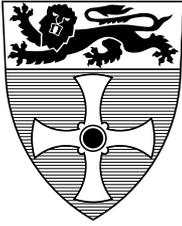


UNIVERSITY OF
NEWCASTLE



COMPUTING SCIENCE

Enhancing Replica Management Services to Cope with Group Failure

Paul D Ezhilchelvan and Santosh K Shrivastava

TECHNICAL REPORT SERIES

No. [CS-TR-665]

October 1999

Contact: Paul Ezhilchelvan
[Paul.Ezhilchelvan@ncl.ac.uk]

A short version of this paper appeared in the Proceedings of the Symposium on Object Oriented Realtime Computing (ISORC99), St Malo, April 99.

Enhancing Replica Management Services to Cope with Group Failures

Paul D Ezhilchelvan and Santosh K Shrivastava

Department of Computing Science
University of Newcastle upon Tyne
Newcastle upon Tyne, NE1 7RU, UK

Abstract

In a distributed system, replication of components, such as objects, is a well known way of achieving availability. For increased availability, crashed and disconnected components must be replaced by new components on available spare nodes. This replacement results in the membership of the replicated group 'walking' over a number of machines during system operation. In this context, we address the problem of reconfiguring a group after the group as an entity has failed. Such a failure is termed a group failure which, for example, can be the crash of every component in the group or the group being partitioned into minority islands. The solution assumes crash-proof storage, and eventual recovery of crashed nodes and healing of partitions. It guarantees that (i) the number of groups reconfigured after a group failure is never more than one, and (ii) the reconfigured group contains a majority of the components which were members of the group just before the group failure occurred, so that the loss of state information due to a group failure is minimal. Though the protocol is subject to blocking, it remains efficient in terms of communication rounds and use of stable store, during both normal operations and reconfiguration after a group failure.

Keywords – system availability, object groups, group failures, network partitions, membership views, membership services.

1. Introduction

In a distributed system, component replication (where a component computational entity such as a process, module) is a well known high availability. Equally well-known are the techniques for building services such as membership and message ordering services. We will consider the issue of enhancing the availability of a replica of failures, while preserving the strong consistency property where states of all replicas that are regarded as available be mutually consistent. A voting paradigm is an efficient way of achieving this end [Jajszczyk]. In a network failure partition the replica group into disjoint subgroups maintained only in the partition (if any) that contains a majority of the master subgroup), with the replicas in all other partitions being inactive. That is, the majority partition (if any) forms the master subgroup and provides services previously provided by the group; if a partition forms a majority of the current master subgroup from the rest, then this partition becomes the new master subgroup. Thus, each newly formed master subgroup contains a majority of the previously existed master sub-group. Replication with dynamic voting offers a better way of maintaining system availability than voting that requires a majority of all the members that initially remain connected. The following example illustrates dynamic voting:

Stage 0: Let the group configuration be initially $G_0 = \{C_1, C_2, C_3, \dots, C_n\}$ where C_i is the i^{th} component.

Stage 1: Say, a network partition splits G_0 into $G_1 = \{C_1, C_2, C_3, \dots, C_6, C_7\}$; G_1 now becomes the master subgroup and thereby the new group configuration.

Stage 2: Say, G_1 splits into $G_2 = \{C_1, C_2, C_3\}$ and $G'_2 = \{C_4, C_5\}$; G_2 is the master subgroup and thereby the third group configuration.

The above example indicates how the dynamic voting can preserve the group services even though the original group G_0 got split into islands having less than half the members of G_0 . Availability can be however maintained as long as the master subgroup exists after a failure. Suppose that after stage 2 G_2 detaches from other members. Now, no master subgroup exists and normal services can no longer be provided. We call this a group failure (short). Note that many combinations of failures can lead to a group failure. A group failure after stage 2 can be caused by simultaneous crashing of C_3 and detachment of C_2 and C_1 , and so on. When there are communication delays between components is not known with certainty, g-failures can occur even in the absence of any physical failure in the system. A burst of network traffic, for instance, increases the communication delays between connected components beyond what was considered to be likely, each component may falsely conclude that the other is not responding and hence must be disconnected. Therefore, g-failures should not be regarded as rare events if communication message delays cannot be estimated accurately.

Let us assume that the components have stable states which do not change on node crashes. Given that the component state survives node crashes

preferable to have the replica management service enhanced to cope instead of relying only on cold-start to resume the group services. To achieve this, we propose a configuration protocol that enables the master subgroup prior to a g-failure, to reconstruct the group once those members have recovered and got reconnected. Of course, to ensure that only one such group is formed. The protocol is implemented solely by the services used to build a replica group, in particular the configuration service. To illustrate this, let us continue on the above example below:

Stage 3: C3 crashes before it could record in its stable store that the master subgroup $G_2 = \{C1, C2, C3\}$ has been formed; the remaining members C1 and C2, record in their stable store that G_2 is the latest master subgroup. C3 disconnects from each other.

No master-subgroup now exists and a g-failure has occurred. Next, C3 recovers and reconnects with C4 and C5, and C2 reconnects to C1. C3 and C5 form the 'master subgroup' on the basis that its members form the last group configuration G_1 that is known to all of them, while C1 and C2 form the 'master subgroup' on the same basis that its members are a known configuration G_2 . Now, we have two live master subgroups. To prevent this from happening, we require that (i) a new master subgroup be considered formed only after a majority of the previous master subgroup have recorded in their stable store the composition of the new master subgroup (req1); and (ii) a subgroup constructed after a g-failure include at least a majority of the latest master subgroup formed prior to the g-failure (req2). Requirement req1 ensures that there can be only one group configuration that qualifies to be the latest master subgroup formed before a g-failure (and, in general, at any given time). Requirement req2 permits no more than one master subgroup to emerge after a failure.

We assume that the construction of the replica management system (including the voting component) can avail the use of a group membership service which is an operational component with an agreed set of components that are expected to be functioning and connected. For such a replica management system, the configuration management subsystem - the main contribution of this paper - provides (i) a group view installation service to enable members of a master subgroup to record group membership information on stable store; and (ii) a configuration service that makes use of these stable views to enable reconstruction of the master subgroup after a g-failure as soon as enough number of the components of the system have recovered and reconnected. A prototype version of the replica management service described here has been implemented [Black97]. Our server-side replica management system called Somersault [Murray97]. Our server-side Somersault by providing recovery from group failures.

The paper is structured as follows: section two introduces the system and some definitions and notations; it also specifies the two services of the configuration management subsystem, namely the view installation and the configuration services. The next two sections describe in detail how

provided. Section five compares and contrasts our work with the other published papers in this area, and concludes the chapter.

2. System Overview and Requirements

2.1. Assumptions and System Structure

It is assumed that a component's host node can crash but contain whose contents survive node crashes. Components communicate with passing messages over a network which is subject to transient or long We assume that a partition eventually heals and a crashed node eventually the bound on repair/recovery time is finite but not known with increased availability, we permit new components created on spare group, with no restriction on the number of such joining nodes and joining. Our system leaves to the administrator to decide how available spare nodes should be instructed to join the group, and spares instructed by the administrator are attempting to join the enables them to join with a guarantee that they could compute the component state from the existing members. For simplicity, we assume of a group do not voluntarily leave the group, but are only for crashes or partitions.

2.1.1. View Maker (VM) Subsystem

We assume that our replica management system has been constructed of the services provided by a group membership subsystem. This subsystem in the host node of an active component, say p , constructs membership where a view is the set of components currently believed to be connected to p . We call this subsystem the View Maker, or VM for short the VM of p as VM_p . In delivering the upto-date views constructed, to provide the abstraction of view synchrony or virtual synchrony model is assumed for the underlying communication subsystem. We refer [Babaoglu95, Babaoglu97] and [Schiper94] for a complete list of view synchronous and virtual synchronous abstractions, respectively highlight some of these properties that are considered important for

vs1: p is present in any view constructed by VM_p . (self-inclusion.)

vs2: a message m from another component q is delivered to p only constructed by VM_p prior to the delivery of m contains q . (view-message

vs3: the delivery of constructed views is synchronised with the delivery such that components receive identical set of messages between components that are identical. (view-message synchrony.)

vs4: If VM_p delivers a view v , then for every component q in v , either v or VM_p constructs consecutive view w that excludes q . (view agreement

There are many protocols in the literature which can be used to assume assumed VM subsystem; e.g., [Birman87, Ricciardi91, Mishra9 asynchronous system with the primary-partition assumption, [Moser96, Amir92, Ezhilchelvan95, Babaoglu95] for partitionable

systems. These protocols are not designed to cope with g-failure described below deals with g-failures using the services of the VM

2.1.2. Configuration Management (CM) Subsystem

On top of the VM service exists a configuration management (CM) (Figure 2). CM of component p , denoted as CM_p , carefully records information provided by VM_p in the local stable store. In a time management system, a new view decided by VM_p is usually delivered to our system, it reaches p via CM_p . VM_p regards CM_p as an application every new view it decides.

CM_p of member p essentially provides the following three functions:

(i) it considers each view delivered by VM_p and decides whether a g-failure occurred. If g-failure occurrence is ruled out, CM_p passes on that view. If certain conditions are met which ensure strong consistency,

(ii) if a new view delivered to p contains a spare node attempting to join, CM_p facilitates the spare node (in co-operation with CM of other members in the new view) to compute the most recent component state.

(iii) if a view constructed by VM_p indicates that a g-failure may have occurred, it executes a configuration protocol with CM of connected components. This ensures that if the group is reconfigured, it is the master subgroup that existed just before the g-failure was suspected to have occurred.

It must be emphasized here that CM_p can only suspect, not accurately detect, the occurrence of a g-failure when it inspects a new view from VM_p . Consider the disconnected component $C3$ in figure 1. With the membership being $\{C1, C2, C3\}$, when CM of $C3$, CM_3 , is delivered as a view $\{C3\}$, it cannot know whether the partition has split the group in two (as in Fig. 1(a)), or in two ways (as in Fig. 1(b)) to form the next master sub-group. So, in both cases, CM_3 would suspect a g-failure and execute the configuration protocol. In case of 3-way partition, $C3$ would execute the g-failure master subgroup with, say, $C1$, if it re-connects to $C1$ and executes the protocol. In the second case, when the partition has been 'walked over': $C1$ and $C2$ have formed the new master subgroup and $C3$ will then join the pool of spares. Note that it is also possible for $C3$ to 'walk over' in the first case: if the isolation of $C3$ lasts so long that it can re-connect to $C1$ and form the next master group. Thus, the outcome of reconfiguration attempts by components is decided by the pattern of components recovery and reconnection.



Figure 1. (a) 3-way partitioning

(b) 2-way partitioning

2.2. View Names within the System

Our replica management system (above the communication layer) is shown in layers as shown in fig. 2. Recall that CM_p delivers to p a view c only if certain conditions are met. That is, a view becomes more significant within the system. To reflect this, we call a view differently significant views constructed by VM_p are called the membership views or significant views. VM_p delivers $Mviews$ to CM_p via a queue called $ViewQ_p$ where $Mviews$ are in the order of delivery. CM_p deals with one $Mview$ at a time, as a $Mview$ reaches the head of $ViewQ_p$, which is denoted as $head_p$. CM_p records $head_p$ in the stable store as the new component view, provided a set of conditions is met. The component view of p is called $Cview_p$. Only $Cview_p$ is maintained by component p and provides p with the current membership view. As discussed earlier (see request of Section 1), making $head_p$ as $Cview_p$ involves several stages; $head_p$ is first recorded in stable store as the stabilised view and then installed as $Cview_p$. CM uses a view numbering scheme for numbering the view contents of $Sview_p$ and $Cview_p$.

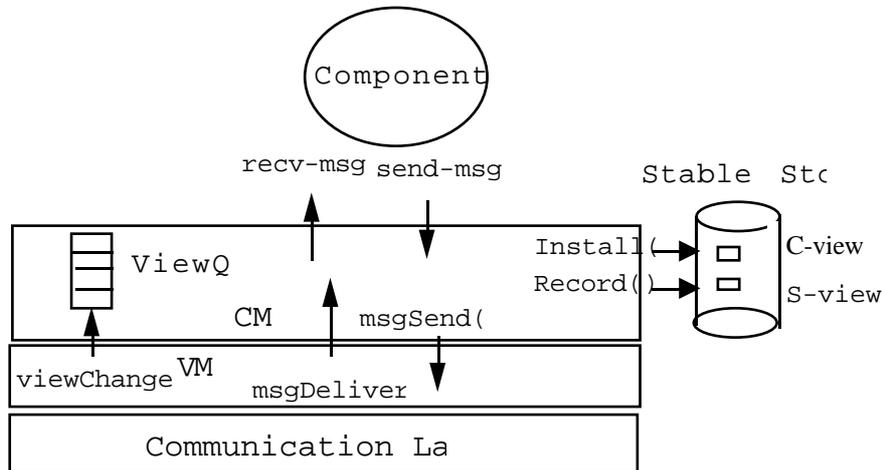


Figure 2. The system Architecture.

2.3. Notations and definitions

Each component p maintains three variables $status_p = (member, spare, waiting)$ (normal, reconfiguration, waiting, joining), and $view_number_p$ (an integer) in its stable store. In addition, it also maintains two view variables $Sview_p$ and $Cview_p$ initialised to null set, if p is spare. $Sview_p$ has a view number $view_number_p$. The view number of $Cview_p$ is indicated by $view_number_p$. $status_p$ is set to $member$ when p considers itself to be a member of the group, or to $spare$ if p is a spare member. p (with $status_p = member$) observes a g -failure and subsequently executes the configuration protocol, it sets its $mode_p$ to $reconfig$. $mode_p$ changes to $normal$ if p succeeds in becoming a member of the group; otherwise p becomes a spare setting $status_p = spare$ and $mode_p = waiting$. $mode_p$ of a spare component p can be either $waiting$ or $joining$; $mode_p = waiting$ means p is waiting to be informed by its VM_p that it has been connected to the group; once connected, p attempts to join the group by setting its $mode_p = joining$. If the join attempt by p succeeds, $status_p$ is set to $member$ and $mode_p$

The variable view-number_p is initialised to -1 at system start time (formed) and whenever p becomes a spare; it is incremented every time a new Cview .

We define the terms survivors and joiners for a pair of Mviews c of a component p . Let $\text{vu}_i, \text{vu}_{i+1}, \dots, \text{vu}_j, j \geq i+1$, be a sequence of VM_p in that order. The set $\text{survivors}(\text{vu}_i, \text{vu}_j)$ is the set of Mviews that survive from vu_i into every Mview constructed upto vu_j : $\text{survivors}(\text{vu}_i, \text{vu}_j) = \text{vu}_{i+1} \cap \dots \cap \text{vu}_j$. The term $\text{joiners}(\text{vu}_i, \text{vu}_j)$ will refer to the set of Mviews which are not in $\text{survivors}(\text{vu}_i, \text{vu}_j)$: $\text{joiners}(\text{vu}_i, \text{vu}_j) = \text{vu}_j - \text{survivors}(\text{vu}_i, \text{vu}_j)$. Finally, we define $\text{M_SETS}(g)$ for a set g of components as the set of all subsets s of g : $\text{M_SETS}(g) = \{s \mid s \subseteq g \wedge |s| > (|g|)/2\}$.

2.4. View Maintenance

When viewQ_p is non-empty, CM_p of member p checks for the occurrence of a g -failure by inspecting the contents of head_p , and by evaluating the condition $\text{survivors}(\text{Cview}_p, \text{head}_p) \in \text{M_SETS}(\text{Cview}_p)$. If this condition is not satisfied, a g -failure is assumed to have occurred. CM_p first sends an Abort message to all components in $\text{joiners}(\text{Cview}_p, \text{head}_p)$, informing the CM of any attempt at recording/installing head_p . We will denote this Abort message (which contains head_p) as $\text{AMsg}_p(\text{head}_p)$. CM_p then sets its variable $\text{recording_condition}$ and executes the configuration protocol to reconfigure itself. If the above condition is met, a copy of head_p is atomically recorded in stable_store as the new Sview_p with the view number = (view-number_p+1) , provided the recording conditions are satisfied. This Sview_p represents the potential next view. If the recording conditions are not met, the CM_p either concludes that a g -failure has occurred and proceeds to execute the configuration protocol to reconfigure itself, or dequeues head_p and proceeds to work with the next view (if any). The recording conditions, the need for them, and how they are satisfied are discussed in the next section.

The newly recorded Sview_p is regarded ready for becoming the next view. If the installation condition is satisfied (again, the need for this condition is verified will be discussed in the next section). In which case, CM_p updates the component view by replacing the current Cview_p by Sview_p , and CM_p discards head_p . The local stable store update operations are indicated by curly braces and are carried out atomically: $\{\text{Cview}_p := \text{Sview}_p; \text{view-number}_p := \text{view-number}_p + 1;\}$ If the installation condition is not met, a g -failure is assumed to have occurred and the configuration protocol is executed.

The view number of Cview_p is indicated in view-number_p . Since Cviews and Sviews are modified along with their view number as an atomic operation, there is at most one Cview_p and one Sview_p associated with a given view number, and view numbers increase monotonically. Further, CM_p installing the Sview_p (and updating Cview_p) can be interrupted only by its suspecting a g -failure; in the absence of a g -failure is suspected, CM_p will not record a new Sview_p until the current Sview_p is installed. Thus, in the absence of g -failure suspicions, either $\text{view-number}_p = \text{view-number}(\text{Cview}_p)$, or the view number of $\text{Sview}_p = \text{view-number}_p + 1$.

of $Cview_{p+1}$, the latter being true while the installation condition upon to be satisfied. Let $Vu_p(k)$ be the Sview or the Cview that view number k ; similarly, let $Vu_{p'}(k')$ be an Sview or a Cview that view number k' , where p and p' may be the same component or distinct say $Vu_{p'}(k')$ is later than $Vu_p(k)$, denoted as $Vu_{p'}(k') \gg Vu_p(k)$, if

2.5. Requirements of the CM subsystem

We now state the two requirements the CM subsystem must meet. The first one is concerned with the "normal service" period during which no g-failure occurs; the second one is concerned with group formation after a g-failure.

Existence of at most one master subgroup at any time is achieved by two components that install Cviews with identical view number, k . Let $Cview_p(k)$ denote the Cview that p installs with view number k . The predicate $installed_p(k)$ is true if p has installed $Cview_p(k)$. $!Cview_p(k)$ is true if $Cview_p(k)$ is unique, i.e., no component q can install a Cview that is different from $Cview_p(k)$:

$$!Cview_p(k) \Rightarrow \forall q: \neg installed_q(k) \vee Cview_q(k) = Cview_p(k).$$

Note that any view installed by a component must contain the instance of the group. So, if $Cview_p(k)$ is unique, then no component outside $Cview_p(k)$ can install a view with view number k ; so, there can be only one k th membership set, hence only one k th master subgroup.

During normal service period, the CM modules of components ensure that Cviews installed are sequentially numbered, and that the k th Cview installed is provided that $(k-1)$ th Cview installed is unique.

Formally, CM subsystem ensures:

Requirement 1:

$$\begin{aligned} \forall k > 0, installed_p(k) &\Rightarrow \exists p': installed_{p'}(k-1); \text{ and,} \\ \forall k > 0, installed_p(k) \wedge !Cview_{p'}(k-1) &\Rightarrow !Cview_p(k). \end{aligned}$$

Section 3 discusses how this requirement is met.

If we assume that $Cview(0)$ is unique when the group is initially formed, then above requirement is met, then there will exist a unique latest Cview. We define this latest view as the last Cview, or simply the last.

Requirement 2: following a g-failure, a set Σ of functioning and correct components with identical Cview, restart-view, should be formed as soon as possible with the following properties:

Uniqueness: $\Sigma \cap last \in M_SETS(last)$. If $last$ is unique before g-failure, then only one Σ that can contain a majority of the last.

Continuity: $restart_view \neq last \Rightarrow view_number(restart_view) = view_number(last)$. The sequentiality of CView numbering is preserved across g-failures. The sequentiality of CView numbering with g-failures is transformed into a view installation of di

Recording Condition 2 (rc2): It is to ensure that all joiners in $head_p$ and stored the component state and also recorded Sview which is tl We will suppose that after CM_j of joiner j has stored the comp recorded an Mview, say vu , as $Sview_j$, it multicasts a Recorded component in vu . This message contains the recorded view vu and $RMsg_j(vu)$. So, the second condition is that CM_p receive an $RMsg_j(h$ joiner j in $head_p$. Formally, $\forall j \in joiners(Cview_p, head_p): recd_p(RM$

If $rc1$ and $rc2$ are met, CM_p atomically records a copy of $head_p$ as with view number = view-number $_p+1$. It then multicasts an $RMsg_p($ components (including itself) in $head_p$. If $rc1$ is met but not r $head_p$ from $ViewQ_p$ but retains a copy to evaluate survivors($Cview_p$ next $head_p$. If $rc1$ is not met, CM_p proceeds to execute the reconf after setting $mode_p$ to reconfiguration. Since no joiner can send R first receiving at least $\Phi(Cview_p)$ State messages, it is not poss without $rc1$.

3.2. Recording Condition for a joining component

The recording condition is verified by CM_j of joiner j (with $mode_j$ as its $head_j$ - the first Mview in $ViewQ_j$ - is constructed by ' designed to become false if it is not possible for CM_j to receive of State messages from members in $head_j$. The design is made somewh the fact that when VM_j delivers an Mview it cannot indicate who i members and who else (except j itself) are joiners. VM_j can obtain only from VMs of member components. Recall that, as far as VM modul components are concerned, the local $Cview$ is transparent and is r variable used by a local application called CM (see figure 2). Mo met but not $rc2$, CM_p of member p dequeues $head_p$, and proceeds to next Mview in $ViewQ_p$; therefore, VM_p cannot even assume that wl Mview reaches the $head_p$, the Mview it delivered immediately befor have been installed as $Cview_p$. So, CM_j cannot rely on VM_j to indic members in $head_j$.

When CM_j does not know $Cview_p$ of member p in $head_j$, its attempt t can result in a deadlock if $head_j$ contains more than one joiner. I member p in $head_j$ crashes before sending the State or the Abort m joiner will wait for ever to receive State messages from other jo essential that CM_j first constructs a reference $Cview$ which can be place of $Cview_p$ of member p in $head_j$ until an $SMsg_p(head_j)$ is rece will contain a copy of $Cview_p$. This reference $Cview$ constructed : $head_j$ is denoted as $RefCview_j(head_j)$ and is initially set to hea discussions are for a given $head_j$, we will refer to $RefCview_j($ $RefCview_j$.) CM_j then sends a Join message to every component in he that it is a joiner. We denote this message as $JMsg_j(head_j)$. When $JMsg_j'(head_j)$, it removes the sender j' of that message from $RefCv$ receives an $SMsg(head_j)$, it irreversibly sets $RefCview_j$ to the Cvi message. No $JMsg(head_j)$ that is received after receiving the first

RefCview_j. The survivors and view-number contained in the received are noted in variables members_j and RefCviewNo_j, respectively.

Once RefCview_j is initialised to head_j, the recording condition starts upon to become true or false. (Verifying the recording condition is not to modifying or irreversibly setting RefCview_j.) This condition is given above for a member:

Recording Condition for joiner (rc_joiner): It verifies whether all distinct components sent their State messages. Formally,

$$|\{q \in \text{RefCview}_j : \text{recd}_j(\text{SMsg}_q(\text{head}_j))\}| \geq \Phi(\text{RefCview}_j).$$

If rc_joiner is met, CM_j atomically records a copy of head_p as its view number = RefCviewNo_j+1 and sets mode_j = joining. It then sends RMsg_j(head_j) to all components (including itself) in head_j. If rc_joiner is not met, an Abort message AMsg(head_j) is received, CM_j dequeues head_j from its queue and discards it.

Recall that CM_p multicasts AMsg_p(head_p) only if survivors(Cview_p, head_p) = M_SETS(Cview_p) when it starts to deal with head_p. So, it sends SMsg_p(head_p) or AMsg_p(head_p), not both, for a given head_p; hence CM_q will receive an AMsg_q(head_j) once rc_joiner is met. Otherwise, this would mean CM_q sent SMsg_p(head_j) without suspecting a g-failure at the start, which is not possible since CM_q has head_q = head_j and survivors(Cview_q, head_q) ≠ M_SETS(Cview_q). This in turn means that Cview_p and Cview_q are not identical which is a contradiction to the induction hypothesis.

To illustrate how certain failure cases that could lead to deadlock are avoided, consider the group {p,q,r} with unique Cview_p; i.e., Cview_p = Cview_q = Cview_r. Let the VM modules deliver an enhanced Mview such that head_p = head_q = head_r = {p,q,r,j,j1,j2,j3} = head_j, where j, j1, j2, and j3 are joiners. Since CM_p multicasts their State messages; if {j, j1, j2, j3} remain correct, CM_p eventually changes to {p, q, r} from its initial value of head_p. CM_c will receive recd_j(SMsg_c(head_j)) will become false for crashed c = p, q, and r since CM_c has not received head_j that rc_joiner cannot be met.

3.3. Installation Conditions

Having recorded head_p as Sview_p, CM_p installs the Sview_p as the new view after verifying that a majority of the existing Cview_p have recorded head_p.

Installation condition (ic): $\{q \in \text{survivors}(\text{Cview}_p, \text{head}_p) : \text{recd}_p(\text{SMsg}_q(\text{head}_p)) = \text{M_SETS}(\text{Cview}_p)\}.$

The CM_j of a joiner j has two installation conditions. The first one is that all joiners of head_j have recorded head_j; the second one is the same as above for member p. Note that CM_j has recorded head_j means that it has received RMsg_j messages from some member p in head_j; so, RefCview_j = Cview_p and survivors(Cview_p, head_p).

Installation condition 1 for joiner (ic1_joiner):

$$\forall c \in \text{head}_j - \text{members}_j: \text{recd}_j(\text{RMsg}_c(\text{head}_j)).$$

Installation condition 2 for joiner (ic2_joiner):

$$\{p \in \text{members}_j: \text{recd}_j(\text{RMsg}_p(\text{head}_j))\} \in \text{M_SETS}(\text{RefCview}_j).$$

If both conditions are met CM_j makes component j a member by atomic $\{\text{Cview}_j := \text{Sview}_j; \text{view-number}_j := \text{RefCviewNo}_j + 1; \text{status}_j = \text{member normal}; \}$. The head_j is then dequeued and discarded. If the first no member p in head_j would have recorded head_j ; so, CM_j 's record is undone by atomically executing: $\{\text{Sview}_j = \text{null}; \text{mode}_j = \text{waiting}; \}$. head_j is then dequeued and discarded. If only the second condition is not met, CM_j goes into a state of reconfiguration and executes the reconfiguration protocol. Observe that mode_j is set to reconfiguration, Cview_j and view-number_j remain unchanged from their initial values which are null and -1 respectively.

3.4. Correctness and Liveness

Correctness: Suppose that CM_p installs head_p as the new Cview_p with view-number $(k+1)$. The majority requirement in the installation condition ensures that a majority of $\text{Cview}_p(k)$ have recorded head_p as their Sview with view-number k . The recording condition (rc2) ensures that every CM of joiner q ($\text{Cview}_q(k+1)$) has also recorded head_p as its Sview with view-number k . CM_q records a new Sview before the existing one is installed. Therefore, $\text{Cview}_p(k)$ is unique, if CM_q of a survivor or joiner q installs $\text{Cview}_q(k+1) = \text{Cview}_p(k+1)$. This means that $\text{Cview}_p(k+1)$ is also unique.

Liveness: CM_p verifying the recording/installation condition requires the evaluation of the predicate $\text{recd}_p(m_q)$ which in turn involves checking whether message m_q has been/can be received from CM_q . Since the node of q is not in Mview_p , m_q can be sent, the evaluation of $\text{recd}_p(m_q)$ must involve checking whether q continues to be present in the subsequent Mviews constructed by VM_p . In this mind, we present an algorithm for evaluating $\text{recd}_p(m_q)$ which does not terminate indefinitely.

Figure 3 shows the ViewQ 's of CM_p and CM_q which, for simplicity, are assumed to be identical. We will also assume that $\text{Cview}_p = \text{Cview}_q = \{p, q, r1, r2, j\}$, $\text{view-number}_p = \text{view-number}_q = k$ (say). That is, $\text{Cview}_p(k)$ is unique. Let Mview_p and Mview_q be constructed by VM_p and VM_q as: $\text{vu1} = \{p, q, r1, r2, j\}$, $\text{vu2} = \{p, q, r2, j\}$, $\text{vu3} = \{p, q, j\}$. vu1 indicates the disconnection of member $r3$ (from Mview_p) on the inclusion of a new component j , vu2 the disconnection of $r2$ from Mview_p , and vu3 the disconnection of $r1$ from Mview_p .

We define $\text{List}_p(\text{Mview})$ as the set of messages which VM_p intends to send between the delivery of Mview and the delivery of the immediate subsequent Mview . $\text{List}_p(\text{vu3})$ is shown to be open and will remain so until a subsequent Mview where vu3 is constructed. $\text{List}_p(\text{vu1})$ and $\text{List}_p(\text{vu2})$, on the other hand, are closed to indicate that no received message can enter these lists any longer. By the message synchrony property of the VM subsystem (see §2.1.1), $\text{List}_p(\text{vu1}) = \text{List}_q(\text{vu1})$, and $\text{List}_p(\text{vu2}) = \text{List}_q(\text{vu2})$.

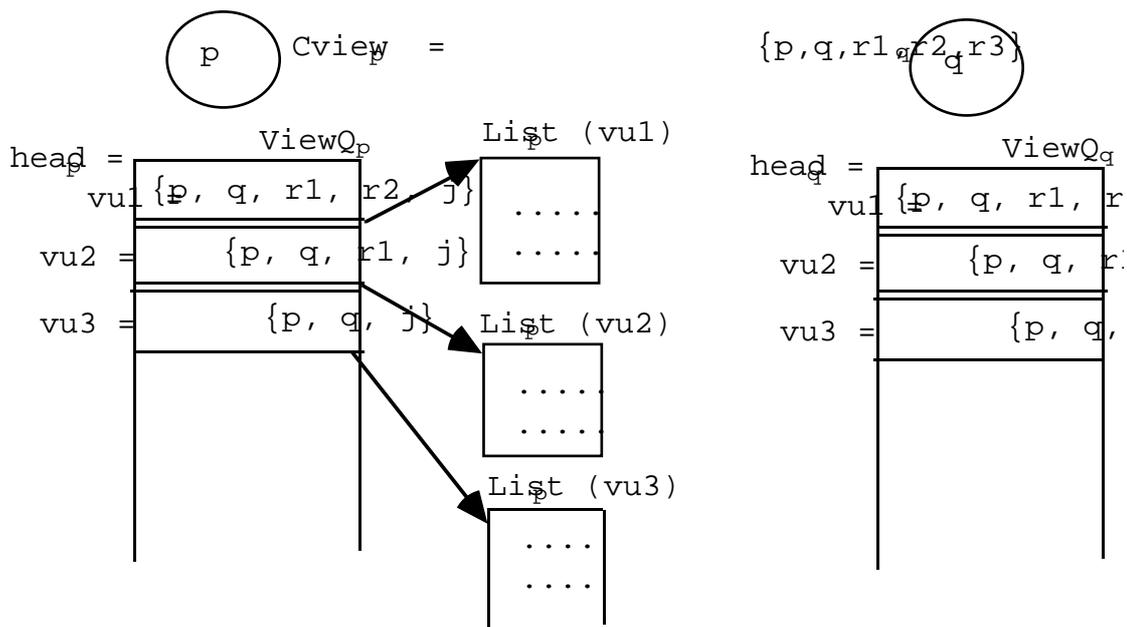


Figure 3. Closed and Open lists of messages delivered by VM after

The algorithm for evaluating $recd_p(m_q)$ is as follows: CM_p waits following two conditions to become true.

Evaluation condition 1 (ec1): $\exists vu \in ViewQ_p: m_q \in List_p(vu) \wedge (vu \in survivors(head_p, vu))$.

Evaluation condition 2 (ec2): $\exists vu \in ViewQ_p: q \notin vu$.

The condition ec1 is true when m_q is present in $List_p(vu)$ for $s \in ViewQ_p$ and q is present in all the views VM_p constructed from head p this vu ; ec2 becomes true when VM_p constructs an $Mview$ without q .

boolean $recd_p(m_q)$

{wait until ec1 \vee ec2; if ec1 then return true else return false}

Recall that CM_p evaluates $recd_p(m_q)$ only for such $q \in head_p$. Suppose CM_p constructs an $Mview$ vu that does not contain q (ec2). By message property of VM, the expected message from q cannot be in $List_p(vu')$ for vu' constructed prior to vu , is closed. If none contains m_q (not ec1), then q crashed or disconnected before $recd_p(m_q)$ is evaluated to be false.

3.5. Examples

We illustrate the working of the view recording and installation examples. The first one is based on Figure 3. We assume that $Cview_number_p = k$; also that p , q , and j remain connected and function VM_j also constructs vu_1 , vu_2 , and vu_3 as shown in the figure. Let $\lfloor (|CV_u|/2) \rfloor + 1$.

Suppose that r_1 crashes after multicasting its State message $SMsg$ (message $RMsg(vu_1)$). CM_c , $c = p, q$, or j , will find their respective installation conditions being met, and install $vu_1 = \{p, q, r_1, Cview\}$. Following the installation of vu_1 , CM_c delivers messages i which will be identical for every c . When CM_c has $head_c = vu_2$, suspected, as $head_c$ contains a majority of components in the current view. Since $head_c$ has no joiner, no recording condition needs to be met. As long as CM_c remains connected and functioning, CM_c will find the installation conditions for vu_2 being met, and install vu_2 as the $(k+1)$ th $Cview$, and deliver messages in $List_c(vu_2)$. CM_c will install vu_3 as the $(k+3)$ th $Cview$ and deliver messages in $List_c(vu_3)$. This shows that when VM_c and $VM_{c'}$ construct an identical sequence of $Cviews$, VM_c and $VM_{c'}$ behave identically; they also deliver an identical set of messages for two consecutive $Cviews$ they install.

Suppose that r_1 and r_2 crash before multicasting their State messages. CM_c , $c = p, q$, or j , will find the recording conditions not being met. CM_p and CM_q will proceed to execute the configuration protocol, while CM_j remains with no change in its status (= spare) and mode (= waiting).

Example with Concurrent Mviews

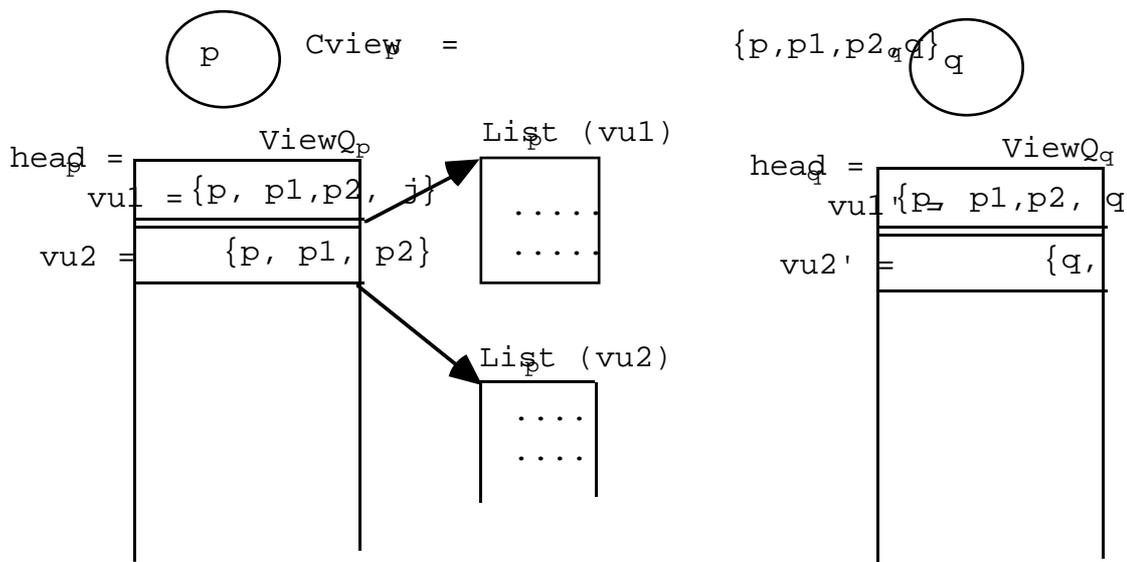


Figure 4. Concurrent and overlapping head views.

The second example is based on Figure 4 and illustrates scenarios where member p to dequeue $head_p$ without recording it, and CM_j of joiner configuration protocol with $Cview_j = null$. As in the previous example, assume that $Cview_p$ is unique and $view-number_p = k$. The figure shows p attempting to join the group $\{p, p_1, p_2, q\}$, and VM_p and VM_q reach view agreement due to the subsequent detachment of $\{p, p_1, p_2\}$ from the group. Suppose that p, p_1 , and p_2 remain connected to each other, and so CM_p and CM_q denote p, p_1 , or p_2 , and c' denote q or j . Every VM_c constructs a view vu_c for p in the figure, and every $VM_{c'}$ constructs $vu_{c'}$ and vu_2' show that VM_c and $VM_{c'}$ have constructed non-identical, overlapping vu_1 and vu_2' .

subsequently construct vu_2 and vu_2' respectively, due to the property (see subsection 2.1.1). Let $\Phi(CVu)$ be defined to be 1, compute the component state by receiving a single State message.

When CM_C has $head_C = vu_1$, it will find rc_1 met; but it will not receive from CM_j and will find rc_2 not being met. Hence, it dequeues head messages in $List_C(vu_1)$ when $Cview_C$ is still $\{p, p_1, p_2, q\}$. Note that CM_C does not contain any application message from j as j does not yet connect to C . Thus, the view-message integrity property (of subsection 2.1.1) is satisfied by CM_C . Since $\Phi(CVu) = 1$, both CM_q and CM_j will find for $head_C$ recording conditions being met, but not the installation condition to execute the reconfiguration protocol, with $Cview_q$ unchanged, $Sview_j$, and $Cview_j = null$.

4. Reconfiguration after a g-failure

The configuration protocol presented here meets Requirement 2 following a g-failure, a unique set of functioning and connected components to become the (first post-g-failure) master subgroup. These component operations with an identical $Cview$ called the restart-view. The restart-view is taken to be either the last view or an $Sview$ SV_u installed last, the latter being the case if and only if a majority of last components installed SV_u before the g-failure occurred. This invariant qualifies the restart-view and ensures the continuity in the numbering of $Cviews$ despite a g-failure. The master subgroup is guaranteed unique by ensuring that it contains a majority of components. The rationale behind the formation of the master subgroup is briefly described in the following section.

4.1. The Rationale

Let R be a set of components that get connected after a g-failure. Following reconfiguration, it needs to be determined whether a subset of components in R become the master subgroup. Let $Sview_r$ and $Cview_r$ denote the $Sview$ and $Cview$ of a component r in R , respectively. (If r is recovering from a crash, $Sview_r$ and $Cview_r$ are restored from its stable store.) Let $presumed_last$ be the latest $Cview$ of all r in R : for every r in R , either $presumed_last = Cview_r$ or $presumed_last \gg Cview_r$. By definition, $presumed_last$ is either the last $Cview$ installed prior to the last.

Let us consider the $Sviews$ recorded by the components of $presumed_last$ (those in $presumed_last \cap R$). One of the following three (mutually exclusive) situations must exist:

- (i) A majority of $presumed_last$ components recorded (at some time) an identical $Sview$ that is later than the $presumed_last$.
- (ii) A majority of $presumed_last$ components never recorded an $Sview$ later than the $presumed_last$.

¹No last component could have installed SV_u ; otherwise the installed version of SV_u which, by definition, is the latest $Cview$ installed by a component.

(iii) Neither 1 nor 2. That is, the number of presumed_last components in Sview that is later than the presumed_last, is at most $\lfloor \text{presumed_last} \rfloor$ and the number of presumed_last components that never recorded an Sview later than the presumed_last, is at most $\lfloor \text{presumed_last} \rfloor / 2$.

Let us first consider situation (1). Suppose that R contains $\lfloor \text{presumed_last} \rfloor / 2$ components with Sview = presumed_last+ and (a) a majority subset of presumed_last+, and (b) a majority subset of presumed_last+. We claim that (a) and (b) are met, restart-view = presumed_last+. The (simple) proof is:

Proof: Suppose that (a) and (b) are met but presumed_last+ \neq restart-view. Then (a) implies that a majority of components in presumed_last+ have installed presumed_last+ as their Cview. By the definition of restart-view, if restart-view \neq presumed_last+, then restart-view \gg presumed_last+. For this to be true, a majority of components in presumed_last+ must have installed presumed_last+ as their Cview before they have proceeded to record an Sview presumed_last++ (say), presumed_last+. None of these components that installed presumed_last+ as their Cview, can be in R, as per the way presumed_last is computed. This cannot be true - a contradiction.

Thus, when (a) and (b) are met, presumed_last+ becomes the restart-view. The components in $R \cap \text{presumed_last+}$ consider themselves to be the master subgroup.

In both situations (2) and (3), a majority of presumed_last have installed a progressive Sview that is later than presumed_last; therefore presumed_last+ is the restart-view. To deduce the existence of (2) and (3) R must contain all presumed_last components.

Observe that deducing which one of the three situations exists, requires at least a majority of presumed_last components with appropriate Sviews in the third situation. So, it is possible that a given R does not contain a majority of presumed_last components. In that case, the attempt to form the master subgroup with R is unsuccessful. Recovery and reconnection of more number of components need to be attempted.

4.2. The Protocol

The protocol is made up of five steps:

Step 0. CM_p sets mode_p to reconfiguration and waits for p to be joined by other components, i.e., for viewQ_p to become non-empty. Say, ViewQ_p is non-empty and $R = \text{head}_p$. (Note: the first Mview in ViewQ_p is only not dequeued.) The remaining four steps are done using R.

Step 1. Send {Sview_p, Cview_p} to every r in R (including itself);
Receive {SviewRecd_r, CviewRecd_r} from every r in R;

Step 2. Determine the presumed_last to be the latest CviewRecd_r;

Step 3. Determine the restart-view if possible; if not possible (i.e., if R is not a majority of presumed_last), discard R and go to step 0.

Step 4. components of $R \cap \text{restart-view}$:
install restart-view and resume group services;

components of R - restart-view:
become spares;

Each step is described in detail in the subsections below.

4.2.1. Step 1: View Exchange

CM_p multicasts a message $msg(Sview_p, Cview_p, mode_p)$ containing $Sview_p$ and $mode_p$. It then evaluates the predicate $recd_p(msg_r(SviewRecd_r, CviewRecd_r, mode_r))$ for every $r \in R$. If this predicate is true for a given r , then p becomes a spare for r whether $CviewRecd_r \gg Cview_p$ and $mode_r = normal$. If this condition is not met, an exception Walked-Over is raised indicating that p has been slow in which time the group is reconfigured without p . This exception is then working with R is given up: terminate the execution with R , $viewQ_p$, and go to step 0. The pseudo-code for step 1 is given below.

```

multicast msg(Sview_p, Cview_p, mode_p) to all r in R;
evaluate recd_p(msg_r(SviewRecd_r, CviewRecd_r, mode_r)) for every r in R;
catch (Walked-Over): { write atomically: {Sview_p := null; state_p := waiting;
                                     mode_p := waiting; }
                    exit; }
if ( $\exists r \in R: \neg recd_p(msg_r(SviewRecd_r, CviewRecd_r, mode_r))$ )
    then {give up on R;}

```

4.2.2. Step 2: Determine presumed_last

$presumed_last$ is computed to be the latest non-null view among the R . If a majority of $presumed_last$ is not in R , then the execution with R is given up.

```

{ presumed_last := CviewRecd_r of some r in R: (presumed_last != null)
  (forall r' in R: presumed_last = CviewRecd_r' > presumed_last)
  if (|presumed_last ∩ R| ≤ (|presumed_last|/2)) then { give up on R; }
}

```

Note that by requiring that $presumed_last$ be non-null, an R of only nodes with $mode = waiting$ are prohibited from forming the master subgroup.

4.2.3. Step 3: Attempt to Determine restart-view

CM_p divides the components in $presumed_last \cap R$ into non-overlapping candidate sets and denoted as CS_v , $v \geq 0$, based on the component $presumed_last+1$, $i \geq 1$, be an $Sview^2$ that is later than $presumed_last$.

² $presumed_last+1$ need not be unique after a g-failure; different $presumed_last$ can be recorded for different progressive $Sviews$, due to their VM modules concluding view agreements at different points. Let, for example, $last = \{1,2,3,4,5\}$. Let $C5$ crash and VM of $C4$ reach agreement on $\{1,2,3,4\}$. If VMs of $C1$, $C2$, and $C3$ suspect $C4$ before they reach agreement on $\{1,2,3\}$, they reach agreement straightaway on $\{1,2,3\}$. If a g-failure occurs after CMs have recorded the $last$ as $\{1,2,3,4\}$ (say, $presumed_last+1$), and $Sview_1 = Sview_2 = Sview_3 = \{1,2,3\}$ (say, $presumed_last+1$).

contains the components in $\text{presumed_last} \cap R$ whose S_{view} is $\text{presumed_last} + v$; CS_0 contains those components in $\text{presumed_last} \cap R$ whose S_{view} is not later than presumed_last . The code for this is given below.

```
(1)  if ( $\exists CS_v: v \geq 1 \wedge CS_v \in M\_SETS(\text{presumed\_last})$ 
       $\wedge \text{presumed\_last} + v \cap R \in M\_SETS(\text{presumed\_last} + v)$  ) then
      {restart-view := presumed_last + v;}
      // existence of situation (1) is deduced

(2)  else if ( $CS_0 \in M\_SETS(\text{presumed\_last})$ ) then
      {restart-view := presumed_last;}
      // existence of situation (2) is deduced

(3)  else if ( $\text{presumed\_last} \subseteq R$ ) then
      {restart-view := presumed_last;}
      // existence of situation (3) is deduced

else {give up on R;}
```

4.2.4. Step 4: Commencing Group Operations

Any p that is not in the restart-view becomes a spare, otherwise it becomes a member in its stable store. The pseudo code is as follows:

```
{  if ( $p \notin \text{restart-view}$ ) then { write atomically:
      { $S_{\text{view}_p}, C_{\text{view}_p} := \text{null}; \text{status}_p := \text{spare}$ 
         $\text{mode}_p := \text{waiting}; \text{view-number}_p := -1;$ 
        exit; }
    write atomically:
      { $\text{view-number}_p := \text{view-number}(\text{restart-view});$ 
         $S_{\text{view}_p}, C_{\text{view}_p} := \text{restart-view}; \text{status}_p := \text{member};$ 
         $\text{mode}_p := \text{normal};$  }
  }
```

4.2 Examples

We explain the working of the protocol with the help of examples and the evolution of S_{views} depicted in Table 1. For simplicity, assume that the changes considered in this discussion are caused by node crashes only, and they occur only when the group is being reconfigured after a g -failure.

Table 1 depicts a possible sequence of membership changes for a group. We adopt the following style to represent the state of the view instances. A component in normal font indicates that it has been installed as the C_{view_p} ; a component that has not been installed is written in mixed fonts: survivors (from the current S_{view_p}) in normal font, joiners in italics and excluded components (in the current C_{view_p} but not in the S_{view_p}) in bold. The S_{view_p} indicates its view-number.

Stage No	Sview of C1, C2	Sview of C3	Sview of C4, C5	Sview of C6, C7	Description
1	$\{1, 2, 3, 4\}^0$	$\{5, 2, 3, 4\}^0$	$\{1, 2, 3, 4, 5\}^0$	---	group initialised, n=5 C6 and C7 are spares
2	$\{1, 2, 3, 5\}^1$	$\{1, 2, 3, 4\}^0$	$\{1, 2, 3, 4\}^0$	---	C4 and C5 crash; CM1 & CM2 record their exclusive view first
3	$\{1, 2, 3, 5\}^1$	$\{1, 2, 3, 5\}^1$	$\{1, 2, 3, 4\}^0$	---	Slow CM3 records Sview(0)
4	$\{1, 2, 3\}^1$	$\{1, 2, 3, 5\}^1$	$\{1, 2, 3, 4\}^0$	---	CM1 & CM2 install Sview(1)
5	$\{1, 2, 3\}^1$	$\{1, 2, 3\}^1$	$\{1, 2, 3, 4\}^0$	---	CM3 install its Sview(0)
6	$\{1, 2, 3, 6, 7\}^2$	$\{1, 2, 3, 6, 7\}^2$	$\{1, 2, 3, 4\}^0$	$\{1, 2, 3, 6, 7\}^2$	C6 & C7 join; all active CM record Sview(2) = $\{1, 2, 3, 6, 7\}$
7	$\{1, 2, 3, 6, 7\}^2$	$\{1, 2, 3, 6\}^2$	$\{1, 2, 3, 4\}^0$	$\{1, 2, 3, 6\}^2$	CM3, CM6, and CM7 install Sview(2)
8	$\{1, 2, 3, 6, 7\}^2$	$\{3, 6, 7\}^3$	$\{1, 2, 3, 4\}^0$	$\{3, 6, 7\}^3$	C1, C2 crash before installing Sview(2); CM3, CM6, CM7 record and then install Sview(2)

Table 1. An evolution of Sviews.

The group is initially formed with $\{C1, C2, C3, C4, C5\}$. At the end of stage 1, each member has $\{1, 2, 3, 4, 5\}^0$ as its (initial) Sview in stable store; this is the end of stage 1. At the end of stage 2, CM₁ and CM₂ have recorded Sview(1) which cannot be installed now as $\{1, 2\} \notin M_SETS(Sview(0))$. The situation changes after stage 3 and CM₁ and CM₂ install Sview(1) in stage 4. In stage 6, the spares C6 and C7 join; CM₁, CM₂, CM₃, CM₆, and CM₇ record the Sview $\{1, 2, 3, 6, 7\}$. In the end of stage 6, C3, C6, and C7 install the Sview, since all components in the old Sview are known to have recorded $\{1, 2, 3, 6, 7\}$. But C1 and C2 crash before they can record their view. In stage 8, C3, C6, and C7 install the Cview with Sview(2). According to Table 1, $\{1, 2, 3, 6, 7\} \gg \{1, 2, 3\}$ and $\{1, 2, 3, 6, 7\} \gg \{1, 2, 3, 4, 5\}^0$ until stage 3, $\{1, 2, 3\}^1$ in stages 4, 5, and 6, and $\{3, 6, 7\}^3$ in stage 8.

Example 0: This example shows that the protocol is safe in not allowing one master subgroup to be formed after a g-failure. Let C1 and C2 be connected after stage 8 but remain partitioned from other components of R. Both C1 and C2 will estimate `presumed_last` to be $\{1, 2, 3\}^1$. Since C3 is already functioning as the group, another master subgroup should not emerge from R even though R contains a majority subset of processes. Components of R will find that they have an identical (proposed) Sview $\{1, 2, 3, 6, 7\}^2 \gg \text{presumed_last}$, and R does not contain a majority subset of $\{1, 2, 3, 6, 7\}^2$. So, none of the conditions in Step 3 of the protocol will be given up.

Example 1: This exemplifies the behaviour of the protocol under the conditions mentioned in section 4.1. Suppose that a g-failure occurs immediately after stage 8.

Here, last is $\{1,2,3\}^1$ and all of the last components have recorded restart-view is $\{1,2,3,6,7\}^2$. Say, $R = \{C1, C2, C4, C6\}$. By step each component r in R determines restart-view to be $\{1,2,3,6,7\}^2$. If restart-view, $C4$ will exit the protocol and join the pool of spare restart-view) install restart-view as their Cview and resume normal. Note that the view R is still the at head of every $ViewQ_r$, $r \in ViewQ$ not empty, CMs of $(R \cap restart-view)$ will execute the view protocol as members and CM of $C4$ as a joiner. Assuming no further disconnections, $C1, C2, C4,$ and $C6$ will get install R as the Cview.

Example 2 considers situation 2 where a majority of presumed-last is an Sview that is later than presumed-last. Say, a g-failure occurs last = $\{1,2,3,4,5\}^0$. Since no last component has recorded a later S also $\{1,2,3,4,5\}^0$; further, since all last components have identical contains any three (majority subset) of the last members will form a subgroup. Say, $R = \{C3, C4, C5\}$. Assuming that R remains connected CM of R computes presumed-last to be $\{1,2,3,4,5\}^0$, i.e., last its of R forms $CS_0 = R$ and decides in step 3.2 of the protocol the presumed-last = $\{1,2,3,4,5\}^0$. After (re-)installing restart-view of $C3, C4,$ and $C5$ subsequently install R as their next Cview = $\{3,4,5\}^1$. Say, after CMs of $C3, C4,$ and $C5$ have installed $\{3, 4, 5\}^1$ in the $C1$ and $C2$ recover and reconnect with $C3, C4,$ and $C5$. While CMs of execute the reconfiguration protocol with $R = \{1,2,3,4,5\}$, CMs of will execute view installation protocol for the delivered view $\{1, C1$ and $C2$ are regarded as joiners. This conflict gets resolved eventually $C4,$ and $C5$ expect CMs of $C1$ and $C2$ to send recorded messages but messages of configuration protocol. They would then respond by sending and Cview to CM_1 and CM_2 which would get Walked_Over exception, spares, and then start executing the view installation protocol as of their $ViewQ$ (still) having $R = \{1,2,3,4,5\}$.

In example 3, we illustrate the need for R to contain all the last circumstances. Let a g-failure occur at the end of stage 2. The $\{1,2,3,4,5\}^0$ which is also the Cview of every member component. We $= \{C1, C2, C4, C5\}$. The presumed_last is the same as last = $\{1, component$ of R knows that a minority of presumed_last (i.e. two) is an Sview that is later than presumed_last; and also that only presumed_last (i.e. two again) are known to have recorded an Sview is later than presumed_last. When Sview of $C3$ is $\{1,2,3,4,5\}^0$ (a restart-view becomes $\{1,2,3,4,5\}^0$. If Sview of $C3$ had been $\{1,2,3,4$ end of Stage 3), then restart-view becomes $\{1,2,3\}^1$. Hence determining view requires that the components of R know the Sview of $C3$. Here, step 3.3 of the protocol which requires R to contain presumed_last situation 1 nor 2 is known to prevail.

5. Related Work

Using the traditional, 2-Phase Commit (2PC) protocol [Gray78] updating membership-related information, [Jajodia90] maintains distinguished partition in a replicated database system. Our CM s

a variation of this traditional 2PC for Cview installation and the by our requirements and efficiency. In the traditional 2PC way of the coordinator - a deterministically chosen member in the new Cview the second (view-installation) phase after learning during the first component of the new Cview has recorded the view. Note that while is in progress, delivery of application messages is put on hold synchrony. Since we only require that at least a majority subset, the current Cview install the next Cview, we can speed up the view having the coordinator initiate the second phase as soon as a majority Cview and all joiners (if any) in the new Cview have recorded the soon as the installation conditions of section 3.3 are met. Further based execution of traditional 2PC are susceptible to co-ordination eliminate this weakness by executing our version of 2PC in a decentralized where every component checks installation conditions.

Since we use a 2PC protocol for view installations, the configuration be non-blocking. This blocking can be removed by using a 3PC protocol. The protocol of [Dolev97] employs the principles of an extended [Keid95] and builds a unique master subgroup after a g-failure. No architecture is remarkably similar to theirs. They differ from our major aspect: a component can have, and may have to exchange, more so more stable information needs to be maintained and message size. Obviously these features of [Dolev97] increase the overhead of the advantage, on the other hand, is that a reconnected set need only majority subset of last, never all last components as we would require when a g-failure occurs during view update (see example 3 of section

The primary partition membership service in [Birman87, Ricciardi91], the assumption that a majority of components in the Cview do not fail and that a functioning component is rarely detected as failed. This does not hold true during periods of network instability caused for network traffic or network congestion. This instability can lead to incorrect which in turn can lead to g-failures. In these circumstances, our [Dolev97]) can provide recovery from g-failures once the network traffic

[Chandra96] establishes the weakest failure detector (denoted as $\diamond S$) the consensus problem. Using this consensus protocol, a primary membership service is designed [Malloth95] and implemented [Feen membership service can construct a totally ordered sequence of views of each view surviving into the next view. It blocks from delivering the periods of g-failures (i.e., when a connected majority does blocking is released as soon as the requirements of $\diamond S$ are realised based membership service provides recovery from g-failures for the requirements of $\diamond S$? The answer appears to be no. The first view which membership service constructs after a g-failure, is what we call requirement 2 in §2.5). Determining the restart-view does not need the new master subgroup exists to restart the group services. To see the following example. Let the current view be $\{p, q, r, s, t\}$ with v_i and t crash before a g-failure occurs. With the $\diamond S$ based membership possible for $R = \{r, s, t\}$ to reconnect and decide that the $(k+1)$

Though the restart-view is now known, the group operations cannot does not contain a majority subset of the restart-view. (Permitti form the master subgroup will lead to two concurrent master subgr operating in a seperate partition.) Group re-configuration with R for either p or q to recover/reconnect. Our approach is differen view cannot be determined until the reconnected set (R) contains a the restart view itself (see section 4.1); in other words, determ straight leads to the re-formation of the group (barring the occ failures). As a future work, we intend to compare these two appr more detailed manner.

6. Conclusions

Group failures can occur even in the absence of any physical fail by sudden bursts in message traffic with potentials to lead to v have designed and implemented a configuration management subsystem provide automatic recovery from group failures, once the real/ disappear and components recover. Our system employs a variatio commit protocol for view updates. Consequently, the recovery prov blocking. On the other hand, it is efficient in terms of message and use of stable store, during both normal operations and reco group failure; it costs only one extra message round to update v failure-free periods. This low, failure-free overhead makes our suited to soft real-time systems where it can be incorporated in in [Hurfin98].

Acknowledgements

The work described here has been supported in part by Hew Laboratories, Bristol. Roger Flemming and Paul Harry from Hew Laboratories provided facilities and resources for implementing c system over Somersault. Discussions with Andre Schiper over \diamond S-ba protocol were useful in clarifying many subtle issues.

References

- [Amir92] Y. Amir, Dolev, D., Kramer, S., and Malki, D., "Membership Algorithm for Communication Groups", Proc. of 6th Intl. Workshop on Dist. Algorithms, pp 29-36, 1992.
- [Babaoglu95] O.Babaoglu, R. Davoli, and A Montresor, "Group Membership and View Maintenance in Partitionable Asynchronous Distributed Systems: Specifications", Technical Report, Dept. of Computer Science, University of Bologna, Italy, Nov 1995.
- [Babaoglu97] O.Babaoglu, A. Bartoli, and G Dini, "Enriched View Synchrony: A New Paradigm for Partitionable Asynchronous Distributed Systems", IEEE ToCS, 46(6), pp.642-658.
- [Birman87] K.Birman and T. Joseph, "Exploiting virtual synchrony in distributed systems", 11th ACM Symposium on Operating System Principles, Austin, November 1987, pp. 1-11.
- [Black97] D. Black, P. Ezhilchelvan and S.K. Shrivastava, "Determining the Largest Consistent Process Group after a Total Failure", Tech. Report No. 602, Dept. of Computing Science, Newcastle upon Tyne.
- [Chandra96] T. D. Chandra, V. Hadzilacos, and S. Toueg, "The weakest Failure Detector Consensus", JACM, 43(4), pp. 685 - 722, July 1996.
- [Ezhil95] P. Ezhilchelvan, R. Macedo and S. K. Shrivastava, "Newtop: a fault-tolerant communication protocol", 15th IEEE Intl. Conf. on Distributed Computing Systems 1995, pp. 296-306.
- [Felber98] P Felber, R Guerraoui and A Schiper, "The implementation of CORBA", Theory and Practice of Object Systems, Vol. 4, No. 2, 1998, pp. 93-105.
- [Gray78] J N Gray. Notes on Database Operating Systems. In Operating Systems: Course, Lecture Notes In Computer Science, Vol 60, pp. 393-481. Springer Verlag 1978.
- [Guerra95] R Guerraoui, M Larrea and A Schiper. Non Blocking Atomic Commitment in an Unreliable Failure Detector. Proceedings of IEEE Symposium on Reliable Distributed Systems, Bad Neunhar, Germany, September 1995, pp. 41-51.
- [Hurfin98] M. Hurfin and M. Raynal, "Asynchronous Protocols to Meet Real-Time Requirements: Really Sensible? How to Proceed?", Proc. of 1st Int. Symp. on Object-Oriented Distributed Computing, (ISORC98) pp. 290-297, April 98.
- [Jajodia90] S Jajodia and D Mutchler. Dynamic Voting Algorithms for Maintenance of a Replicated Database. ACM Transactions on Database Systems, Vol 15, No 2, pp. 230-280
- [Keidar95] I Keidar and D Dolev. Increasing the Resilience of Distributed and Replicated Systems. Journal of Computer and System Sciences (JCSS). 1995.
- [Lotem97] E Y Lotem, I Keidar and D Dolev. Dynamic Voting for Consistent Replication of Components. Proceedings of ACM Symposium on Principles of Distributed Computing, pp. 63-71, 1997.
- [Malloth95] C Malloth and A Schiper, "Virtually Synchronous Communication in Distributed Networks", BROADCAST Third Year Report, Vol 3, Chapter 2, July 1995. (Anonymous URL: broadcast.esprit.ec.org in directory projects/broadcast/reports)
- [Melliars-Smith91] P. M. Melliars-Smith, Moser L.E., and Agarwala, V., "Membership Algorithms for Asynchronous Distributed Systems", Proc. of 12th Intl. Conf. on Distributed Computing, pp. 480-488, May 1991.
- [Mishra91] S. Mishra, L. Peterson and R. Schlichting, "A membership Protocol for Consistent Replication", Proc. IFIP Conf. on Dependable Computing For Critical Applications, Tucson, Arizona, pp. 137-145.
- [Moser96] L.E. Moser, P.M. Melliars-Smith et al, "Totem: a Fault-tolerant communication system", CACM, 39 (4), April 1996, pp. 54-63.
- [Murray97] P. Murray, R. Flemming, P. Harry and P. Vickers, "Somersault software", Hewlett-Packard Technical Report, 1997.

[Ricciardi91] A. Ricciardi and K P Birman, "Using Process Groups to Implement Fault-Tolerant Asynchronous Environments", In Proceedings of ACM symposium on PoDC, pp. 480-488.

[Schiper94] A Schiper and A Sandoz, "Primary-Partition Virtually Synchronous Consensus is Harder Than Consensus", Proc. of the 8th International Workshop on Distributed Computing (1994), Sept. 94, LNCS 857, Springer Verlag. (Also in BROADCAST Second Year Report, October 1994).

[Skeen 81] D. Skeen, "Non-Blocking Commit Protocols", ACM SIGMOD, pp.133 - 142.