attempt to join/re-join. In such cases, previously suggested virtual synchrony algorithms, e.g., [8, 17], can have the current invocation of the membership and virtual synchrony proceed to termination without adding the joining processes, and then immediately start an attempt to add them. This strategy results in overhead (e.g., increased network load) because applications react to such outdated views just as they do to any other view, e.g., by re-synchronizing with the new members. Moreover, this strategy precludes situations when applications may rely on virtual synchrony to avoid the costly re-synchronizations all together. For example, consider a transient failure when a process p is unsuspected right after an attempt to remove p

from the membership has started. Existing algorithms typically deliver a view excluding p and then re-invoke the algorithm to allow p to re-join. Since p does not move into the resulting view from the same view as the rest of the processes, these processes cannot rely on the key property of virtual synchrony to avoid re-synchronizing with p. In contrast, our algorithm never delivers views that reflect a membership that is already known to be out of date.

Traditionally, virtual synchrony semantics were implemented by algorithms that were integrated with group membership algorithms (e.g., in [8, 9, 2]). In contrast, our group multicast service is designed for a client-server architecture in which a small set of dedicated membership servers maintains client membership information (i.e., which clients are members of which group). This architecture was designed to provide scalable membership services in wide area networks (cf. [3]). Our virtual synchrony algorithm acts as the client of an *external* membership service.

Introducing the client-server design poses a major challenge: One has to define an interface by which a membership server interacts with its clients, in a way that would allow for simple and efficient implementations of both group membership (by the membership servers), and virtual synchrony (by service end-points at the clients). Such an interface has to provide sufficient level of synchronization to allow the virtual synchrony algorithm to reach agreement upon the set of messages delivered in the old view in *parallel* with the servers' agreement on views. At the same time, the virtual synchrony algorithm should avoid imposing limitations on the membership's choice of views (as explained above). In addition, one has to try to minimize the communication overhead induced by the client-server interaction.

We have designed an interface that addresses the challenges above. Our interface consists of two types of messages sent from membership servers to their clients: When a server engages in a view change, it sends its clients a `start_change` message. Each `start_change` message contains a locally unique identifier. This identifier is *not* globally agreed upon: `start_change` messages sent to different processes can contain different identifiers. Once the server agrees upon the new view with the other servers, it sends a `view` message to its clients. The view contains information that maps clients to the last `start_change` identifiers they received before receiving this view. A similar view structure is suggested in [17], for the purpose of not having concurrent views intersect. The servers do not need to hear from their clients in order to complete the algorithm.

Our interface allows for straightforward and efficient implementations of both membership and virtual synchrony. The algorithm we present in Section 5 exploits this interface to achieve virtual synchrony in a single round. We have implemented this algorithm (in C++) using the scalable one-round membership algorithm of [13]. The virtual synchrony round and the membership round are conducted in parallel: once the end-points receive the `start_change` notifi-

cations, they send each other special synchronization messages which allow them to agree upon the set of messages to be delivered before moving to the new view. We are not aware of any other algorithm that implements virtual synchrony in one communication round without pre-agreement upon a globally unique identifier while also not imposing restrictions on the membership's choice of the next view.

Throughout this paper we use the I/O automaton formalism (cf. [15], Ch. 8) to provide rigorous specifications and algorithm descriptions. Previously suggested I/O automaton-style specifications of group communication systems (e.g., [7, 10]) used a single abstract automaton to represent multiple properties of the same system component and presented a single algorithm automaton that implements all of these properties. Thus, no means were provided for reasoning about a subset of the properties, and it was often difficult to follow which part of the algorithm implements which part of the specification. We address this shortcoming by specifying separate properties as separate abstract automata, and by incrementally constructing the algorithm that implements them – in each step adding support for an additional property – using a novel *inheritance*-based construct, recently introduced to the I/O automaton model [12]. This paper informally argues the algorithm's correctness; a formal correctness proof by simulation is included in the full paper [11].

## 2. Formal Model and Notation

In the I/O automaton model (cf. [15], Ch. 8), a system component is described as a state-machine, called an *I/O automaton*. The transitions of this state-machine are associated with named actions, which are classified as either *input*, *output*, or *internal*. Input and output actions model the component's interaction with other components, while internal actions are externally-unobservable.

Formally, an I/O automaton is defined as the following five-tuple: a signature (input, output and internal actions), a set of states, a set of start states, a state-transition relation (a cross-product between states, actions, and states), and a partition of output and internal actions into *tasks*. Tasks are used for defining fairness conditions.

An action $\pi$ is said to be *enabled* in a state s if the automaton has a transition of the form $(s, \pi, s')$; input actions are enabled in every state. An *execution* of an automaton is an alternating sequence of states and actions that begins with its start state and in which every action is enabled in the preceding state. An infinite execution is *fair* if, for each task, it either contains infinitely many actions from this task or infinitely many occurrences of states in which no action from this task is enabled; a finite execution is *fair* if no action is enabled in its final state. A *trace* is a subsequence of an execution consisting solely of the automaton's external actions. A *fair trace* is a trace of a fair execution.

When reasoning about an automaton, we are only interested in its externally-observable behavior as reflected in its

traces. There are two types of trace properties: *safety* and *liveness*. Safety properties usually specify that some particular bad thing never happens. In this paper we specify safety properties using centralized (global) I/O automata that generate the legal sets of traces; for such automata we do not specify task partitions. Each external action in such a centralized automaton is tagged with a subscript which denotes the process at which this action occurs. An algorithm automaton *satisfies* a specification if all of its traces are also traces of the specification automaton. Liveness properties usually specify that some good thing eventually happens. An implementation automaton satisfies a liveness property if the property holds in all of its *fair* traces.

The *composition operation* defines how automata interact via their input and output actions: It matches output and input actions with the same name in different component automata; when a component automaton performs a step involving an output action, so do all components that have this action as an input one. When reasoning about a certain system component, we compose it with abstract specification automata that specify the behavior of its environment.

I/O automata are conveniently presented using the *precondition-effect* style: In this style, typed state variables with initial values specify the set of states and the start states. A variable type is a set (if $S$ is a set, the notation $S_\perp$ refers to the set $S \cup \{\perp\}$). Transitions are grouped by action name, and are specified as a list of triples consisting of an action name (possibly with parameters), a `pre :` block with preconditions on the states in which the action is enabled, and an `eff :` block which specifies how the pre-state is modified *atomically* to yield the post-state.

We use a novel *inheritance*-based formal concept, recently introduced into the I/O automaton model [12]. A *child* automaton is specified as a modification of the parent automaton's code. When presenting a child we first specify a *signature extension* which consists of new actions (labeled new) and modified actions (a modified action is labeled with the name of the action which it modifies as follows: modifies `parent.action(parameters)`). We next specify the *state extension* consisting of new state variables added by the child. Finally, we describe the *transition restriction* which consists of new preconditions and effects added by the child to both new and modified actions. For modified actions, the preconditions and effects of the parent are appended to those added by the child. New effects added by the child are performed before the effects of the parent, all of them in a single atomic step. The child's effects are not allowed to modify state variables of the parent, to ensure that the set of traces of the child, when projected onto the parent's signature, is a subset of the parent's set of traces.

## 3. Environment Specification

Our service is implemented in an asynchronous message-passing environment. Processes can crash, communication links may fail and may later recover, possibly
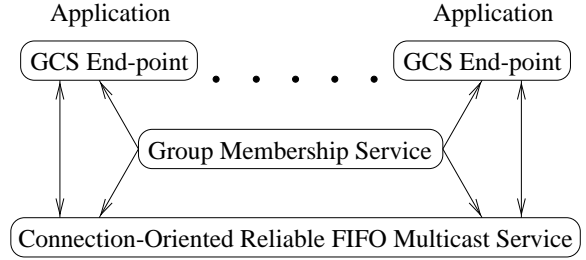


**Figure 1. The client-server architecture.**

causing network partitions and merges. In [11], we also model process recovery.

The service is implemented by *group communication service (GCS) end-points* running as clients of an external membership service whose specification appears in Section 3.1. The end-points communicate with each other using a reliable FIFO multicast service which we describe in Section 3.2, as depicted in Figure 1. We use the words "process" and "end-point" interchangeably.

### 3.1. The membership service

In Figure 2 we specify an external membership service whose interface consists of two output events:

`start_change`$_p$`(cid, set)` notifies process p that the membership service is attempting to form a new view with the members of `set`; `cid` is a local identifier.

`view`$_p$`(v)` notifies process p of the new view v. A view v is a triple consisting of an identifier v.id, a set of members v.set, and a function v.startId that maps members of v to start_change identifiers. Two views are the same if they consist of identical triples.

The membership specification captures two basic membership properties, which are fulfilled by virtually all group membership services (e.g., [2, 8, 4, 13, 17]): *Self Inclusion* requires every view delivered to an end-point p to include p as a member, and *Local Monotonicity* requires that view identifiers delivered to p be monotonically increasing.

In addition, the MBRSHP automaton specifies, using the mode[p] variable, that the membership service must precede every view v sent to end-point p with at least one `start_change` notification to p. It also requires that, for every view v sent to p, the start_change identifier v.startId(p) be the same as the `cid` of the latest `start_change` issued to p before the view, and that v.set be a subset of the `set` suggested in that `start_change`. Note that this specification does allow the membership service to add new processes while it is reconfiguring, as long as a new `start_change` is sent to the clients. Also note that the specified service is *partitionable* [19, 4], i.e., allows several disjoint views to exist concurrently.

The specification allows for simple and efficient distributed implementations, e.g., [13], as well as many other

**Type** StartChangeId: Total-order; $cid_0$ is smallest.
       ViewId: Partial-order; $vid_0$ is smallest.
       View: ViewId x SetOf(Proc) x
                              x (Proc $\rightarrow$ StartChangeId)
**Def**  $v_p = \langle vid_0, \{p\}, \{(p \rightarrow cid_0)\}\rangle$

**Signature:**
  Output: start_change$_p$(cid, set), Proc p,
                StartChangeId cid, SetOf(Proc) set
          view$_p$(v), Proc p, View v

**State:**
  ($\forall$ Proc p) View  mbrshp_view[p], initially $v_p$
  ($\forall$ Proc p) (StartChangeId x SetOf(Proc))
          start_change[p], initially $\langle cid_0, \{\}\rangle$
  ($\forall$ Proc p) mode[p] $\in$ {normal, change_started},
          initially normal

**Transitions:**
  OUTPUT  **start_change$_p$**(cid, set)
  pre: cid > start_change[p].id
       p $\in$ set
  eff: start_change[p] $\leftarrow$ $\langle$cid, set$\rangle$
       mode[p] $\leftarrow$ change_started


  OUTPUT  **view$_p$**(v)
  pre: p $\in$ v.set $\wedge$ v.id > mbrshp_view[p].id
       v.set $\subseteq$ start_change[p].set
       v.startId(p) = start_change[p].id
       mode[p] = change_started
  eff: mbrshp_view[p] $\leftarrow$ v
       mode[p] $\leftarrow$ normal

**Figure 2. Membership service safety spec.**

existing membership algorithms (e.g., [2]) which could be easily extended to provide the specified interface and semantics. In a possible implementation, a small number of servers could support a large number of clients, communicating with them asynchronously via FIFO ordered channels. Fault-tolerant implementations that support client migration are also possible if the server name is included in the start_change identifier to guarantee its local uniqueness.

## 3.2. The reliable FIFO multicast service

The group communication end-points communicate with each other using an underlying multicast service that provides reliable FIFO communication between every pair of connected processes. Many existing group communication systems (e.g., [9, 4]) implement virtual synchrony over similar underlying reliable communication substrates. In our implementation, we currently use the service of [18].

Figure 3 presents a centralized automaton CO_RFIFO which specifies a multicast service appropriate for our group communication algorithm. CO_RFIFO maintains a FIFO queue channel[p][q] for every pair of end-points. An input action send$_p$(set, m) models the multicast of message m from end-point p to the end-points listed in the set by ap-

**Signature:**
 Input:
   send$_p$(set,m), Proc p, SetOf(Proc) set, Msg m
   reliable$_p$(set), Proc p, SetOf(Proc) set
   live$_p$(set), Proc p, SetOf(Proc) set
 Output:   deliver$_{p,q}$(m), Proc p, Proc q, Msg m
 Internal: lose(p,q), Proc p, Proc q

**State:**
  ($\forall$ Proc p, Proc q) SeqOf(Msg) channel[p][q],
                            initially empty
  ($\forall$ Proc p) SetOf(Proc) reliable_set[p], init. {p}
  ($\forall$ Proc p) SetOf(Proc) live_set[p], init. {p}

**Transitions:**
  INPUT **send$_p$**(set, m)
  eff: ($\forall$ q $\in$ set) append m to channel[p][q]


  OUTPUT **deliver$_{p,q}$**(m) hidden parameter live_set[p]
  pre: m = First(channel[p][q])
  eff: dequeue m from channel[p][q]


  INPUT **reliable$_p$**(set)
  eff: reliable_set[p] $\leftarrow$  set

  INTERNAL **lose**(p, q)
  pre: q $\notin$ reliable_set[p]
  eff: dequeue last message from channel[p][q]

  INPUT **live$_p$**(set)
  eff: live_set[p] $\leftarrow$ set

**Tasks:**
  1. ($\forall$ p)($\forall$ q $\in$ live_set[p]) {deliver$_{p,q}$(m)}
  2. {dummy()} $\cup$ {deliver$_{p,q}$(m) | q $\notin$ live_set[p]}
     $\cup$ {lose(p,q)}

**Figure 3. CO_RFIFO service specification.**

pending m to the queues channel[p][q] for every end-point q in set. The deliver$_{p,q}$(m) action removes the first message from channel[p][q] and delivers it to q.

An end-point p may use the action reliable$_p$(set) to require CO_RFIFO to maintain reliable (gap-free) FIFO connections to the end-points in set. For every process q not in this set, CO_RFIFO may lose an arbitrary suffix of the messages sent from p to q, as modeled by the action lose(p, q).

In specifying liveness of CO_RFIFO, we require that messages sent to live and connected processes eventually reach their destinations. We formulate this property by defining every deliver$_{p,q}$(m) to be a task if and only if q is a member of live_set[p], a special variable periodically set by input actions live$_p$(set). The live$_p$(set) inputs are assumed to reflect the real state of the network, that is, the set of processes which are really alive and connected to p. Notice that we could not use the variable reliable_set[p] in this formulation because it is controlled by the client and thus does not necessarily reflect the *real* network situation.

## 4. GCS Specifications

We present the safety and liveness properties satisfied by our group communication service in Sections 4.1 and 4.2 respectively. These properties have been proven to be useful for many distributed applications (see [19]).

### 4.1. Safety properties

We present our safety specifications in four steps, as four automata: In Section 4.1.1 we specify a simple group communication service that provides reliable FIFO multicast within views. In Section 4.1.2 we extend the specification of Section 4.1.1 to also require that processes moving together from view v to view v′ deliver the same messages in view v. In Section 4.1.3 we specify a service which provides *transitional sets* (first presented as part of *Extended Virtual Synchrony (EVS)* [16]). In Section 4.1.4 we specify the Self Delivery property which requires processes to deliver their own messages. The specified services are partitionable.

#### 4.1.1 Within-view reliable FIFO multicast

```
AUTOMATON WV_RFIFO : SPEC

Signature:
 Input:  send_p(m), Proc p, AppMsg m
 Output: deliver_p(q, m), Proc p, Proc q, AppMsg m
         view_p(v), Proc p, View v

State:
 (∀ Proc p, View v) SeqOf(AppMsg) msgs[p][v],
                    initially empty
 (∀ Proc p, Proc q) Int last_dlvrd[p][q], init. 0
 (∀ Proc p) View current_view[p], init. v_p

Transitions:
 INPUT  send_p(m)
 eff: append m to msgs[p][current_view[p]]

 OUTPUT  deliver_p(q, m)
 pre: m=msgs[q][current_view[p]][last_dlvrd[q][p]+1]
 eff: last_dlvrd[q][p] ← last_dlvrd[q][p]+1

 OUTPUT  view_p(v)
 pre: p ∈ v.set ∧ v.id > current_view[p].id
 eff: (∀ q) last_dlvrd[q][p] ← 0
      current_view[p] ← v
```

**Figure 4.** WV_RFIFO **service specification.**

In Figure 4 we present the within-view reliable FIFO (WV_RFIFO) service specification. The specification uses centralized queues msgs[p][v] of application messages for each sender p and view v. The action send_p(m) models the multicast of message m from process p to the members of p's current view by appending m to msgs[p][current_view[p]]. The deliver_p(q, m) action models the delivery to process p of message m sent by process q while in p's current view. The specification enforces gap-free FIFO ordered delivery of

messages by using the variable last_dlvrd[q][p] to index the last message from q delivered to p in p's current view.

This specification captures the following properties:

- Views delivered to the application satisfy Local Monotonicity and Self Inclusion (cf. Sec. 3.1).

- Messages are delivered in the same view in which they were sent. This property is useful for many applications (cf. [8, 19]) and appears in several systems and specifications (e.g., see [1, 2, 16, 7]). A weaker property that requires each message to be delivered in the same view at every process that delivers it, but not necessarily the view in which it was sent, is typically implemented on top of an implementation of within-view delivery (see [19]).

- Messages are delivered in gap-free FIFO order (within views). This is a basic property upon which one can build services with stronger ordering guarantees (e.g., causally or totally ordered multicast).

#### 4.1.2 Virtual synchrony

In this section we specify a virtually synchronous reliable FIFO multicast service, VS_RFIFO, as a child of the presented above WV_RFIFO automaton. The VS_RFIFO specification consists of the code given in both Figures 4 and 5.

```
AUTOMATON VS_RFIFO : SPEC   MODIFIES WV_RFIFO : SPEC

Signature Extension:
 Output:   view_p(v) modifies wv_rfifo.view_p(v)
 Internal: set_cut(v, v′, c), View v, View v′,
                        (Proc → Int)_⊥ c new

State Extension:
 (∀ View v, v′) (Proc→Int)_⊥ cut[v][v′], init. ⊥

Transition Restriction:
 OUTPUT  view_p(v)
 pre: cut[current_view[p]][v] ≠ ⊥
 (∀ q) last_dlvrd[q][p]=cut[current_view[p]][v](q)

 INTERNAL  set_cut(v, v′, c)
 pre: cut[v][v′] = ⊥
 eff: cut[v][v′] ← c
```

**Figure 5.** VS_RFIFO **service specification.**

In addition to the properties inherited from WV_RFIFO, the VS_RFIFO specification also requires that processes moving together from view v to view v′ deliver the same set of messages in v. To enforce this property, the specification introduces internal actions set_cut(v, v′, c) that nondeterministically fix the set of messages to be delivered in view v by every process that moves from v to v′. This set is represented by the index of the last message to be delivered in v from each sender. Note that a process that delivers messages beyond an already established cut is not allowed to move into the view associated with the cut.

This property is commonly provided (e.g., [16, 8, 19, 17, 10]) and is often called *Virtual Synchrony* by itself. It is especially useful for applications that implement data replication using the state machine approach (e.g., [7, 14]). Such applications may exploit Virtual Synchrony to avoid sending costly synchronization messages among processes that continue together from one view to the next.

### 4.1.3 Transitional set

While Virtual Synchrony is a useful property, a process that moves from view $v$ to view $v'$ cannot locally tell which of the processes in $v.set \cap v'.set$ move to view $v'$ directly from view $v$, and which move to $v'$ from some other view. In order for the application to be able to exploit the Virtual Synchrony property, application processes need to be told which other processes move together with them from their old views in to their new views. The set of such processes is called a *transitional set*. The notion of a transitional set was first introduced as part of a special transitional view in the EVS [16] model. In our formulation (as in [19]), transitional sets are delivered to the applications together with (regular) views, as an additional parameter $T$. The delivery of transitional sets satisfies the following property (cf. [19]):

**Property 4.1** *The transitional set delivered by a process* $p$ *when it moves from view* $v$ *to view* $v'$ *is a subset of* $v.set \cap v'.set$ *that includes (a) all the processes (including* $p$*) that move directly from* $v$ *to* $v'$ *and (b) no member of* $v'.set$ *that moves to* $v'$ *from any view other than* $v$.

Note that processes that move to the same view from different views deliver different transitional sets.

AUTOMATON TRANS_SET : SPEC

**Signature :**
```
Output: view_p(v,T), Proc p, View v, SetOf(Proc) T
Internal: set_prev_view_p(v), Proc p, View v
```

**State :**
```
(∀ Proc p)  View current_view[p], initially v_p
(∀ Proc p, View v) View⊥ prev_view[p][v], init. ⊥
```

**Transitions :**
```
 OUTPUT  view_p(v, T)
 pre: prev_view[p][v] = current_view[p]
      (∀ q ∈ v.set ∩ current_view[p].set)
           prev_view[q][v] ≠ ⊥
      T = {q ∈ v.set ∩ current_view[p].set |
           prev_view[q][v] = current_view[p]}
 eff: current_view[p] ← v

 INTERNAL  set_prev_view_p(v)
 pre: p ∈ v.set
      prev_view[p][v] = ⊥
 eff: prev_view[p][v] ← current_view[p]
```

**Figure 6. Transitional set specification.**

Figure 6 contains an automaton (without inheritance) specifying the Transitional Set property. Before $p$ can

move from view $v$ to $v'$, each member $q$ of $v.set \cap v'.set$ must execute $set\_prev\_view_q(v')$ to "declare" the view from which it intends to move to $v'$; this action sets $prev\_view[q][v']$ to $q$'s current view. The transitional set delivered by $p$ with $v'$ is then computed to consist of those $q$ in $v.set \cap v'.set$ for which $prev\_view[q][v'] = v$.

#### 4.1.4 Self delivery

In Figure 7 we modify the WV_RFIFO specification automaton (Fig. 4) to capture the Self Delivery property by forbidding an end-point $p$ to move from view $v$ to $v'$ before delivering all its own application messages sent in $v$.

AUTOMATON SELF : SPEC   MODIFIES WV_RFIFO : SPEC

**Signature Extension :**
```
 Output: view_p(v) modifies wv_rfifo.view_p(v)
```

**Transition Restriction :**
```
 OUTPUT  view_p(v)
 pre: last_dlvrd[p][p] =
       = LastIndexOf(msgs[p][current_view[p]])
```

**Figure 7. Self Delivery property specification.**

This *safety* property, when accompanied by the liveness property of Section 4.2, implies the Self Delivery *liveness* properties of [19] and [16], which require processes to *eventually* deliver their own messages.

### 4.2. Liveness property

In a fault-prone asynchronous model, it is not feasible to require that a group communication service be live in every execution. The only way to specify useful liveness properties without strengthening the communication model is to make these properties *conditional* on the underlying network behavior (as specified, e.g., in [7, 19]). Since our GCS uses an external membership service, we condition its liveness on the behavior of the membership service (which itself is assumed to satisfy some meaningful liveness properties, e.g., those of [13]). Provided the membership eventually delivers the same view to all the view end-points and does not deliver any subsequent views (i.e., stabilizes), we require the end-points to eventually deliver this view and all the messages sent in this view to their applications.

**Property 4.2** *Let* $v$ *be a view with* $v.set = S$. *Let* $\alpha$ *be a fair execution of a group communication service* GCS *in which, for every* $p \in S$, MBRSHP.view$_p(v)$ *action occurs and is followed by neither* MBRSHP.view$_p$ *nor* MBRSHP.start_change$_p$. *Then at each* $p \in S$, GCS.view$_p(v)$ *eventually occurs. Furthermore, for every* GCS.send$_p(m)$ *that occurs after* GCS.view$_p(v)$, *and for every* $q \in S$, GCS.deliver$_q(p, m)$ *also occurs.*

It is important to note that although our liveness property requires GCS to be live only in *certain* executions, any

implementation which satisfies this property has to attempt to be live in *every* execution because of its inability to test the external condition of the membership becoming stable. Also note that, even though membership stability is formally required to last forever, in practice it only has to hold "long enough" for GCS to reconfigure.

## 5. The Group Multicast Algorithm

Our group communication service is implemented by a collection of GCS end-points, each running the same algorithm. Figure 8 (a) shows the interaction of a GCS end-point with its environment, MBRSHP and CO_RFIFO (see Sec. 3). The end-point interacts with its application client by accepting the client's send-requests and by delivering application messages and views to the client. The end-point uses CO_RFIFO to send messages to other GCS end-points and to receive messages sent by other GCS end-points. When necessary, the end-point uses the action `reliable` to inform CO_RFIFO of the set of end-points to which CO_RFIFO must maintain reliable (gap-free) FIFO connections. The GCS end-point also receives `start_change` and `view` notifications from the membership service.



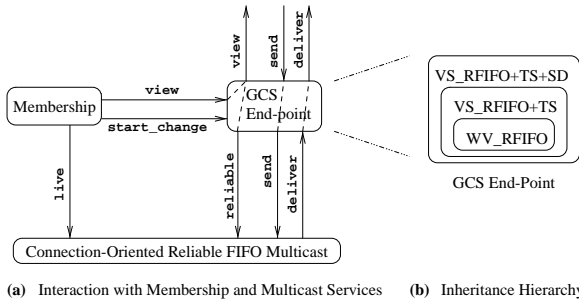(a)  Interaction with Membership and Multicast Services    (b)  Inheritance Hierarchy

### Figure 8. GCS end-point and environment.

The algorithm running at each end-point is constructed in steps, at each step adding support for a new property:

- First, we present an algorithm $WV\_RFIFO_p$ for an end-point of a within-view reliable FIFO multicast service.

- Then, in Section 5.2, we add support for the Virtual Synchrony and Transitional Set properties. We present a child $VS\_RFIFO+TS_p$ of $WV\_RFIFO_p$, and argue that the service built from $VS\_RFIFO+TS_p$ end-points satisfies safety specifications VSRFIFO : SPEC and TS : SPEC, and liveness Property 4.2.

- Finally, in Section 5.3, we add support for Self Delivery. The resulting automaton $VS\_RFIFO+TS+SD_p$ models a GCS end-point. Due to the use of inheritance, the service built from these end-points satisfies WV_RFIFO : SPEC, VSRFIFO : SPEC, and TS : SPEC. We argue that it also satisfies safety specification SELF : SPEC and liveness Property 4.2.

In the presented automata, each locally controlled action is defined to be a task by itself, which means that, if it becomes and stays enabled, it would eventually get executed.

When composing automata into a service, actions of the type $MBRSHP.start\_change_p(id, set)$ are linked with $CO\_RFIFO.live_p(set)$, and $MBRSHP.view_p(v)$ with $CO\_RFIFO.live_p(v.set)$. This way, the $live\_set[p]$ at CO_RFIFO matches the MBRSHP's perception of which processes are alive and connected to p. (We assume that every permanently disconnected end-point is eventually excluded by either a `start_change` or a `view` notification.) Also, in the composed system, all output actions except the application interface are reclassified as internal.

We present our algorithm at a level that would be easy to follow and then supplement this presentation with a discussion of some important, practical optimizations.

### 5.1. Within-view reliable FIFO multicast algorithm

In this section we present algorithm $WV\_RFIFO_p$ for an end-point p of a service that interacts with MBRSHP and CO_RFIFO services and satisfies the WV_RFIFO : SPEC safety specification and liveness Property 4.2.

The MBRSHP and CO_RFIFO services by themselves already provide most of the properties required by the WV_RFIFO : SPEC specification: MBRSHP generates views that satisfy Local Monotonicity and Self Inclusion, and CO_RFIFO provides gap-free FIFO communication. Since $WV\_RFIFO_p$ can just forward to its application the views generated by MBRSHP and can use CO_RFIFO to multicast application messages to other end-points, it only needs to ensure that messages are delivered in the same views in which they were sent. This can be done simply by tagging messages with the views in which they were sent and by allowing delivery of a message when its view tag matches the end-point's current view.

As an optimization of this idea, instead of tagging each message with a view in which it was sent, our algorithm sends a single, special $view\_msg(v)$ to all members of view v before sending them application messages in that view. An end-point can deduce the view in which an application message is sent from the latest $view\_msg(v)$ received from the application message sender.

The algorithm is captured in the $WV\_RFIFO_p$ automaton of Figure 9. Note that, instead of blindly relying on CO_RFIFO regarding which messages get delivered in a given view, WV_RFIFO allows processes to forward application messages on behalf of other processes. The code of $WV\_RFIFO_p$ does not specify a particular forwarding strategy – it allows for non-deterministic forwarding of messages. Without this, more refined versions and extensions of WV_RFIFO would not be able to introduce a specific forwarding strategy (as we do in VS_RFIFO+TS by adding a precondition on the action that sends forwarded messages).

There is also another place where the code leaves a non-deterministic choice: it is in handling of the `reliable_set`

**Signature:**
```
Input:   send_p(m), AppMsg m
         co_rfifo.deliver_q,p(m), Proc q,
              (AppMsg + ViewMsg + FwdMsg) m
         mbrshp.view_p(v), View v

Output:  deliver_p(q, m), Proc q, AppMsg m
         co_rfifo.send_p(set, m), SetOf(Proc) set,
              (AppMsg + ViewMsg + FwdMsg) m
         view_p(v), View v
```

**State:**
```
// Variables for handling application messages
(∀ Proc q, View v) SeqOf(AppMsg_⊥) msgs[q][v],
                              initially empty
Int last_sent, initially 0
(∀ Proc q) Int last_rcvd[q],  initially 0
(∀ Proc q) Int last_dlvrd[q], initially 0

// Variables for handling views and view messages
View current_view, initially v_p
View mbrshp_view,  initially v_p
For all Proc q: View view_msg[q], initially v_q

SetOf(Proc) reliable_set, initially v_p.set
```

**Transitions:**
```
INPUT  mbrshp.view_p(v)
eff: mbrshp_view ← v


OUTPUT  view_p(v)
pre: v =  mbrshp_view ≠ current_view
eff: current_view ← v
     last_sent  ← 0
     (∀ q) last_dlvrd[q] ← 0


OUTPUT  co_rfifo.reliable_p(set)
pre: current_view.set ⊆ set
eff: reliable_set ← set


OUTPUT  co_rfifo.send_p(set, tag=view_msg, v)
pre: view_msg[p] ≠ current_view
     current_view.set ⊆ reliable_set
     set =  current_view.set - {p}
     v = current_view
eff: view_msg[p] ← current_view


INPUT  co_rfifo.deliver_q, p(tag=view_msg, v)
eff: view_msg[q]  ← v
     last_rcvd[q] ← 0
```

```
INPUT  send_p(m)
eff: append m to msgs[p][current_view]


OUTPUT  deliver_p(q, m)
pre: m = msgs[q][current_view][last_dlvrd[q]+1]
     (q = p)  ⇒  (last_dlvrd[q] < last_sent)
eff: last_dlvrd[q] ← last_dlvrd[q] + 1


OUTPUT  co_rfifo.send_p(set, tag=app_msg, m)
pre: view_msg[p] = current_view
     set =  current_view.set - {p}
     m = msgs[p][current_view][last_sent + 1]
eff: last_sent ← last_sent + 1


INPUT  co_rfifo.deliver_q,p(tag=app_msg, m)
eff: msgs[q][view_msg[q]][last_rcvd[q]+1] ←m
     last_rcvd[q] ← last_rcvd[q] + 1


OUTPUT  co_rfifo.send_p(set,tag=fwd_msg,r,v,m,i)
pre: m = msgs[r][v][i]


INPUT  co_rfifo.deliver_q,p(tag=fwd_msg,r,v,m,i)
eff:  msgs[r][v][i] ← m
```

**Figure 9. Within-view reliable** FIFO **multicast end-point automaton.**

of CO_RFIFO. The code allows it to be an arbitrary superset of `current_view.set`. This set is further restricted in a child VS_RFIFO+TS$_p$ of WV_RFIFO$_p$.

The correctness of WV_RFIFO follows from the use of ordered message queues, the safety and liveness properties of CO_RFIFO, and the safety properties of MBRSHP. A formal proof is given in [11].

Also note that the presented code never removes messages from its buffers. An actual implementation can and should employ some sort of a garbage collection mechanism, for example discard messages sent in older views when moving in to a new view.

## 5.2. Virtual Synchrony and Transitional Sets

The WV_RFIFO service presented above guarantees that in each view v every member delivers some prefix of the FIFO ordered messages sent by each end-point in v.

The VS_RFIFO+TS service presented in this section extends WV_RFIFO to also guarantee that those end-points which transition directly from view v to the same view v′ deliver not just "some" prefixes but "the same" prefixes of the FIFO ordered messages sent by each end-point in view v (cf. Sec. 4.1.2). Moreover, every view delivery is accompanied by a transitional set T that satisfies the Transitional Set property of Sec. 4.1.3.

In order to satisfy these two properties, an end-point moving from a view v to a view v′ must first learn which other end-points may transition from v to v′ and must agree with them on the lengths of the prefixes they need to deliver. In a nutshell, here is how the VS_RFIFO+TS service accomplishes this: Each time an end-point p is notified via MBRSHP.start_change$_p$(cid, set) of the MBRSHP's attempt to form a new view, p reliably sends to set a synchronization message tagged with cid. When MBRSHP.view$_p$(v′) is delivered to p, p uses the v′.startId

mapping to determine which synchronization message to use from each end-point q in v.set ∩ v'.set; it uses the one tagged with v'.startId(q). As a result, all end-points that move from view v to v' use the same set of synchronization messages for computing the transitional set and the set of messages to be delivered to their application clients before v'. Notice that, by enriching views with the startId mapping, we eliminate the need to pre-agree on a common tag for identifying which synchronization messages to consider for a given view.

### 5.2.1 Algorithm details and safety argument

Figure 10 presents the VS_RFIFO+TS$_p$ automaton as a child of WV_RFIFO$_p$. While there are no view changes, VS_RFIFO+TS$_p$ does not modify the behavior of WV_RFIFO$_p$. During a view change, VS_RFIFO+TS$_p$ sends and handles synchronization messages, and also restricts the delivery of application messages according to the synchronization messages associated with the new view.

Upon receiving a start_change$_p$(cid, set) notification from MBRSHP, end-point p stores ⟨cid, set⟩ in the variable start_change, tells CO_RFIFO to maintain reliable communication to the end-points in current_view ∪ set, and then sends a synchronization message tagged with cid to every end-point in set. The synchronization message contains p's current view v and a cut, which is a mapping from processes to indices; cut(q) is the index of the last message from q that p *commits* to deliver before delivering any view v' with v'.startId(p) = cid.

End-point p stores the synchronization message from q tagged with cid in sync_msg[q][cid]. Until p receives a view from MBRSHP, it does not know which synchronization messages from others to consider, so it restricts delivery of application messages to only those identified in its own latest cut. When a MBRSHP view v' is delivered to p, the v'.startId mapping tells p to use the synchronization messages sync_msg[q][v'.startId(q)] from q ∈ v'.set. The members of p's transitional set for view v' are those end-points q whose sync_msg[q][v'.startId(q)].view is the same as p's current view v. After receiving view v' from MBRSHP, p allows delivery of application messages identified by cuts in the synchronization messages from the processes that are already known to be members of the transitional set. The delivery of view$_p$(v', T) to p's application is enabled only after p has received the synchronization messages from all the potential members of T and after it has delivered all application messages committed to by the cuts of the members of T. Since all the end-points that move from v to v' use the same set of synchronization messages, the Virtual Synchrony and Transitional Set safety properties are satisfied.

End-point p is guaranteed to eventually receive all the application messages sent by the members of its transitional set T. However, p may fail to receive some of the application messages sent by disconnected end-points (not in T)

although certain cuts of members of T commit to deliver these messages. Such messages need to be forwarded to p by the members of T that have them. These members of T deduce from the p's cut that p lacks these messages and use a ForwardingStrategyPredicate to compute which of them have to forward which missing messages to p. We describe some of the many possible such predicates in [11].

### 5.2.2 Liveness of VS_RFIFO+TS

We show that, in a fair execution of VS_RFIFO+TS in which the same view v' is delivered to all its members as their last MBRSHP event, the three preconditions on the view$_p$(v', T$_p$) delivery are eventually satisfied for every p ∈ v'.set:

1. Condition v'.startId(p) = start_change.id remains true since by the assumption there are no subsequent start_change events at p.

2. End-point p eventually receives synchronization messages tagged with the "right" cid from everybody in v.set ∩ v'.set because they keep taking steps towards reliably sending these synchronization messages to p (by low-level fairness of the code) and because CO_RFIFO eventually delivers these messages to p (by the liveness assumption on CO_RFIFO).

3. End-point p eventually receives and delivers all the messages committed to in the cuts of the members of the transitional set T$_p$ because for each such message there is at least one end-point in T$_p$ that has the message in its msgs buffer and that would reliably forward it to p (according to the ForwardingStrategyPredicate) if so necessary. Also, p never delivers any messages beyond those committed to in the cuts of the members of T$_p$ because of the precondition on application message delivery.

### 5.2.3 Optimizations

Notice that end-point p does not need to send its current view and its cut to end-points which are not in current_view.set because p cannot be included in their transitional sets. Nevertheless, these end-points may wait to hear from p as p may still be in their current views. Therefore, in our algorithm, p sends synchronization messages to all the end-points in start_change.set. As an optimization, p could send a smaller synchronization message to processes in start_change.set − current_view.set, containing its start_change.id only (but neither a view nor a cut). The recipients of this message would know not to include p in their transitional sets for views v' with v'.startId(p) = p's start_change.id. When using this optimization, p also does not need to include its current view in the synchronization messages sent to current_view.set − start_change.set, since the view information can be deduced from p's view_msg.

AUTOMATON VS_RFIFO+TS$_p$  MODIFIES  WV_RFIFO$_p$

**Signature Extension:**
```
Input:   send_p(m)  modifies wv_rfifo.send_p(m)
         mbrshp.start_change_p(id, set), StartChangeId id, SetOf(Proc) set  new
         co_rfifo.deliver_q,p(m), Proc q, SyncMsg m   new

Output:  deliver_p(q, m)  modifies wv_rfifo.deliver_p(q, m)
         view_p(v, T), SetOf(Proc) T   modifies wv_rfifo.view_p(v)
         co_rfifo.reliable_p(set), SetOf(Proc) set  modifies wv_rfifo.co_rfifo.reliable_p(set)
         co_rfifo.send_p(set, m), SetOf(Proc) set, SyncMsg m  new
         co_rfifo.send_p(set, m)  modifies wv_rfifo.co_rfifo.send_p(set, m), FwdMsg m
```

**State Extension:**
```
(StartChangeId × SetOf(Proc))_⊥ start_change, initially ⊥
For all Proc q, ViewId id:  (View v, (Proc→Int) cut)_⊥ sync_msg[q][id], initially ⊥
SetOf(FwdMsg) forwarded_set, initially empty
```

**Transition Restriction:**
```
INPUT  mbrshp.start_change_p(id, set)
eff: start_change ← ⟨id, set⟩


OUTPUT  co_rfifo.reliable_p(set)
pre: start_change = ⊥ ⇒ set = current_view.set
     start_change ≠ ⊥ ⇒ set = current_view.set ∪ start_change.set


OUTPUT  co_rfifo.send_p(set, tag=sync_msg, cid, v, cut)
pre: start_change ≠ ⊥
     start_change.set ⊆ reliable_set
     ⟨cid, set⟩ = ⟨start_change.id, start_change.set - {p}⟩
     sync_msg[p][cid] = ⊥ ∧      v = current_view
     (∀ q ∈ current_view.set) cut(q) = LongestPrefixOf(msgs[q][v])
eff: sync_msg[p][cid] ← ⟨v, cut⟩


INPUT  co_rfifo.deliver_q,p(tag=sync_msg, cid, v, cut)
eff: sync_msg[q][cid] ← ⟨v, cut⟩


OUTPUT  deliver_p(q, m)
pre:  if (start_change ≠ ⊥ ∧ sync_msg[p][start_change.id] ≠ ⊥) then
         if start_change.id ≠ mbrshp_view.startId(p) then
            last_dlvrd[q]+1 ≤ sync_msg[p][start_change.id].cut(q)
         else  let S = {r ∈ mbrshp_view.set ∩ current_view.set |
                        sync_msg[r][mbrshp_view.startId(r)].view = current_view}
            last_dlvrd[q]+1 ≤ max_{r ∈ S} sync_msg[r][mbrshp_view.startId(r)].cut(q)


OUTPUT  view_p(v, T)
pre: v.startId(p) = start_change.id     // to prevent delivery of obsolete views
     (∀ q ∈ v.set ∩ current_view.set) sync_msg[q][v.startId(q)] ≠ ⊥
     T = {q ∈ v.set ∩ current_view.set | sync_msg[q][v.startId(q)].view = current_view}
     (∀ q ∈ current_view.set) last_dlvrd[q] = max_{r ∈ T} sync_msg[r][v.startId(r)].cut(q)
eff: start_change ← ⊥


OUTPUT  co_rfifo.send_p(set, tag=fwd_msg, r, v, m, i)
pre: (∀ q ∈ set) ⟨q, r, v, i⟩ ∉ forwarded_set
     ForwardStrategyPredicate(⟨set, r, v, i⟩, current_state)
eff: (∀ q ∈ set) add ⟨q, r, v, m, i⟩ to forwarded_set
```

**Figure 10. Virtually synchronous reliable FIFO multicast and transitional set end-point automaton.**

Another optimization can be used to minimize synchronization message sizes if we strengthen the membership specification to require a `MBRSHP.start_change` to be sent every time the membership changes its mind about the next view. In this case, the latest `MBRSHP.start_change` has the same membership as the delivered `MBRSHP.view`, and therefore the synchronization messages do not need to include information about messages delivered from end-points in `start_change.set` ∩ `current_view.set` because the synchronization message from each of these end-points can terminate a stream of application messages that this end-point would deliver in its current view.

## 5.3. Self Delivery

As a final step in constructing the automaton that models an end-point of our group communication service, $GCS_p$, we add support for Self Delivery to the $VS\_RFIFO+TS_p$ automaton presented above. Self Delivery requires each end-point to deliver to its application all the messages the application sends in a view, before moving on to the next view.

AUTOMATON $GCS_p = VS\_RFIFO+TS+SD_p$  MODIFIES $VS\_RFIFO+TS_p$

**Signature Extension :**
```
Input:   block_ok_p() new

Output:  block_p() new
         view_p(v,T) modifies vs_rfifo+ts.view_p(v,T)
         co_rfifo.send_p(set, m) modifies
            vs_rfifo+ts.co_rfifo.send_p(set,SyncMsg m)
```

**State Extension :**
```
block_status ∈ {unblocked, requested, blocked},
                  initially unblocked
```

**Transition Restriction :**
```
OUTPUT  block_p()
pre: start_change ≠ ⊥
     block_status = unblocked
eff: block_status ← requested

INPUT  block_ok_p()
eff: block_status ← blocked

OUTPUT  co_rfifo.send_p(set, tag=sync_msg, cid, v, cut)
pre: block_status = blocked

OUTPUT  view_p(v,T)
eff: block_status ← unblocked
```

**Figure 11. $GCS_p$ end-point automaton.**

In order to implement Self Delivery and Virtual Synchrony together in a live manner, each end-point must *block* its application from sending new messages while a view change is taking place (as proven in [8]). Therefore, we modify $VS\_RFIFO+TS_p$ to have an output action `block` and an input action `block_ok`, and we assume that the application at end-point p has the matching actions and that it eventually responds to every `block` request with a `block_ok`

response and subsequently refrains from sending messages until a `view` is delivered to it. In [11], we model this assumption with an abstract application automaton.

The $GCS_p$ automaton appears in Figure 11. After receiving the first `start_change` notification in a given view, end-point p issues a `block` request to the application and awaits receiving a `block_ok` response before sending a synchronization message to other members of `start_change.set`. The `cut` sent in the synchronization message commits to all the messages p received from its application in the current view.

Since the application is required to respond with `block_ok` and is then blocked from sending further messages, and since the p's cut commits to all the messages the application has sent in the current view, the set of messages agreed upon based on the `cuts` includes all of p's messages. Therefore, p delivers all these messages before moving on to a new view, and Self Delivery is satisfied. Since end-point p has its own messages on the `msgs[p][p]` queue, the modification does not affect the liveness property of $VS\_RFIFO+TS$. Finally, we note that due to the use of inheritance, the $GCS_p$ automaton satisfies all the properties we have specified in Secion 4.

## 6. Conclusions

We have constructed a virtually synchronous group multicast algorithm which exchanges one round of synchronization messages during reconfiguration, in parallel with the execution of a group membership algorithm. In contrast to previously suggested virtual synchrony algorithms, our algorithm does not require processes to conduct an additional communication round in order to pre-agree upon a globally unique identifier and does not impose restrictions on membership service's choice of views. We are not aware of any other algorithm that has both of these features.

These features are achieved by virtue of a simple yet powerful idea: Membership service issues a *locally* unique start-change identifier every time it has new information about client membership. The inclusion of such identifiers in views eliminates the need to tag clients' messages with a common (globally unique) identifier.

The start-change interface is an important aspect of the design of a client-server oriented group communication service which decouples membership maintenance from group multicast in order to provide scalable group membership services in WANs. Maestro [6] and the service of [17] also separate the maintenance of membership from group multicast. Unlike Maestro [6], in our design, the client does not wait for the membership to agree upon a globally unique identifier before starting the virtual synchrony algorithm, and the membership service does not wait for responses from clients asserting that virtual synchrony was achieved before delivering views. Unlike [17], our service does not impose restrictions on the membership service's choice of

views, thereby allowing its applications to benefit from Virtual Synchrony in more cases (as explained in Introduction).

In [11] we show that the service presented in Section 5 also provides meaningful and correct semantics in the environment where GCS end-points can crash and recover. In particular, it allows the recovered GCS end-points to continue running the algorithm under their original identity (in contrast e.g., to Isis [5] which requires recovered processes to assume new identities). Furthermore, GCS end-points do not need to store any information on stable storage.

Our service is implemented as part of a novel architecture for scalable group communication in WANs. After testing its current scalability limits, we intend to explore ways to improve the scallability further by incorporating a two-tier hierarchy into our algorithm, as suggested by Guo et al. [9]. With this approach, processes would send synchronization messages to their designated leaders, who would in turn exchange only the cumulative information among themselves. The framework in which we presented our algorithm allows us to incorporate extensions such as this one.

In [11] we formally prove the correctness of our algorithms. In particular, we prove the safety properties by defining simulation relations from the algorithm automata to the specification automata. The incremental way in which we have constructed our algorithms and specifications allows us to also construct the simulation proof incrementally. For example, in order to prove that VS_RFIFO+TS simulates VS_RFIFO+TS : SPEC we extend the simulation relation from WV_RFIFO to WV_RFIFO : SPEC and reason solely about the extension, without repeating the reasoning about the parent components. This reuse is justified by the Proof Extension theorem of [12].

## Acknowledgments

## References

[1] ACM. *Commun. ACM 39(4), special issue on Group Communications Systems*, April 1996.

[2] Y. Amir, L. E. Moser, P. M. Melliar-Smith, D. A. Agarwal, and P. Ciarfella. The Totem single-ring ordering and membership protocol. *ACM Trans. Comp. Syst.*, 13(4), Nov. '95.

[3] T. Anker, G. Chockler, D. Dolev, and I. Keidar. Scalable group membership services for novel applications. In M. Mavronicolas, M. Merritt, and N. Shavit, editors, *Networks in Distributed Computing*, vol. 45 of *DIMACS*, pp. 23–42. American Mathematical Society, '98.

[4] Ö. Babaoğlu, R. Davoli, and A. Montresor. Partitionalbe Group Membership: Specification and Algorithms. TR UBLCS97-1, Comp. Sci., Univ. of Bologna, Jan. '97.

[5] K. Birman. *Building Secure and Reliable Network Applications*. Manning, 1996.

[6] K. Birman, R. Friedman, M. Hayden, and I. Rhee. Middleware support for distributed multimedia and collaborative computing. In *Multimedia Computing and Networking (MMCN98)*, 1998.

[7] A. Fekete, N. Lynch, and A. Shvartsman. Specifying and using a partionable group communication service. In *16th ACM Symposium on Principles of Distributed Computing (PODC)*, pp. 53–62, August 1997.

[8] R. Friedman and R. van Renesse. Strong and Weak Virtual Synchrony in Horus. TR 95-1537, Dept. of Computer Science, Cornell University, August 1995.

[9] K. Guo, W. Vogels, and R. van Renesse. Structured virtual synchrony: Exploring the bounds of virtual synchronous group communication. In *7th ACM SIGOPS European Workshop*, September 1996.

[10] J. Hickey, N. Lynch, and R. van Renesse. Specifications and proofs for ensemble layers. In *5th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, LNCS, Mar. 1999.

[11] I. Keidar and R. Khazan. A client-server approach to virtually synchronous group multicast: Specifications, algorithms, and proofs. MIT-LCS-TR-794, MIT, Nov. 1999.

[12] I. Keidar, R. Khazan, N. Lynch, and A. Shvartsman. An inheritance-based technique for building simulation proofs incrementally. MIT-LCS-TR. 2000. In preparation.

[13] I. Keidar, J. Sussman, K. Marzullo, and D. Dolev. A Client-Server Oriented Algorithm for Virtually Synchronous Group Membership in WANs. In *20th Intern. Conf. on Distr. Comp. Syst. (ICDCS)*, Apr. 2000. Full version: MIT-LCS-TM-593.

[14] R. Khazan, A. Fekete, and N. Lynch. Multicast group communication as a base for a load-balancing replicated data service. In *12th International Symposium on DIStributed Computing (DISC)*, pp. 258–272, September 1998.

[15] N. Lynch. *Distributed Algorithms*. Morgan Kaufmann, '96.

[16] L. E. Moser, Y. Amir, P. M. Melliar-Smith, and D. A. Agarwal. Extended virtual synchrony. In *14th Intern. Conference on Distr. Comp. Systems (ICDCS)*, pp. 56–65, June 1994.

[17] A. Schiper and A. Ricciardi. Virtually synchronous communication based on a weak failure suspector. In *23rd IEEE Fault-Tolerant Comp. Symp. (FTCS)*, pp. 534–543, June '93.

[18] I. Shnaiderman. Implementation of Reliable Datagram Service in the LAN environment. Lab project. The Hebrew University of Jerusalem, January 1999. http://www.cs.huji.ac.il/~transis/publications.html.

[19] R. Vitenberg, I. Keidar, G. V. Chockler, and D. Dolev. Group Communication Specifications: A Comprehensive Study. TR CS99-31. The Hebrew Univ. of Jerusalem. Sept. 1999.