

Network control as a distributed object application

This article has been downloaded from IOPscience. Please scroll down to see the full text article.

1998 Distrib. Syst. Engng. 5 19

(<http://iopscience.iop.org/0967-1846/5/1/003>)

View [the table of contents for this issue](#), or go to the [journal homepage](#) for more

Download details:

IP Address: 38.107.179.214

The article was downloaded on 20/02/2012 at 07:50

Please note that [terms and conditions apply](#).

Network control as a distributed object application

Huw Oliver[†], Søren Brandt[‡], Andrew Thomas[§] and Nathalie Charton^{||}

Hewlett-Packard Laboratories, Filton Road, Stoke Gifford, Bristol BS12 6QZ, UK
and
Alcatel Alsthom Recherche, Route de Nozay, 91460 Marcoussis, Paris, France

Received 25 February 1997

Abstract. We describe our experiences and performance figures from an open-switching approach to broadband (ATM) connection management. The control system described is the lowest layer of a general connection management architecture such as TINA or Xbind.

1. Introduction

The following work was carried out as part of the ACTS Project ReTINA (AC048).

The aim of ReTINA is to develop and demonstrate an industrial-strength open distributed processing environment (DPE). This DPE will support distributed real-time multi-media applications over emerging broadband networks.

In terms of technology, the project will deliver a CORBA-compliant distributed object-based DPE with real-time and quality-of-service management capabilities.

We describe our experiences and performance figures from broadband (ATM) connection management using a distributed (CORBA-based) control system. The control system described is the lowest layer of a general connection management architecture such as TINA [5] or Xbind [3, 4].

1.1. ATM/ISDN Forum style signalling

The TINA or Xbind architectures are an attempt to improve upon the current approaches to connectivity services in broadband ISDN/ATM [1, 2, 13].

Existing work on connectivity services in broadband ISDN/ATM standards bodies has concentrated on traditional signalling techniques. These come from the telephony world. For example, the ATM standard for signalling (SVC) is now Q.2931 which is based on the narrowband ISDN (digital telephony) Q.931 standard.

While there are suggested enhancements to Q.2931 still under study, there are several outstanding problems or, more accurately, limitations to this signalling paradigm that will greatly compromise the potential success of the new public networks.

The limitations with this model can be summarized as follows.

(1) They are geared to running a single connection within a single call. We need much more flexibility. We need complex topologies (and we need some notion of a session to encompass these multiparty connections).

(2) They assume that calls last for minutes and that set-up latencies are correspondingly non-critical. Now we want near-instant set-up times. An example is a user switching between video channels.

(3) The implementations of UNI signalling have too large a footprint. All code needed to run the protocol has to be built into the terminal. The switch also needs to run the entire protocol leading to a huge amount of code running on the switch. The code is also all mixed up with session level logic about what the end parties can do.

(4) More than the code, the specification of the ATM-Forum UNI 3.0 requires 143 textbook pages to define the connection management signalling. This is because it must define the delivery semantics and encoding at the bit level. An equivalent CORBA-based interface for slightly more functionality requires three pages and is mainly self-documenting. The expression 'CORBA-based' automatically defines the delivery semantics and encoding.

We have been motivated, therefore, to look at connection management approaches that overcome the above problems and that meet some new requirements. We can summarize the requirements as follows.

- Complex topologies: we need the ability to establish connections of more complex topologies beyond simple point-to-point connections. We will want multipoint-to-multipoint connections. We will want to provide for connectivity of multiple types of media flow. We will also need higher level abstractions for constructs to describe these operations. Note that this is nowadays possible as the customer premises equipment becomes at least as intelligent as the switching equipment.

[†] E-mail address: heo@hplb.hpl.hp.com

[‡] E-mail address: sbb@hplb.hpl.hp.com

[§] E-mail address: dat@hplb.hpl.hp.com

^{||} E-mail address: Nathalie.Perdigues@aar.alcatel-alsthom.fr

- We want to separate the connectivity software from the switching hardware. The intention is to allow the access to connectivity services by a third party. The hope is to achieve service creation, deployment and management as easily as it currently is in the personal computer and Internet worlds.

This separation should be achieved via an open connection management interface.

Connection management needs access to interfaces that can provide information about the state of the network elements and their current utilization. Information about the network topology and current link utilization facilitates the set-up process. To this end, each connection controller should maintain a view of the network to enable it to calculate a complete provisional route that then must be verified.

This need for access to resource information also applies to call admission algorithms. They need knowledge of the state of the links in the net (i.e. the traffic load) in order to perform.

- Connection management should provide for connectivity of multiple types of media ‘flow’. Of these, certain traffic flows will need special handling. Connection management needs to be able to respond to sender requests for special handling of real-time services. This is important to support applications that require some degree of consistent throughput, delay and/or jitter.

Following from the network’s need to deliver to a particular flow a quantitatively specified quality of service, it is usually necessary for the network to set aside certain resources such as a number of buffers for that flow. Thus, we need to be able to create and maintain resource reservations on each network element along the target path.

- Fast Set-up: for many services the speed of call set-up and modification is critical. Examples are very demanding scenarios (including those, for example, when integrated televisions and computers appear). For these, the call set-up processing will be at least as demanding as the current call processing demands on telecommunications switch control points. Other services, for example video on demand, will have undemanding set-up latencies (10 s for a set-up will be acceptable).

- Management: connection management should integrate seamlessly with multimedia general network management covering configuration, performance, fault management, accounting and security.

There are many other requirements not detailed here such as scalability, interoperability, simplicity, monitoring, flow synchronization, ownership and joining mechanisms.

1.2. Connection management architectures

Connection management architectures such as TINA [5] and Xbind [4, 6] separate connection management into two layers:

- (1) Link-by-link set-up (bottom layer).
- (2) Routing (top layer).

Link-by-link set-up is concerned with the creation of connection through single network elements such as ATM switches. Routing is concerned with the determination

of a particular path through the network to affect a connection with certain characteristics (i.e. quality-of-service guarantees). Sections 1.2.1 and 1.2.2 provide short introductions to routing in TINA and Xbind.

1.2.1. Routing in TINA. In TINA, routing is achieved by the decomposition of trails into smaller requests to connection performers down a CP hierarchy as illustrated in figure 2. The pattern is repeated until the switch agents in the fabric are directed to set up their part of the trail. The routing, then, is done as part of the recursive edge creation done by successive network management layer connection performers (NML-CPs) hierarchy calls. No part of the network has a complete knowledge of the internal configuration of its subnetworks.

The layer network† coordinator (LNC) receives requests for connectivity in a layer network. It has federation capabilities with other LNCs of other domains in the network. Because these capabilities are immature and are the subject of further study in TINA, we have only looked at single-domain layer networks.

Underneath the LNC there is a hierarchy of NML-CPs. Each NML-CP attempts to satisfy requests to create, modify and release connections within its own subnetwork. The NML-CP may pass these requests on to other NML-CPs in charge of various partitions of the subnetwork and this may continue recursively.

Ultimately the request will be passed to an object that has access to an agent on a real switch. These element management layer connection performers (EML-CPs) attempt to satisfy requests to create, modify and release real connections between termination points of the subnetwork.

The critical interface, for this discussion, is the subnetwork connection (SNC) interface to the NML-CP. This is based on an abstract model of the network resources called the NRIM [12]. Essentially it is concerned with the creation, deletion and modification of edges. It is the specialization (in object-oriented terms) of this interface that is supported by the EML-CP.

1.2.2. Routing in Xbind. In Xbind, the aim is to make the network programmable by exposing the network state. The various binding algorithms—Xbind does not specify these—operate upon a set of interfaces known as the binding information base (BIB). Xbind seeks to populate the BIB so that interesting algorithms can be built on top of the information that the BIB holds. Connection management is then an algorithm which runs on top of the BIB.

The BIB offers a wide range of object types: some deal with terminal capabilities, which might be used in constructing the terminal connection session manager (TCSM) of TINA; others deal with switching hardware and there are objects representing resources that are used to determine whether QoS constraints can be met.

With respect to network control, then, the Xbind BIB is currently populated with objects for controlling ATM

† A layer network is one that consists of a single technology.

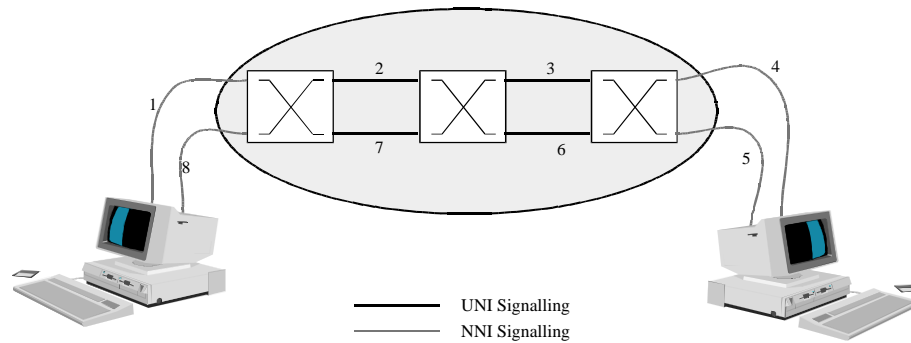


Figure 1. Traditional signalling.

switches. These objects allow the construction of QoS-constrained physical links, algorithms external to the switch construct end-to-end connectivity out of the primitive services offered by the switch objects. This flexibility for example allows third-party connection set-up, leaf initiated joins and so on. Routing algorithms in Xbind are again applications running over the BIB.

In summary, with Xbind the programmer is free to implement any algorithm over the BIB. The connection management framework in TINA implies hierarchical routing. However, this approach is not prescriptive for implementing connectivity in TINA. The programmer can implement any algorithm (hierarchical or not) assuming he/she respects the upper level interface of the connection management which is specified as TINA connectivity provider reference point (ConS-RP).

1.3. Open switching

The key element to support architectures such as TINA and Xbind is that each switch advertises its resources, making them accessible to algorithms running outside the switch. Third party providers can then develop hardware-independent connectivity services where their software utilizes low-level open interfaces for switch control and configuration.

Figure 2 shows two terminals making and receiving requests to/from an off-switch connection controller. The connection controller coordinates the set-up requests to the switches involved in the path.

We have experimented with several different ways of implementing the interface between the switches and the off-switch controller. Broadly, we have used a Fore ASX-100 switch (running SunOS) and two different controller platforms NT and HP-UX. We have been impressed by the use of the Java virtual machine (JVM) on the different controllers to achieve platform portability and also the use of the JVM on the switch to build the server software. The use of ‘Java everywhere’ has allowed us to drop useful pieces of software onto the VM such as a CORBA ORB. Without Java it would have been a substantial piece of platform-dependent work to embed an ORB on the switch.

We have experimented with different platforms for the controller (HP-UX and NT), using the same IIOF interface to the switch. Figure 3 shows two examples of the software

pieces. The central block represents the switch while the right and left blocks show two alternative controllers. The IIOF interface corresponds to the protocol between interfaces 4 and 5 in figure 2.

2. Background specific to our experiment

One of the important aims of ReTINA is to provide a set of DPE services in order to support telecommunication applications [12]. One of these DPE services is a TINA connectivity service. In the context of ReTINA (release 1), this service is supported by an implementation of TINA connection management provided by the Alcatel† Alcin CM software.

The Alcin software has been deployed by Telenor as a central component of the ReTINA broadband private virtual network (BVPN) demonstrator [8]. The BVPN demonstrator investigates how to open up connection management in public networks, in order to allow third-party providers and customers to create ‘leased lines’ on demand, without the need to contact the telecommunications operator.

The experiments described in the following sections were prompted by the belief that use of Simple Network Management Protocol (SNMP) constitutes a performance bottleneck in the first BVPN version, which uses SNMP to do link-by-link set-up. SNMP was suggested as a possible explanation to the inferior performance experienced by Telenor during BVPN experiments. In this paper, we therefore compare SNMP performance with the performance of an alternative switch interface developed at Hewlett-Packard laboratories. The overall conclusion of these measurements is that SNMP performance can indeed be improved by an order of magnitude by means of an alternative switching interface. However, our measurements also suggest that SNMP performance is unlikely to be the primary BVPN bottleneck.

The physical set-up used during performance measurements is illustrated in figure 4.

Connections through the ASX-100 ATM switch were created by means of either SNMP or our open switching interface. The latter will be described in more detail in section 3.1. As illustrated in figure 4, all workstations,

† Alcatel and Telenor are partners in the ReTINA project.

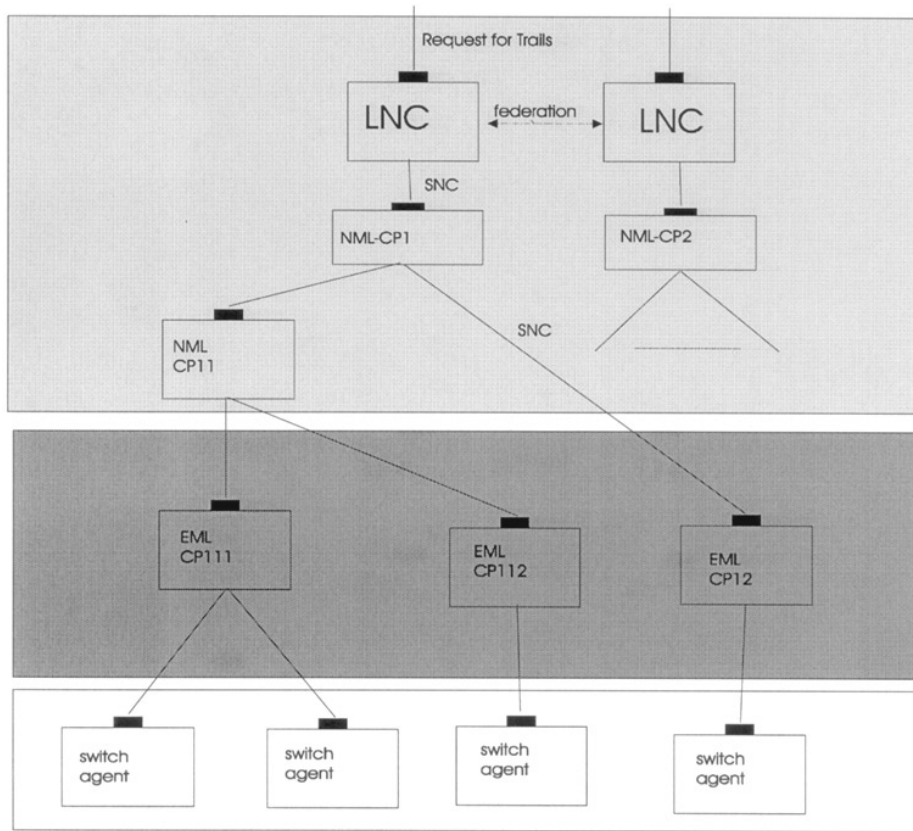


Figure 2. TINA network control structure.

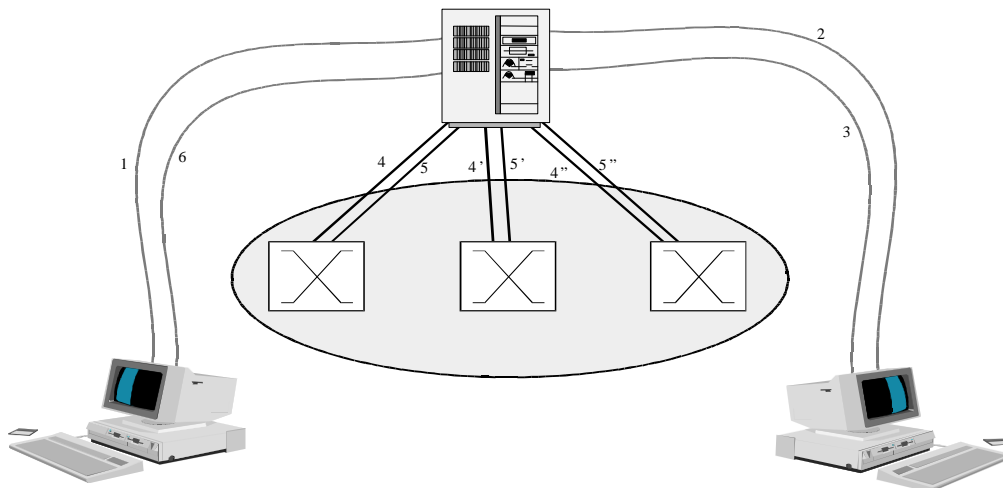


Figure 3. Message 1 is a call initiation request. 2 is a set-up request to the end terminal. 3 is a call acceptance from the end terminal. 4 are resource allocation messages issued in parallel to network switches. 5 are returned acknowledgments. 6 is call accept.

including the switch itself, were connected by means of an Ethernet. The SNMP protocol and the open switching protocols both ran over the Ethernet, and did not use the ATM connections. Hence, the actual ATM connections were not used during the performance measurements described in the following section. However, to test that connections were actually created and deleted, we used two

KNET AVA-200 boxes. The AVA boxes sample analogue audio and video inputs, and feed ATM AAL5 streams onto the ATM network. The resulting streams were routed through the switch to software decoded video and audio sinks running on a HP-UX workstation (torino) equipped with an ATM card.

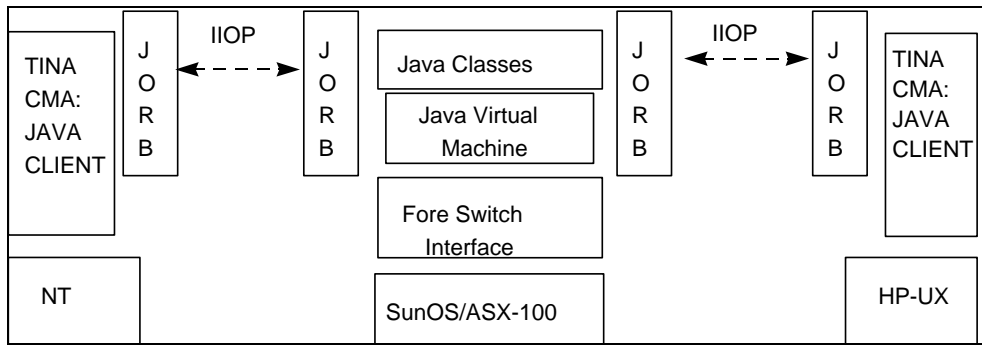


Figure 4. Java allows easy porting of the client software.

3. Performance figures

This section describes a series of experiments and performance measurements carried out to assess the achievable performance of connection creation and deletion on an ATM switch. One of the comparisons we make in this section is between the use of SNMP as a switch control interface and a more direct control interface developed at Hewlett-Packard laboratories. The overall conclusion of these measurements is that SNMP performance can indeed be improved by an order of magnitude by means of the alternative switching interface.

However, the performance measurements described in this section also provide input to other discussions, as listed below.

- (1) Performance of SNMP versus alternative link-by-link set-up mechanisms.
- (2) Performance of different Java implementations.
- (3) Performance of C versus Java in I/O bound applications.
- (4) Performance of CORBA versus raw sockets.
- (5) Deployment of TINA CM computational objects.

Item (5) needs a brief explanation. We saw in section 1.2.1 that the TINA CM framework consists of a number of computational objects, each described by CORBA IDL interfaces. At the bottom layer, so-called element management layer connection performer (EML-CP) objects represent the physical network components, such as individual switches. However, the TINA framework does not specify where to actually deploy these objects, leaving this decision to actual implementations. At one extreme, all objects could be instantiated in a central computer, communicating with switches by means of whatever protocol is available, including SNMP. At the other extreme, an object representing a physical entity might run directly on the entity represented. For example, an EML-CP could run directly on the switch, providing a CORBA interface to the outside world. Performance measurements in this paper suggest that the latter solution, of running the EML-CP on the switch, will result in inferior performance, mostly due to the performance of the switch hardware itself. The use of another switch with an up-to-date control station (e.g. a Pentium) would probably lead to another conclusion.

3.1. An alternative switch interface

As an alternative to SNMP, a 'connection server' process was hosted on a Fore Systems ASX-100 ATM switch running SunOS 4.1.1. The server process allows clients on other hosts to create and delete connections by means of RPC calls. The interface of the connection server is shown below:

```
interface ConnectionServer {
    add_connection(in_port,in_vp,in_vc,
                  out_port,out_vp,out_vc);
    delete_connection(in_port,in_vp,in_vc,
                     out_port,out_vp,out_vc);}
```

The connection server is implemented on top of the Vince [11] software, which provides an interface to the switch table by means of memory mapping. Hence, internal to the connection server, the creation or deletion of a connection amounts to a few memory accesses to update the memory mapped switch table. Since the amount of work done in the server to create a connection is therefore negligible, the performance of the RPC mechanism becomes crucial to the overall performance. To quantify this claim, we measured the processing time used internally in the server to create or delete a connection. When calls to add or delete a connection were carried out from inside the server process itself, on the order of 10000 calls/s were achieved with the server implemented in Java. That is, the time spent on work proper inside the server is on the order of 0.1 ms/connection created or deleted, and can be ignored as compared with communication overheads.

The negligible cost of work internal to the server meant that we could run the server in dummy mode during performance measurements. That is, the server did not actually create or delete connections, but simply returned success to the calling client. This allowed us to run the server on machines other than the switch itself, in order to compare switch performance with the performance of other workstations.

3.2. Switch interface test configurations

A number of different Java implementations of the connection server were created, allowing comparison of several: (1) Java implementations, (2) communication mechanisms, and (3) hardware platforms. It was possible

Table 1.

Machine	Model
asx100	FORE Systems ASX-100 ATM switch
torino	HP PA 9000/755
brandt-s-1	HP Vectra Xu 6/200 (Pentium Pro)
vecheo	HP Vectra Xu 6/200 (Pentium Pro)

to run Java programs on the switch by means of a port of the Java virtual machine to the switch, as was illustrated in figure 3. Finally, an implementation of the connection server in C allows us to compare the performance of Java with C.

3.2.1. Hardware. The hardware platforms involved in performance measurements are listed in table 1.

All machines, including the switch, were interconnected by an Ethernet LAN as illustrated in figure 4.

3.2.2. Languages. The connection server and its client were implemented in Java. To assess the performance differences between different Java implementations, several Java virtual machines were used, as listed in table 2.

Furthermore, to compare Java with C, a version of the connection server was implemented in C.

3.2.3. Communication mechanisms. Physically, network connections were provided by an Ethernet LAN. Logically, communication was achieved by means of two distinct mechanisms:

- (1) a CORBA/IIOP interface;
- (2) a raw interface based on TCP/IP sockets.

The connection server can be configured at start-up to provide either or both of the two interfaces, allowing clients to choose the means of communication. During measurements, the server was configured to only provide one of the two possible network interfaces at a time, although provision of multiple interfaces is unlikely to affect server performance.

Measurements on pure CORBA solutions suggested that CORBA is too heavy-weight for communication with the physical switch. This realization prompted implementation of a raw socket interface for comparison. However, to be plugged directly into the TINA framework, the server is required to provide a CORBA interface. Therefore, a hybrid solution was tested, using a CORBA proxy running on a fast workstation, and forwarding requests to the switch by means of the raw socket interface.

Two different ORBs were used, both written entirely in Java, and both providing a full CORBA 2.0 implementation:

- (1) JORB (Java ORB) from HP, in alpha 3 release.
- (2) OrbixWeb2.0 from Iona, providing the Java equivalent of Orbix 2.1.

3.3. Measurements

Table 3 lists performance measurements of Java-based solutions.

As shown in row 1, JORB communication with the asx100 switch allows around 14 RPC/s. This number is clearly limited by three main factors:

- (1) the performance of the switch itself;
- (2) the performance of the JVM;
- (3) the performance of the JORB implementation.

The impact of switch performance can be seen by comparing rows 1 and 3, respectively rows 2 and 4[†]. The performance gain from running the server on brandt-s-1, as opposed to asx100, is 7.6 using JORB, and 5.6 using sockets. Hence, performance of the switch is a bottleneck. That performance gain is larger in the JORB case must be attributed to the larger overhead of parameter marshalling, resulting in more time spent in the JVM.

The impact of virtual machine performance can be seen by comparing rows 9 and 11. Both rows reflect JORB performance with the client and server running on the same machine, in this case a 200 MHz Pentium Pro PC. In row 9, the Symantec virtual machine is used, which includes a just-in-time (JIT) compiler. In row 11, the JDK 1.02 virtual machine is used. The performance gain from the JIT compiler is a factor of 2.6 with JORB communication. Similarly, comparing rows 10 and 12 we get a factor 1.8 with raw socket communication. The difference in speed-up must again be attributed to the different times spent in the JVM, with the socket solution being the most I/O bound of the two solutions. In any case, similar speed-ups can in principle be achieved on the switch by applying a better JVM, although one is currently not available.

The impact of using CORBA is apparent from comparing any odd-numbered row with its following even-numbered row. For example, the difference between rows 1 and 2 is a speed-up of 7.6 from replacing JORB communication by socket communication. However, this overhead also depends on the CORBA implementation. For example, as can be seen by comparing rows 9 and 15, OrbixWeb performs roughly 50% better than JORB, and would therefore allow around 20 RPC/s against the switch.

It will be interesting to compare these two CORBA implementations to the future ReTINA real-time DPE based on the Chorus-Cool platform when it is available.

3.4. A hybrid solution

Given the above measurements, it might be worthwhile trying to combine the best of the two worlds: the speed of plain sockets with the plugability of CORBA/IIOP. To that end, we tried to insert a proxy object between the client and the server. By means of JORB, the client makes calls to a proxy running on a fast PC on top of a fast Java implementation. The proxy then forwards the request to the real server, using the raw socket interface. Measurements are shown in table 4.

As can be seen, this approach does in fact combine the best of the two solutions. Combined with a fast JIT-based

[†] This assumes that JDK 1.01 and JDK 1.02 exhibit similar performance. This has not been tested.

Table 2.

Virtual machine	JDK version	Used on	Notes
JDK 1.01	1.01	asx100	SunOS port by Hewlett-Packard laboratories
JDK 1.01	1.01	torino	HP-UX 9.x port (alpha)
JDK 1.02	1.02	brandt-s-1/vecheo	Official Sun release
Kaffe 0.6.0	1.02	asx100	Public Domain [10]
Symantec Café 1.51	1.02	brandt-s-1/vecheo	Just-In-Time compiler

Table 3.

	Network interface	Client	Server	RPC/s
1	JORB	torino/JDK 1.01	asx100/JDK 1.01	14
2	Socket	torino/JDK 1.01	asx100/JDK 1.01	106
3	JORB	torino/JDK 1.01	brandt-s-1/JDK 1.02	107
4	Socket	torino/JDK 1.01	brandt-s-1/JDK 1.02	597
5	JORB	brandt-s-1/Symantec	asx100/JDK 1.01	16
6	Socket	brandt-s-1/Symantec	asx100/JDK 1.01	121
7	JORB	brandt-s-1/Symantec	vecheo/Symantec	401
8	Socket	brandt-s-1/Symantec	vecheo/Symantec	1330
9	JORB	brandt-s-1/Symantec	brandt-s-1/Symantec	418
10	Socket	brandt-s-1/Symantec	brandt-s-1/Symantec	1818
11	JORB	brandt-s-1/JDK 1.02	brandt-s-1/JDK 1.02	159
12	Socket	brandt-s-1/JDK 1.02	brandt-s-1/JDK 1.02	998
13	JORB	brandt-s-1/Symantec	asx100/Kaffe	10
14	Socket	brandt-s-1/Symantec	asx100/Kaffe	105
15	OrbixWeb2.0	brandt-s-1/Symantec	brandt-s-1/Symantec	600

Table 4.

	Client	Proxy	Server	RPC/s
15	brandt-s-1/Symantec	brandt-s-1/Symantec	asx100/JDK 1.01	97

Java implementation on the switch, around 200 connection creations/s should be comfortably within reach.

3.5. A server implementation in C

Finally, we implemented the server in C. The performance of the server written in C was measured in two different configurations: with the client connecting directly to the server using sockets, and with the client connecting via JORB to a proxy server, which then connects to the server using sockets. The C server does not provide a CORBA interface. The performance is shown in table 5.

Comparing rows 15 and 16, we see that rewriting the connection server in C provides a speed-up factor of 2.6 in a proxy-based solution. With the client calling the server directly by means of raw sockets, the C server provides a speed-up of a factor of 5, as can be seen by comparing rows 6 and 17. It should be noted that this speed-up would decrease to between 2 and 3 if a faster virtual machine could be used on the switch. However, this relatively small speed-up reflects the I/O bound nature of the server. Hence, if the computational complexity of the server was larger, a larger speed-up could be expected.

3.6. SNMP

The performance of connection creation and deletion by means of SNMP was measured by Matthieu Goutet of Hewlett-Packard laboratories. Using a Smalltalk client and communicating with the switch by means of SNMP, around 50 connection creations or deletions/s was achieved. Assuming that Smalltalk exhibits performance similar to Java (this is ignoring the fact that many Smalltalk systems include a compiler), this number is comparable with row 17, providing around 600 connection creations/s. Hence, using a connection server instead of SNMP, a performance of an order of magnitude better can be achieved.

Although we have only performed measurements on the ASX-100, others quote SNMP performance around two connection creations or deletions on other switches. This seems to suggest that there are large differences between the performance of SNMP on different switches, and that SNMP implementations can therefore in general be improved with respect to performance. We would expect to see optimized implementations of SNMP on switches where the interface is being used for switch control.

Another optimization would be an attempt to overcome the SNMP feature that multiple get/set attribute operations have to be made in multiple calls to the interface. One obvious way to address this would be to have an object on

Table 5.

	Client	Proxy	Server	RPC/s
16	brandt-s-1/Symantec	brandt-s-1/Symantec	asx100/C	251
17	brandt-s-1/Symantec	Direct (no proxy)	asx100/C	600

the switch that could handle multiple attribute operations as a bundle and that could then make the multiple SNMP calls locally.

To be fair to SNMP, we should not forget that it is significantly more general than our protocol, it applies to most or all switches for private networks (switches for public networks may support the CMIP protocol), and was furthermore not conceived as an alternative to connection set-up by means of signalling. Nevertheless, although it can be improved, SNMP performance on the ASX-100 is probably sufficient for many purposes, maybe even including the implementation of ‘open switching’ approaches such as TINA and Xbind. In any case, as we will discuss below, SNMP is unlikely to be a major bottleneck in the BVPN demonstrator.

4. Discussion

4.1. SNMP performance

Our measurements demonstrate that SNMP performance is indeed inferior to the connection server used in this experiment. On the test equipment used, SNMP achieves around 50 connection creations/s, which should be compared with around 600 connection creations/s achieved by the connection server implemented in C.

However, only a few switches actually allow the implementation of a connection server by means of an open switching interface such as that used within Hewlett-Packard laboratories. Our ability to put this interface onto the ASX-100 is thus an exception.

4.2. Deployment of TINA computational objects

The inferior performance of the switch hardware as compared with a modern PC or UNIX workstation suggests that it would be unwise to run a full EML-CP on the switch itself, although this is indeed possible. A more viable approach is to run the EML-CP on a high-performance PC or workstation, only communicating with the switch itself when necessary to actually create or delete physical connections. This separation of concerns also makes sense from a software engineering point of view, since it allows different switches to be controlled by the same (powerful) controller hosting the corresponding EML-CP object implementations.

The nature of the communication from the EML-CP object to the physical switch can then be achieved by means of SNMP or other mechanisms such as GSMP [9] or the connection server used in this experiment. The latter may achieve a performance of an order or magnitude better than SNMP.

4.3. The TINA CM software bottleneck

Given our SNMP performance measurements, it is our conclusion that the overhead experienced by users of the CM software is not caused by SNMP.

The Alcin CM software has been experienced mainly by Telenor within the BVPN demonstrator. This software includes computational objects: connection coordinator (CC) that controls the whole underlying network, layer network coordinator (LNC) and network management layer–connection performer (NML-CP) described in section 1.2. The EML-CP objects that control the Fore switches have been provided by Telenor.

The performance of the connection creation by means of Alcin was measured within the Alcatel laboratory. A client object, the Alcin package and one ‘dummy’ EML-CP were executed on top of ORBIX 2.0 MT on one Ultra1 Sparc station (running Solaris 2.5 with 64 Mb Ram). The EML-CP is called ‘dummy’ because it is not interconnected to the actual corresponding switch (see figure 5).

A (point-to-point) connection was created in around 28 ms, which would correspond to 35 creations/s. This figure has to be compared with the 50 creations/s provided by SNMP. This result can be explained by the number of remote operations required to establish one connection: nine are required between the client and the Alcin package and four are required between Alcin and the EML-CP object. Indeed, a remote interaction between two ORBIX objects on an Ultra1 Sparc station takes approximately 2 ms. Therefore approximately 90% of the time is spent in ORBIX transport. Furthermore, we noted that this result had remained valid when the client, the Alcin package and the EML-CP were executed on different computers.

In order to validate this conclusion, we carried out the same experiment with the Alcin package and the EML-CP co-located in the same address space (one process). With such a configuration a connection was created with nine invocations in around 20 ms (50 creations/s). Again approximately 90% of the time was spent in transport. Therefore the first conclusion is that the use of Alcin on top of SNMP to control one switch may not be seen as a real critical bottleneck (factor of 1.4).

However, Alcin allows the control of a global layer network, i.e. of a set of switches. In its current version, the Alcin software controls the set of EML-CPs in a sequential manner. Thus, the duration of an end-to-end connection is proportional to the number of involved switches. Since four interactions (8 ms) are required to establish a sub-connection on each EML-CP, we may estimate the duration of the creation for an end-to-end connection crossing three EML-CPs to approximately 44 ms (24 ms plus 20 ms for the interactions between the client and Alcin). This ‘sequential’ behaviour is the main cause of the overhead

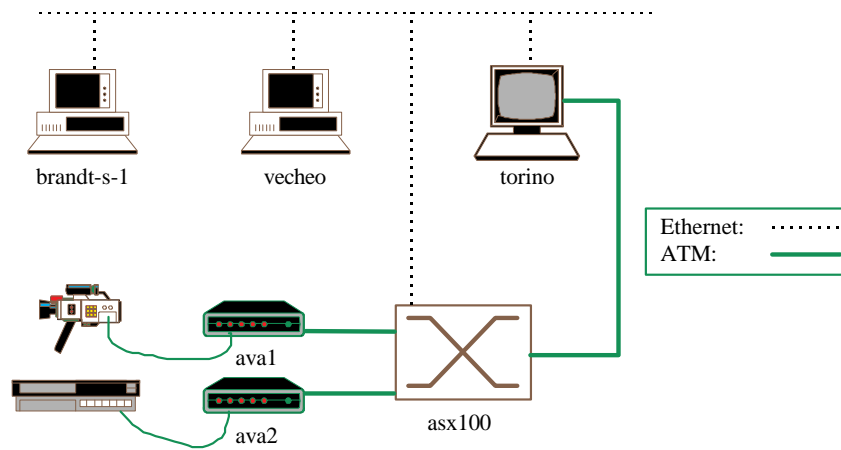


Figure 5. Experiment set-up.

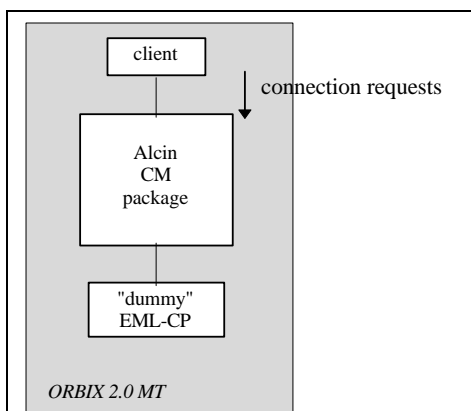


Figure 6. Experiment configuration.

experienced by Alcin users and creates a bottleneck in the BVPN demonstrator.

Alcin was provided as the ReTINA CM first release. It is important to note that no specific effort has been made to support non-functional requirements such as latency of connections establishment or real-time control. The improvement of the CM non-functional behaviour is one of the objectives for the second release. This will be done through the re-engineering of the CM computational objects, reducing the number of needed invocations, processing multiple requests concurrently and parallelizing the establishment of sub-connections.

Further experiments with this new CM should be carried out in order to measure the overheads imposed by the use of the CM software, both for the control of one switch and for the establishment of an end-to-end connection. However, the quantitative results presented above allow us to be optimistic. Indeed, if we could reduce the number of required invocations to two (one between the client and Alcin, and one between Alcin and the EML-CP), the average time for the creation of a connection on a switch could be around 5 ms, or maybe 10 ms if the process time was not so negligible. Thanks to this approach and the introduction of parallelism, the use of ReTINA CM with an

effective switch control interface such as, for example, the proxy-based hybrid solution would become conceivable.

5. Future directions

We now have a fairly complete understanding of the link-by-link CM issues and the next phase of the project will address the wider and more difficult issue of end-to-end connections and the QoS attributes of these connections. Within ReTINA we will also implement a complete application that can use the mechanisms provided by the QoS manager. This application will be a BVPN. The BVPN demonstrator addresses the problem of how to open up connection management in public networks in order to allow third-party providers and customers to create 'leased lines' on demand without the need to contact the telecommunications operator.

Two of the design aims of ATM were to support multiple traffic classes simultaneously and to make all resources in the network available to the client services. ATM uses a statistical multiplexing approach to the simultaneous support of multiple connections. Also, connections do not reserve the network resources they need for the entire duration of their lifetime. These efficiency-motivated aspects have two important effects in an ATM network:

- congestion is possible when excessive amounts of traffic contend for limited buffer resources (this results in lost cells);
- variable cell delays occur as the cells traverse the network. These are caused mainly by random queuing delays at each switch and multiplexer (the queues themselves are there to avoid excessive cell loss).

Our challenge will be how to protect the QoS of those multiple traffic classes, while simultaneously optimizing the statistical sharing of the network resources.

5.1. ReTINA quality-of-service management

ReTINA's focus is on real-time services and so we are concentrating on the quality of service required in the context of high-speed networks and distributed multimedia

applications. The objective of the QoS manager is to assure the quality of connections in accordance with the requirements of the session customer on whose behalf the connection is being made. ReTINA will adopt a contractual approach to the establishment of connections with requested QoS attributes.

Under our architecture a user application makes a connection request based on a particular traffic type and a set of parameters characterizing the kind of quality of service desired. The request goes to a connection manager and this request is hierarchically partitioned across the network until the decision is taken locally at each switch involved in the path. At each switch local CAC and resource reservation decisions will be made. The strategies differ from switch to switch. We will develop simple mathematical models of switch behaviour and a calculus for coordinating the local QoS capacities into end-to-end estimates. We are also developing algorithms for the aggregation of the figures over our subnetwork hierarchy. These QoS capacity estimates will be used in high-level admission control and in QoS-based routing.

Acknowledgments

Peter Graubmann of Siemens AG produced the first requirements paper for the QoS work. Steve Vogelsang of Fore Systems answered queries on ATM switch capabilities.

References

- [1] Prycker M de 1995 *Asynchronous Transfer Mode* (Englewood Cliffs, NJ: Prentice-Hall)
- [2] Chen T M and Liu S S 1995 *ATM Switching Systems* (London: Artech)
- [3] Lazar A, Bhonsle S K and Lim K-S 1995 *A Binding Architecture for Multimedia Networks* (New York: Center for Telecommunications Research, Columbia University)
- [4] 1996 *Project Xbind* Columbia University in the City of New York <http://comet.ctr.columbia.edu/opensig>
- [5] TINA <http://www.tinac.com>
- [6] Lazar A and Marconcini F 1996 *Towards an Open API for ATM Switch Control* (New York: Comet Group, Center for Telecommunications Research, Columbia University)
- [7] Dahle E and Solbakken H 1996 *BVPN Service Enterprise Specification* (Telenor)
- [8] 1996 *IP Switching: The Intelligence of Routing, The Performance of Switching* (Ipsilon)
- [9] <http://www.tjwassoc.demon.co.uk/kaffe/kaffe.htm>
- [10] Research Networks Section, Naval Research Laboratory (ftp site: [hsdndev.harvard.edu/pub/mankin](ftp://hsdndev.harvard.edu/pub/mankin))
- [11] 1994 *Network Resource Information Model Specification* No TB.LR.001.2.1.95, TINA-C
- [12] 1996 *DPE Services Specifications (Telecom Services)* ReTINA ACTS Project AC048 Deliverable RT/TR-96-10.1/WP3
- [13] Ginsburg D 1996 *ATM Solutions for Enterprise Networking* (Reading, MA: Addison-Wesley)