

Autonomous Multicast Reflectors over Active Networks

Lidia Yamamoto; Guy Leduc

University of Liège, Research Unit in Networking; Institut Montefiore, B28, B-4000 Liège, Belgium
yamamoto@run.montefiore.ulg.ac.be; leduc@montefiore.ulg.ac.be

Abstract

We present a reflector service that seeks to maintain application-level connectivity in the presence of network-level multicast failures. The service is based on the dynamic deployment of autonomous reflectors based on mobile code, on top of an active network infrastructure. It is able to repair multicast tree failures by building a self-organising tree of reflectors, which will be connected to each other via unicast. We focus on the basic decision mechanisms related to code mobility during the tree construction and destruction phases, namely: cloning, migration, merging and termination. To assist with the decisions a market-based mechanism is employed. We show some preliminary simulation results that confirm the viability of the approach and settle directions for further research.

1 Introduction

The demand for multimedia group communication is growing, and multicast is widely recognised as an important service to enable efficient group communication. However, multicast protocols are still not widely available, and the experimental Multicast Backbone (MBone) is slow to take off. The network conditions are unstable, it is difficult to monitor traffic and to detect points of failure, etc. The result is poor quality for the users. Besides that, some users have no MBone access at all, providers are reluctant to allow IP Multicast in their networks, and firewalls can block multicast traffic.

MBone sessions frequently fail due to multicast problems in some sections of the network. It is very difficult to diagnose a failure and eliminate it during the lifetime of the session. As a fallback solution, MBone content transmitters often establish multicast reflectors to serve users that have no multicast access. Experience shows that reflectors are still a key element for a successful MBone session transmission. A reflector is a user-level gateway application that acts as a proxy between a multicast-enabled network and a set of unicast users. It forwards packets from the multicast group to all unicast clients, and from every unicast client to the multicast group and to all other unicast clients. This guarantees that connectivity is maintained within the session, in spite of the fact that some participants have no access to IP multicast, or in the presence of failures in the multicast tree.

Existing reflector software must typically be installed manually, so that the location of a reflector and its configuration must be decided beforehand and informed to all session participants. Therefore, establishing reflectors represent an extra burden for the session organisers, and also an extra burden for the users that need to manually chose when to switch from multicast to the reflector. Besides that, the session users seldom have enough knowl-

edge about the current network conditions in order to be able to choose an optimum location for a reflector. What happens then is that reflectors are typically placed close to the main session source, and all multicast disabled participants must connect to it as clients. This generates an amount of redundant traffic which is proportional to the number of reflector clients, and therefore obviously does not scale to large sessions where potentially large portions of the network might need unicast.

It would be interesting to be able to dynamically install reflectors when there are connectivity failures or administrative restrictions to multicast traffic. The location of reflectors should be automatically determined according to the network conditions observed during the session, and the reflector software should be dynamically installed and started at the chosen locations. At the client side, software agents acting on behalf of one or more group participants should be able to detect when their multicast feed is down, and react by installing a suitable reflector or connecting to an existing one.

We have designed an autonomous reflector based on mobile code that runs on top of an active network execution environment. The candidate locations for such reflector agents are active network (AN) (Tennenhouse et al., 1997) or active server (AS) nodes (Amir, 1998; Fry and Ghosh, 1999). These nodes run an execution environment (EE) capable of downloading and executing the reflector's code, and of discarding it when the session is finished.

Our reflectors are autonomous and decide when to migrate to other nodes, clone in order to cope with increasing demand, merge with other reflectors, or disappear when no longer needed. The decisions are based solely on local knowledge available at the terminals or active nodes where they run. This guarantees that reflector code is deployed only where needed and when needed, and that after the session finishes all the reflectors will be automatically eliminated. The idea is that reflectors

progressively move from the receivers towards the session's main source, until failure points are successfully bypassed. Additionally, reflectors that do not receive sufficient demand die out, and those which are overloaded spawn others to less loaded nodes. Using such a scheme, a loosely connected network of reflectors emerges as a result of failure detection, and disappears by itself when the failure is repaired.

It should be noticed that although our reflectors will tend to cluster around failure points in order to bypass them, they are obviously not able to diagnose nor repair them. Their objective is only to maintain application-level connectivity in the presence of network-level multicast failures. Network management mechanisms to detect and repair such failures are orthogonal and outside the scope of this work.

2 Background

Several commercial and non-commercial reflector systems are available, e.g. Highfield; Live Networks, Inc.; Kirstein and Bennett (2000) (more references on p.27 of Kon et al. (2000)). These systems are typically software packages that must be manually installed at the sites that will provide the reflector service. Therefore the location of reflector sites must be decided beforehand, and cannot be easily changed during the lifetime of the multimedia session. Changes in reflector location or configuration generally lead to temporary service disruption.

A multicast reflector has been mentioned as an example application over the ALAN environment (Fry and Ghosh, 1999). It has been used to transmit Mbone sessions via unicast to sites not connected to the Mbone, and has the potential to move across ALAN nodes. However, dynamic reflector placement algorithms seem not to have been proposed in this context yet.

Kon et al. (2000) propose a distributed framework to manage a network of reflectors, based on dynamic code distribution and (re)configuration. Reflectors are used to support network and terminal heterogeneity. Within this framework, it is possible to manage networks containing a large number of reflectors. Each reflector has a limited degree of autonomy, such as reconfiguring the neighbour nodes to bypass a reflector that failed. For most other operations, however, the reflector elements need to be managed by a privileged user, the reflector administrator, who decides where to install new reflectors or to remove reflectors from the network. To take such decisions, the administrator needs to have a global view of the network topology and characteristics, which is not trivial to obtain from the wide-area Internet or the Mbone. Another difficulty is to cope with receivers dynamically joining and leaving the session, and the corresponding reflector tree reconfiguration in order to maintain a tree which always tracks the optimum. Kon et al. (2000) report that in one experiment they were forced to deny approximately one

million connection requests due to lack of bandwidth on the reflector sites. This could have been avoided if reflectors were able to automatically clone themselves to other sites in order to cope with the additional demand.

In the context of agents and active networking, a number of proposals for self-deploying services have been made. Shehory et al. (1998) propose a framework in which agents deal with overload by cloning, passing tasks to others, merging, or dying. Agents decide when to clone according to the loads of the different resources they use, such as memory, processing and communication resources. The possible decisions that the agent can take are described by a decision tree, and the optimum decision is calculated via dynamic programming. Tschudin (1999a) implemented an election service based on active packets, that deploys itself to every reachable node. In a later work (Tschudin, 1999b) the same author addresses the security issues involved with a necessary self-destructing mechanism for such kind of services.

Najafi (2001) has proposed a cost model for active networks that takes into account the cost of processing a flow in the active nodes as well as the transmission costs. He includes an algorithm that converges to the minimum cost for a flow that may be transformed in several active nodes before reaching its destination. He has also proposed an agent positioning algorithm in which an agent can decide to reposition itself in the network in order to reduce the session cost.

Roadknight and Marshall (2000) address the issue of quality of service differentiation by using a distributed genetic algorithm inspired by the behaviour of bacteria. They show that the amount of servers and their location in the network evolve according to the user demand for a given type of service and a requested trade-off between latency and packet loss. The potential of genetic techniques such as the ones proposed by Roadknight and Marshall (2000) is the increased variability to find new solutions and adapt to new situations not envisaged at the beginning. However, in an environment where nodes and links are heterogeneous, propagating successful rules ("genes") to neighbouring nodes might not necessarily be a good idea, since a rule that is successful in one node might fail completely in another node due to different resource constraints. In the context of self-organising systems, biologically-inspired and market-based techniques seem complementary, and an interesting research challenge would be to combine the best of both worlds to obtain new adaptation mechanisms.

3 Autonomous Reflectors

In this section we describe our autonomous reflector scheme. Since this is work in progress, only part of the mechanisms described here have been implemented at the time of this writing. Sections 4, 5, and 6 provide more details on the implemented mechanisms as well as some

initial simulation results.

Two assumptions are essential for this scheme to work properly. First of all, we assume that some basic default unicast routing service interconnecting all active nodes is available, such that at any moment it is possible for an active application to obtain the next hop to a given destination. This service can be either provided by the EE itself, or installed as active extension code with a well-defined interface exported to the active applications that need it. By default it can simply map directly to IP routing, but more sophisticated techniques such as application-layer routing (Ghosh et al., 2000) could also be available to provide optimised paths according to specific criteria.

The second assumption is that there is only one main source of content in the session (e.g. the lecturer's site), although all session participants can potentially generate data to the session, as it is generally the case with RTP sessions (Schulzrinne et al., 1996), that nowadays are widespread on the Mbone. The reflectors start at the receiver nodes and then progressively move, clone and merge with other reflectors along their respective unicast paths towards the main source. This receiver-initiated approach is similar to filter placement schemes based on RSVP, such as the AMNet prototype described by (Wittmann et al., 1998).

For the time being, the reflectors we propose are only able to repair multicast trees using unicast: multicast tree failures will cause a tree of reflectors to be formed, which will be connected to each other via unicast. Since multicast routing and active network unicast routing are independent, the unicast tree of reflectors might not coincide with the corresponding multicast subtree for a given set of receivers. One can imagine that it would be interesting to use unicast only to bypass "broken" segments of a tree, using multicast everywhere else. We leave this possibility open for future study.

We distinguish two types of reflectors: terminal and intermediate. Terminal reflectors are placed at the receiver side and do not move, while intermediate reflectors are placed on other active nodes in the network, and might move from one node to another according to the network conditions.

A terminal reflector ideally serves one local client running at the same machine, which is the user application that handles the corresponding media (typically Mbone tools such as `vic` or `rat`). A terminal reflector must be installed at each receiver host that wishes to make use of autonomous reflectors. Alternatively it can be installed as close as possible to the receiver host (or set of hosts) to be served. The terminal reflector works as a proxy between the actual multicast group and the user application, so that the direct use of multicast or the use of reflectors is hidden from the application. This allows the use of reflectors based on mobile code without requiring any change to the existing media tools. Terminal reflectors are intermediate reflectors that have their migration rules disabled. They must be fixed because the existing media tools are

not able to detect moving peers.

The tree of reflectors organises itself in a client-server hierarchy. Each intermediate reflector serves a number of downstream reflectors (clients). Here the server will sometimes be called parent reflector, and the clients child reflectors. Most intermediate reflectors will be both servers for a number of clients and clients of an upstream reflector, except a reflector that succeeds bypassing the failure point. This becomes the root of its reflector tree, and therefore plays only the role of server. Actually not only one tree but several ones might arise in response to a failure, in case multiple reflectors cross a failure point at different nodes. This can happen, for example, if the failure "point" is not a single link but a whole network with multicast capabilities disabled for some reason.

When multicast failure is detected, the terminal reflector sends another reflector to the next active hop towards the main source. This operation is called *upstream cloning*. The clone reflector spawned in this way is not an exact copy of the original reflector, but has the same goal of repairing the multicast failure by finding the main session source. It is an intermediate reflector, since it acts as a proxy between the multicast group and the terminal reflector. The intermediate reflector listens to the multicast group for a while, and if no multicast activity is detected it migrates to the next active hop towards the main source (*upstream migration*). This process continues until the reflector reaches a node where data from the main source is received (either via multicast or via another reflector). It then informs the downstream reflector (which in this case is the terminal reflector that originally spawned it) about its current location address. The terminal reflector then starts listening to unicast packets coming from the new upstream address. A simple local selector ensures that no duplicate packets are forwarded downstream, and that packets coming from downstream are forwarded to the selected channel (either multicast or parent reflector).

When two reflectors belonging to the same session meet at the same active node, they merge into a single reflector. A reflector that has more than one client might decide to clone upstream instead of migrating, therefore adding one more hierarchy level to the tree. A reflector might decide to merge with an upstream reflector when it is serving only one client and the upstream connection is unicast, because in this case it is more interesting to establish shortcuts bypassing the node where the reflector currently is. This ensures that most of the time reflectors will only be present at branch points in the reflector tree. Finally, a reflector that runs out of clients automatically terminates itself.

When multicast connectivity is restored, reflectors start receiving data from the main source via the multicast channel, and disconnect from their upstream reflectors. The latter will eventually die out due to lack of clients.

In order to improve reaction time, soft-state variables can be left in the active nodes. Consider a reflector that decides to migrate after having unsuccessfully listened

to the multicast group. Before migrating, it can leave a “message” (in the form of soft-state) in the node, to inform future new-coming reflectors, and avoid them to repeat the unsuccessful experience. Since network conditions change, the reliability of such message must decay with time. Such an indirect communication mechanism among reflectors resembles the stigmergy mechanism used in ant algorithms (Dorigo and Gambardella, 1997).

We are currently working on the decision mechanisms to either clone, merge, migrate or terminate reflectors. These decisions are driven by market-based mechanisms (Clearwater, 1996), and the current state of our research on this issue is summarised in Section 4, with some initial results shown in Section 5.

4 Market-based reflector decisions

Reflectors can use resource control based on market mechanisms (Clearwater, 1996) to make decisions to either clone, merge, migrate or terminate themselves. Such mechanisms can also be used to dynamically decide on the maximum number of clients to accept at a given machine, in order not to cause link or CPU overload. Another usage is to make downstream cloning decisions: for example, in the case of an overloaded reflector, a number of clones can be sent downstream to handle part of the clients.

In this paper we concentrate on the basic decision mechanisms related to code mobility during tree construction and destruction, that is, cloning, migrating, merging and terminating. We are currently working on the additional mechanisms related to tree reshaping and load control.

Each reflector has costs associated with its consumption of node and network resources. The tree of reflectors is organised in a client-server hierarchy, such that server reflectors sell session data to their clients, and buy data from their server reflector. A reflector uses the revenues that come from its clients to pay for resource usage in the active nodes and for the services of the upstream reflector.

4.1 Resource usage costs

Each reflector has associated fixed costs and variable costs for the use of node resources. The fixed costs do not vary with the number of clients that a reflector has, and correspond to the costs of using the mobile code platform. They represent the minimum processing plus storage cost that the mobile agent incurs, even when no clients are connected. Note however that the fixed costs are not constant in general, as they may vary as a function of the load level of the resource in question (CPU, memory).

The variable costs increase with the number of clients, and correspond to the link transmission costs to all clients, plus the processing costs for all packets. These costs may

also vary according to the total load of the corresponding resource (link bandwidth or processing).

From the cost point of view, having many clients is good for a reflector because the fixed costs are shared among all the clients, but if the number of clients becomes too large, the demand for one or more resources might exceed the supply (congestion situation), leading to an increase in processing and link prices, with possible packet losses and consequent degradation in quality for the end user.

Every time a new reflector is added to the tree, there is an increase in costs corresponding to the resources that the new reflector needs. However, this increase might be compensated by a decrease in costs for other reflectors, e.g. because their load is alleviated.

The cost of processing at an active node is also related to the delay penalty imposed to the end user due to the use of a tree of reflectors. The delay penalty is the ratio between the actual delay experienced by an end user and the delay that would be experienced if the receiver could connect directly to the multicast session without the help of reflectors. If the processing power were infinite, the extra delay imposed by a reflector would be null. On the other hand, a very low processing power would incur a high additional delay. The same is valid for a machine with high processing power but which is overloaded, such that the processing time available to a reflector is very low. Therefore if a reflector tries to choose nodes that have low processing costs, it is likely to be moving towards a lower delay penalty for its users.

4.2 Definitions

We begin by providing some definitions of terms that will be used later in this section:

- n_i : a reflector that runs at a given node i
- nc : number of clients of reflector n_i
- $r_{i,j}$: reflector j , the j -th client of reflector n_i , for $j = 1..nc$
- $R_{i,j}$: data sending rate of reflector $r_{i,j}$ to n_i
- nk : number of terminal reflectors in session
- S : data sending rate of the main source

SR : total rate of the session (session bandwidth)

$$SR = S + \sum_{1 \leq k \leq nk} R_k \quad (1)$$

for all terminal reflectors r_k with sending rate R_k .

cf_i : fixed costs at node n_i : costs that don't vary with the number of clients of reflector n_i .

$cv_i(nc)$: variable costs at node n_i : costs that vary with the number of clients (nc) at node n_i .

$cvp_i(nc)$: processing costs at n_i , depend on the amount of data treated per second.

$cvl_i(nc)$: total link costs at n_i : represent the costs associated with the total amount of bandwidth emitted by reflector n_i to each link that leads to clients of n_i , and to the parent reflector if any.

$ct_i(nc)$: total cost at node n_i when nc clients are present: the sum of fixed plus variable costs, as follows:

$$\begin{aligned} ct_i(nc) &= cf_i + cv_i(nc) \\ &= cf_i + cvp_i(nc) + cvl_i(nc) \end{aligned} \quad (2)$$

4.3 Estimating costs

In order to make a decision to either clone or to migrate, a reflector first needs to estimate the costs that would result from choosing either option. A simple decision strategy would then be just to choose the configuration with the lowest cost. However, there are a number of difficulties in obtaining such estimation. Actually this is a typical problem of making decisions in the presence of risks, and decision analysis could be applied here as in Shehory et al. (1998). In this section we present a first simplified approach to the problem. Further research is necessary in order to extend it to a more general case.

One of the main difficulties is that, at the beginning, when the reflector still hasn't reached the main source (directly via multicast or via another server reflector), it is not able to measure the actual resource consumption that will result when it reaches it. When that happens, it goes into full operation mode, but at this moment, it is too late to revise its previous decisions concerning cloning or migrating. Especially, if the reflectors underestimate the aggregate sending rates of all the session members beyond the multicast failure point while building the tree, several points of congestion might appear as soon as the tree becomes fully operational.

A solution to this problem would be to rely on an estimation of the total rate of the session (session bandwidth), that must be available somehow before the session starts. In practice it is possible to obtain such information by looking at the media types in the SDR session announcements. Additionally, if RTCP is used (Schulzrinne et al., 1996), and assuming that only a limited number of session members send significant amounts of data to the group, the session bandwidth grows very little with the total session size.

Using such an upper bound, resources could be reserved at the active nodes along the path in order to guarantee that enough resources are available when the reflectors reach the main source. However, resource reservation might not be available at all nodes, and most of the nodes might not even be active. Besides that, if the session bandwidth is overestimated, too many costs might

incur with little extra benefit for the end user. We adopt a simple solution that relies on an upper bound on the session bandwidth to simplify the cost calculations, but does not reserve resources on the nodes.

Now we try to quantify each cost component in our context. We begin with the processing costs.

4.3.1 Processing costs

Network packets constitute the bulk of the data treated by a reflector. Therefore the processing costs during a given interval increase with the number and size of the packets treated. Every packet received is reflected to everyone else. Thus every packet from the upstream channel (reflector or multicast) is copied to every client, and every packet from a client is copied to the upstream channel plus all the other clients except itself. For a reflector that has already reached the main source (directly via multicast or via another server reflector), the total number of bits per second treated at n_i is:

$$st_i(nc) = sp_i(nc) + sc_i(nc) \quad (3)$$

where:

$sp_i(nc)$ is the data rate sent from the parent to all child reflectors of n_i

$sc_i(nc)$ is the data rate sent from all child reflectors of n_i to all others and to the parent.

$$sp_i(nc) = nc \cdot (SR - \sum_{1 \leq j \leq nc} R_j) \quad (4)$$

$$\begin{aligned} sc_i(nc) &= \sum_{1 \leq j \leq nc} ((nc - 1) \cdot R_j + 1 \cdot R_j) \\ &= nc \cdot \sum_{1 \leq j \leq nc} R_j \end{aligned} \quad (5)$$

Substituting equations 4 and 5 in 3, we have:

$$st_i(nc) = nc \cdot SR \quad (6)$$

Assuming constant prices, and processing costs that increase linearly with the data rate treated, we have:

$$\begin{aligned} cvp_i(nc) &= pp_i \cdot st_i(nc) \\ &= pp_i \cdot nc \cdot SR \end{aligned} \quad (7)$$

where:

pp_i is the (constant) processing price per bit per second at node n_i .

4.3.2 Link costs

The link usage costs include the costs for bandwidth and queueing. Here we consider only the bandwidth costs for simplification. There are only costs associated with the

transmission of packets, not with the reception of packets. Thus the link costs are the sum of the costs to reflect a packet from the parent reflector to all child reflectors, and from each child to every other child plus the parent.

Assuming constant link prices, the total link cost for n_i can be written as:

$$cvl_i(nc) = pl_{i,p} \cdot \sum_{1 \leq j \leq nc} R_j + \sum_{1 \leq j \leq nc} (pl_{i,j} \cdot (SR - R_j)) \quad (8)$$

where:

$pl_{i,j}$ is the price per unit of bandwidth on the link in n_i that leads to the client $r_{i,j}$.

$pl_{i,p}$ is the price per unit of bandwidth on the link in n_i that leads to the parent reflector of n_i (or the candidate parent in case a decision to clone or to migrate is about to be made).

If the link price is the same for all clients and equal to $pl_{i,l}$, or when all clients of n_i share the same link l , we can rewrite the link cost as:

$$cvl_i(nc) = pl_{i,p} \cdot \sum_j R_j + pl_{i,l} \cdot (nc \cdot SR - \sum_j R_j) \quad (9)$$

4.3.3 Cost of the cloning configuration

If we are going to send a clone from the origin node n_i to an upstream destination node n_d , the cost of the resulting clone configuration can be calculated as:

$$\begin{aligned} ctc_{i,d} &= ct_i(nc) + ct_d(1) \\ &= cf_i + cvp_i(nc) + cvl_i(nc) \\ &\quad + cf_d + cvp_d(1) + cvl_d(1) \end{aligned} \quad (10)$$

Here $ct_i(nc)$ represents the total cost of running the agent at the current node when the agent is fully operational, while $ct_d(1)$ is the cost of a new agent running at the upstream node n_d with a single node (n_i) as a client.

4.3.4 Cost of the migration configuration

When migrating to an upstream destination n_d , a reflector n_i carries its client list along with it. Assuming symmetric unicast routing paths, the traffic will continue to go through node i , therefore consuming the same amount of bandwidth resources at the links leading to each client reflector. Since the reflector itself will disappear from node n_i , there are no fixed nor processing costs associated with it anymore at this node. Therefore the cost of resulting configuration after migration can be calculated as:

$$\begin{aligned} ctm_{i,d} &= cvl_i(nc) + ct_d(nc) \\ &= cvl_i(nc) + cf_d + cvp_d(nc) + cvl_d(nc) \end{aligned} \quad (11)$$

4.4 Making a decision

We would like to make a decision to either clone or migrate based on the total costs of resources for each configuration.

A simple decision strategy is to choose the configuration with the lowest cost:

if $ctm_{i,d} > ctc_{i,d}$ then clone else migrate.

In order to simplify the calculations, we rewrite the above rule as:

if $ctm_{i,d} - ctc_{i,d} > 0$ then clone else migrate.

The costs of each configuration are given by equations 11 and 10, therefore we have:

$$\begin{aligned} ctm_{i,d} - ctc_{i,d} &= cvp_d(nc) - cvp_d(1) - cvp_i(nc) \\ &\quad + cvl_d(nc) - cvl_d(1) - cf_i \end{aligned} \quad (12)$$

Note that the fixed costs at n_d , as well as the link costs at n_i have disappeared since they are the same in both configurations. With symmetric unicast paths and no multipath, all the traffic from n_d to n_i will go through a single link. Thus, assuming linear costs for link resources, we can use equation 9 to calculate the terms $cvl_d(nc)$ and $cvl_d(1)$. Assuming linear costs also for processing resources, equation 7 can be used to calculate $cvp_d(nc)$, $cvp_d(1)$, and $cvp_i(nc)$. After these operations we obtain:

$$\begin{aligned} ctm_{i,d} - ctc_{i,d} &= \\ &SR \cdot ((nc - 1) \cdot (pp_d + pl_{d,i}) - nc \cdot pp_i) - cf_i \end{aligned} \quad (13)$$

Which can also be written as:

$$\begin{aligned} ctm_{i,d} - ctc_{i,d} &= SR \cdot nc \cdot (pp_d + pl_{d,i} - pp_i) - cf_i \\ &\quad - SR \cdot (pp_d + pl_{d,i}) \end{aligned} \quad (14)$$

The first line of the right side of the equation 14 represents the increase in costs associated with the migration configuration, while the last line corresponds to the increase associated with the cloning configuration. When nc is high, the migration configuration tends to become more expensive than the cloning configuration. Therefore the cloning configuration will generally be preferred for high nc , unless the fixed or processing costs at n_i are prohibitive.

With this result we obtain an easy way to make a decision, by using equation 13 to choose the cheapest configuration. As discussed earlier, an estimation on the upper bound of SR is considered available before the session starts. The number of clients, nc , is known at n_i , as well as the local cost cf_i and price pp_i . Consequently, before making a decision, the agent needs to obtain the following information from its neighbour n_d :

pp_d : processing price per unit

$pl_{d,i}$: link price per unit for the outgoing interface from n_d to n_i .

The information above is collected by a Prospecting capsule that is sent to the destination node before the cloning or migration action actually takes place.

This is the strategy adopted to obtain the simulation results shown in this paper. Although it seems a bit too simplistic, this strategy already takes into account an important criterion which is the delay penalty for the receiver which is imposed by the use of the reflectors instead of native IP multicast. As discussed earlier, this delay penalty is implied within the processing costs.

In classical multicast algorithms such decision dilemma usually doesn't apply, because only link resources are typically taken into account. In this case we can make $pp_d = pp_i = cf_i = 0$, and our calculations reduce to:

$$ctm_{i,d} - ctc_{i,d} = SR \cdot pl_{d,i} \cdot (nc - 1) \quad (15)$$

In the above we have:

For $nc > 1$: $ctm_{i,d} - ctc_{i,d} > 0$ and we choose to clone.

For $nc \leq 1$: $ctm_{i,d} - ctc_{i,d} \leq 0$ and we choose to migrate.

These observations confirm that when bandwidth is the only scarce resource, cloning is the default choice except in the trivial case ($nc \leq 1$), since migration always implies duplicating packets on the link from n_d to n_i when $nc > 1$, and therefore causes the total costs to increase.

4.5 Merging

When two reflectors belonging to the same session meet at the same active node, they merge into a single reflector. This operation involves the union of both client lists and any other necessary other information. Since this might result in resource overload, a preliminary negotiation between both agents is desirable to achieve favourable configurations. For instance, when sending the Prospecting capsule to an upstream neighbour to find out about costs, the capsule could also be programmed to look for the presence of another reflector for the same session, and check its current resource consumption. An outcome of the negotiation could be that server delegates some of its own clients to the prospecting reflector, in order to balance the load and reduce costs.

This is related to the tree reshaping problem (splitting operation) that we have not addressed yet. Currently the agent takes cloning or migration decisions independently on the fact that another reflector might already be present in the upstream node, therefore the merging operation is always carried out.

4.6 Terminating

Since reflectors must "pay" for resource usage in the active nodes, and their clients are their sole source of in-

come, they will be automatically eliminated by the active platform when there are no further clients. However, there is a risk that sudden changes in load make prices increase in unpredictable ways, causing fully operational reflectors to die out prematurely.

In our current implementation this problem is still not solved, and in our view it can only be solved with the help of load control operating at shorter time scales than the ones in which the reflectors operate. This requires adaptive (elastic) flows or transcoding, and here we are assuming that the reflectors merely repair connectivity failures, and don't interfere with the session data contents.

5 Simulation results

We have performed some simulation experiments using ns-2, in order to visualise the tree construction and destruction mechanisms. The topology for the simulations is illustrated in Figure 1. All links have a fixed capacity of 10Mbit/s and a propagation delay of 10ms. The root of the multicast tree is node S where the main session source is located. The leaves of the tree contain terminal reflectors that join the session at random times from $t=0s$ to $t=10s$. All nodes are active and have the same prices for resources: $pp = 1$, $pl = 1$, and $cf = 1 \cdot 8 \cdot C$ where C is the size of the mobile code in bytes, and is currently set to 50000, which is the current approximate size of the bytecode in our Java prototype.

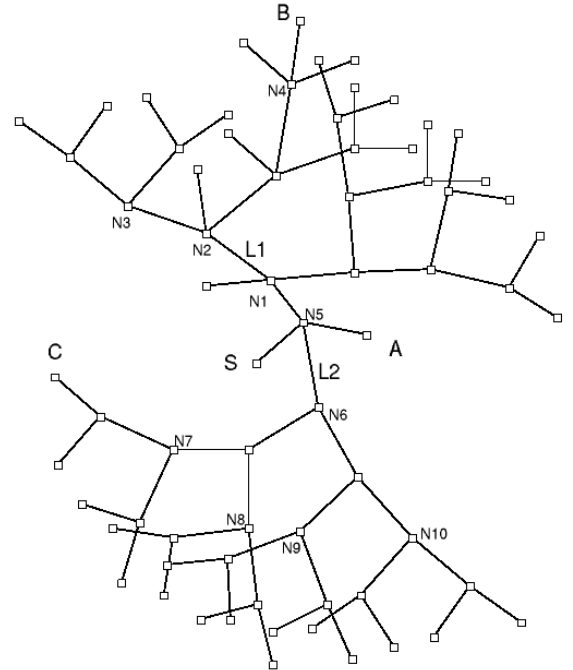


Figure 1: Topology used in the simulations.

The multicast communication via links L1 and L2 is interrupted at $t=20s$. As a result two reflector trees appear.

Both trees starts at around $t=24s$, but the tree on the upper side of the topology is ready at $t=29s$, while the second one is only ready at $t=37s$. After this construction phase all terminal reflectors affected by failures are served by an intermediate reflector. The resulting trees are shown in Figure 2.

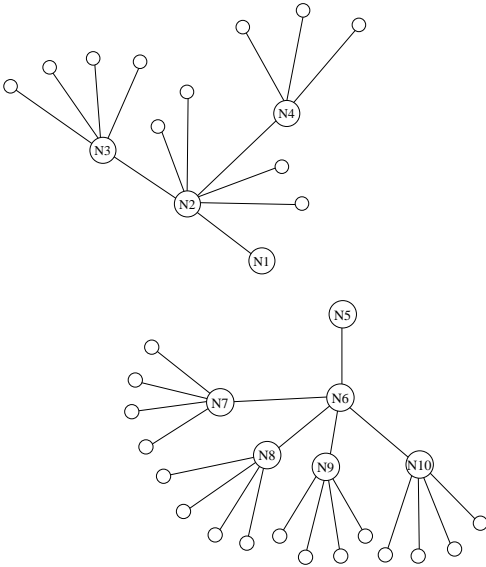


Figure 2: Two sample reflector trees. Top: tree that results from the failure of L1, rooted at N1. Bottom: tree that results from the failure of L2, rooted at N5.

At $t=70s$ the multicast communication via L1 and L2 is restored. Most reflectors detect this a couple of seconds later, and disconnect from their parent reflectors, which die out between $t=77s$ and $t=80s$.

Figure 3 shows the aggregate session data rate received by three sample session participants: A, B, and C, whose location in the tree can be observed in Figure 1. The main source rate is 500kbps while all the other session members send around 10kbps each.

Participant A happened to join the group at around $t=10s$, while B joined right at the beginning $t=0s$. During the failure period, although the multicast feed to node A is up and running, it receives less aggregate traffic until the reflector tree is fully operational, since during this period it doesn't receive the multicast packets coming from the participants that have stayed on the other side of the failure point. Receivers B and C suffer from the failures until about $t=30s$ and $t=38s$, respectively. After that, their respective level of reception becomes about the same as the one of A, as if they were also unaffected by the failure. When the multicast feed is restored, sudden peaks of traffic arrive at B and C, due to duplicate packets sent once via multicast and again via the reflector. These packets are eliminated by the terminal reflectors before being sent to the applications. We can verify this by looking at the sequence numbers received by the decoder connected to C, on Figure 4.

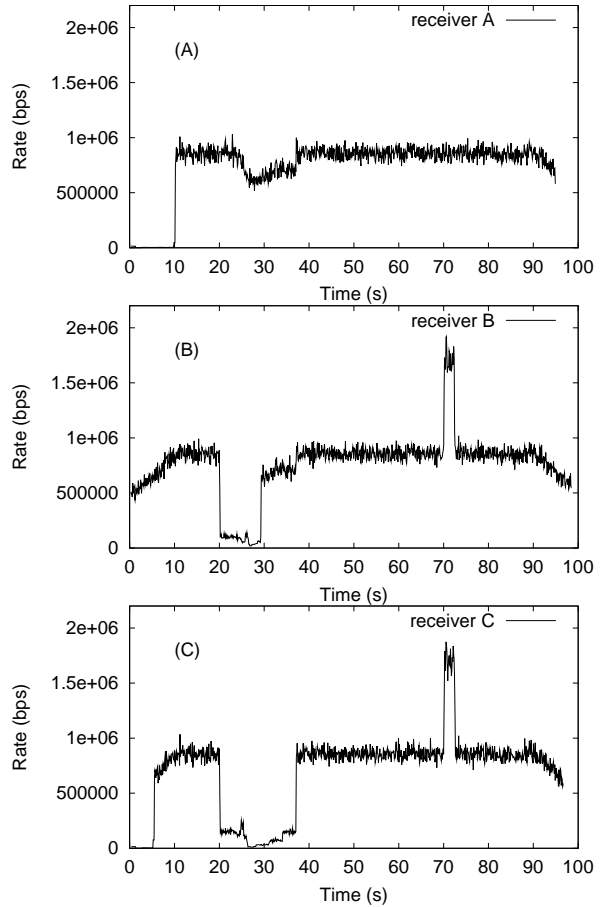


Figure 3: Data rate of three sample participants. A: unaffected by multicast failure. B and C: affected by failures of L1 and L2 respectively.

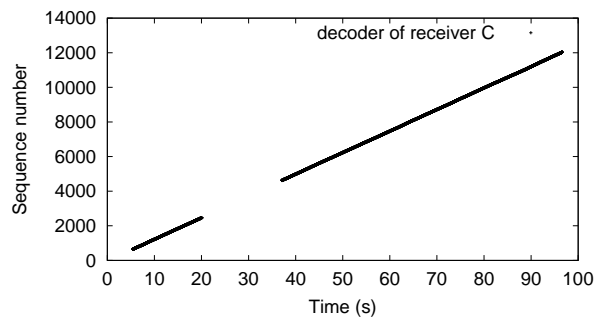


Figure 4: Sequence numbers received by the decoder of participant C.

In order to visualise the dynamics of the mobile code operations of migrating, moving, merging and terminating, we describe them in Figure 5 for the upper reflector tree on the topology. For compactness, we number events in time with integer numbers starting from 1, followed by the code of the operation performed. From $t=24s$ to $t=26s$ all terminal reflectors spawn upstream clones. Since the

order in which each terminal reflector sends a clone won't have any influence on the subsequent operations, we assign event number 1 to all. This is indicated as "1C" in the figure. The next event is event 2, and it's a cloning operation from node N4 towards its upstream neighbour. This is indicated as "2C". By following the sequence of events in this way, it is possible to track the main actions that lead to the tree configuration shown in Figure 2, and to its subsequent destruction ("T" operation). Although the merging operations are not indicated, they can be deduced as well, since they occur whenever a reflector arrives at a node where another one is already present.

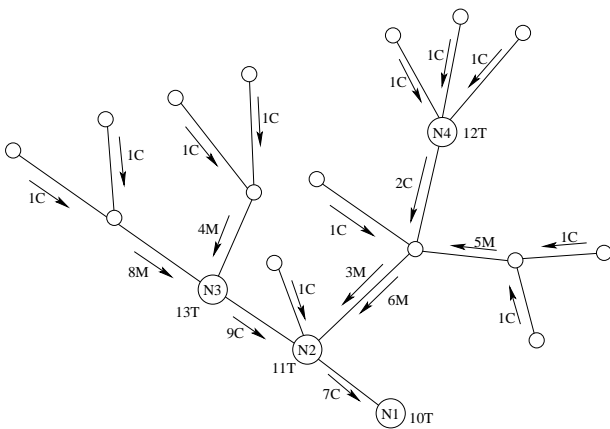


Figure 5: Dynamics of mobile code operations. The event numbers are shown beside the arrows or the node names, followed by the code of the operation: C (clone), M (migrate), T (terminate).

The results above are intended to illustrate the basic behaviour of our autonomous reflectors in ideal conditions. They are not intended to show a realistic picture of a real network. The network here is unloaded, all nodes are active, the delays are short and the paths for unicast and multicast traffic coincide. In our simulations we have noticed little impact of increased network latencies or moderate load on the results, even when the propagation delay on each link is increased to the order hundreds of milliseconds. However, we have often observed much larger latencies for joining live Mbone sessions, as well as variable loss patterns. Thus we can expect higher reaction times for our reflectors in such a situation.

6 Implementation

We are currently implementing the mobile reflectors in Java using an architecture that allows the code to be easily ported to any EE that supports active extensions with minor modifications. The architecture is organised in three planes: data plane, monitoring plane, and control plane. This structure roughly follows the one suggested by Blair et al. (1999), although no computational reflection capabilities are included yet.

The data plane is responsible for the blind forwarding of multicast and unicast data. Its core is inspired by the Mug reflector (Highfield): in Mug, a node that sends a UDP packet to the reflector is added to its client list; a client that stays idle (i.e. sends no more packets) for some time is removed from the client list. A selector is attached to the Mug-like core in order to switch between the multicast and the unicast upstream channels. The selector is controlled by the control plane. The monitoring plane keeps track of the current resource usage and performance parameters of the data plane. The control plane uses the data available in the monitoring plane to make decisions. Its core is a state machine with transitions triggered by events generated at the monitoring plane.

This architecture allows new strategies to be easily added to the control plane without affecting the other planes. It also maps naturally to the Bond platform (Bölöni and Marinescu, 1999), which opens up future possibilities for dynamic updates to the state machine through "agent surgery" (Bölöni and Marinescu, 1999). The communication mechanism among neighbouring reflectors takes the form of capsules such as in ANTS (Wetherall et al., 1998). Alternatively, the communication could be made via an existing agent message passing mechanism (e.g. Bond, see Bölöni and Marinescu (1999)). Both offer extra flexibility for enhancements and preclude the need to specify application-specific message formats and develop the corresponding parsers.

7 Conclusions and Future Work

We have described a decentralised scheme based on mobile code, to build a loosely connected network of autonomous reflectors that seeks to maintain session connectivity in the presence of multicast failures. The self-organising nature of the scheme ensures its robustness, scalability and autonomy properties, which make it suitable for sessions of any size, while minimising the necessary amount of human intervention.

For the moment each reflector treats only one media stream (e.g. either audio, or video, or whiteboard). In order to deal with several media, we plan to group multiple physical reflectors (each treating one media type) into a single logical reflector for cloning and migration purposes. In a near future, experiments over the Mbone can be envisaged in the framework of the European COST Action 264, and with the help of existing active network overlays such as the ABONE.

The work presented in this paper is part of a larger effort on the use of active networks for adaptive applications. We plan to integrate it with previous work on congestion control (Yamamoto and Leduc, 2000b,a) such that reflectors also perform application-oriented filtering and/or transcoding of data in the presence of congestion, in a network which is likely to be only sparsely populated by active nodes.

Acknowledgements

This work has been carried out within the TINTIN project funded by the Walloon region in the framework of the programme “*Du numérique au multimédia*”. Part of this work was performed while the main author was a visiting researcher at Lancaster University. We would like to thank David Hutchison, Steven Simpson, Mark Banfield, Laurent Mathy, Stefan Schmid, and Katia Saikoski for their helpful support.

References

- E. Amir. *An Agent-based Approach to Real-time Multimedia Transmission over Heterogeneous Environments*. Ph.D. dissertation, University of California at Berkeley, 1998.
- G. S. Blair, A. Andersen, L. Blair, and G. Coulson. The Role of Reflection in Supporting Dynamic QoS Management Functions. In *IEEE/IFIP International Workshop on Quality of Service (IWQoS)*, London, UK, June 1999.
- L. Bölöni and D. C. Marinescu. A Multi-Plane State Machine Agent Model. Technical Report CSD-TR 99-027, Purdue University, September 1999. Also a poster at the Fourth International Conference on AUTONOMOUS AGENTS (Agents 2000) Barcelona, Spain, June 2000.
- S. Clearwater, editor. *Market-based Control: A Paradigm for Distributed Resource Allocation*. World Scientific Publishing, 1996.
- M. Dorigo and L. M. Gambardella. Ant Colony System: A Cooperative Learning Approach to the Traveling Salesman Problem. *IEEE Transactions on Evolutionary Computation*, 1(1):53–66, 1997.
- M. Fry and A. Ghosh. Application level active networking. *Computer Networks*, 31(7):655–667, 1999.
- A. Ghosh, M. Fry, and J. Crowcroft. An Architecture for Application Layer Routing. In *Second International Working Conference on Active Networks (IWAN 2000)*, Springer LNCS 1942, pages 71–86, Tokyo, Japan, October 2000.
- J. Highfield. Mug multicast packet reflector, URL <http://www.stile.lboro.ac.uk/cojch/mug/mug.html>.
- P. T. Kirstein and R. Bennett. RE 4007 MECCANO Project Final Report, June 2000. URL <http://www-mice.cs.ucl.ac.uk/multimedia/projects/meccano/deliverables/>.
- F. Kon, R. Campbell, and K. Nahrsted. Using Dynamic Configuration to Manage A Scalable Multimedia Distribution System. *Computer Communication Journal* (Special Issue on QoS-Sensitive Distributed Systems and Applications), 2000. Elsevier Science Publisher, Fall 2000.
- K. Najafi. *Modelling, Routing and Architecture in Active Networks*. Ph.D. dissertation, University of Toronto, Canada, 2001.
- C. Roadknight and I. W. Marshall. Differentiated Quality of Service in Application Layer Active Networks. In *Second International Working Conference on Active Networks (IWAN 2000)*, Springer LNCS 1942, pages 358–370, Tokyo, Japan, October 2000.
- H. Schulzrinne, S. L. Casner, R. Frederick, and V. Jacobson. RTP: A Transport Protocol for Real-Time Applications, January 1996. Internet RFC 1889 (update in progress).
- O. Shehory, K. Sycara, P. Chalasani, and S. Jha. Agent Cloning: An Approach to Agent Mobility and Resource Allocation. *IEEE Communications Magazine*, pages 58–67, July 1998.
- D. L. Tennenhouse et al. A Survey of Active Network Research. *IEEE Communications Magazine*, 35(1):80–86, January 1997.
- C. F. Tschudin. A Self-Deploying Election Service for Active Networks. In *Proc. 3rd International Conference on Coordination Models and Languages (COORDINATION’99)*, Springer LNCS 1594, pages 183–195, Amsterdam, The Netherlands, April 1999a.
- C. F. Tschudin. Apoptosis - The Programmed Death of Distributed Services. In J. Vitek and C. Jensen, editors, *Secure Internet Programming - Security Issues for Mobile and Distributed Objects*, Springer LNCS 1603, pages 253–260, July 1999b.
- D. J. Wetherall, J. V. Guttag, and D. L. Tennenhouse. ANTS: A Toolkit for Building and Dynamically Deploying Network Protocols. In *Proceedings of IEEE OPENARCH’98*, San Francisco, USA, April 1998.
- R. Wittmann, K. Krasnodembski, and M. Zitterbart. Heterogeneous Multicasting based on RSVP and QoS Filters. In *SYBEN’98*, Zürich, Switzerland, May 1998.
- L. Yamamoto and G. Leduc. An Active Layered Multicast Adaptation Protocol. In *Second International Working Conference on Active Networks (IWAN 2000)*, Springer LNCS 1942, pages 180–194, Tokyo, Japan, October 2000a.
- L. Yamamoto and G. Leduc. An Agent-Inspired Active Network Resource Trading Model Applied to Congestion Control. In *Proceedings of the MATA 2000 Workshop*, Springer LNCS 1931, pages 151–169, Paris, France, September 2000b.