

Flow Labelled IP: A Connectionless Approach to ATM*

Peter Newman, Tom Lyon, and Greg Minshall[†]

Ipsilon Networks Inc.[‡]

Abstract

A number of proposals for supporting IP over ATM are under discussion in the networking community including: LAN emulation, classical IP over ATM, routing over large clouds, and multiprotocol over ATM. Each of these proposals hides the real network topology from the IP layer by treating the data link layer as a large, opaque, network cloud. We argue that this leads to complexity, inefficiency and duplication of functionality in the resulting network.

We propose an alternative in which we discard the connection oriented nature of ATM and integrate fast ATM hardware directly with IP, preserving the connectionless nature of IP. We use “soft” state in the ATM hardware to cache the IP forwarding decision. This enables further traffic on the same IP flow to be switched by the ATM hardware rather than forwarded by IP software. We claim that this approach combines the simplicity, scalability, and robustness of IP with the speed, capacity, and multiservice traffic capabilities of ATM.

1. Introduction

Asynchronous Transfer Mode (ATM) has recently received much attention because of its high capacity, its bandwidth scalability, and its ability to support multiservice traffic. However, ATM is connection oriented whereas the vast majority of modern data networking protocols are connectionless. This mismatch has led to complexity, inefficiency, and duplication of functionality in attempting to apply ATM technology to data communication.

The Internet Protocol (IP) has also seen very rapid growth in the last several years. Current research suggests that given a suitable implementation, IP is no less capable of supporting real-time and multimedia applications than ATM [1, 2]. Much attention is being focused on the use of IP multicast for multimedia and conferencing applications [3]. Furthermore, many believe that the connectionless model on which IP is based, with the addition of soft-state for traffic management, is a much more robust and flexible basis on which to construct an integrated services network.

In this paper we investigate the implementation of IP directly on top of ATM hardware while preserving the connectionless model of IP. We discard the connection oriented nature of the ATM protocol stack and couple fast ATM switching hardware directly to IP. This has the particular advantage of not requiring a signalling protocol, or any address resolution protocol, and requiring only the standard IP routing protocols — protocols that have been well debugged and heavily tested. It also directly supports IP multicast which is currently incompatible with the ATM implementation of multicast. Of course, it rather assumes that IP, instead of ATM, be the underlying universal, ubiquitous, integrated services protocol, but some would claim that this is already the case.

2. ATM under IP

2.1 Why IP?

Because it's there!

The Internet is currently connected to approximately 5 million hosts on 45 thousand interconnected networks covering 86 countries and continues to grow

*© 1996 IEEE. Published in Proc. IEEE Infocom, San Francisco, March 1996. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

[†]{pn,pugs,minshall}@ipsilon.com

[‡]<<http://www.ipsilon.com>>

exponentially. The traditional design of packet switch (router) on which the Internet is based is beginning to run out of steam. Routers are expensive and of limited throughput when compared to switches. To support the continued traffic demand of the Internet, IP needs to go faster and cost less. To support the increasing demand for real-time and multimedia applications IP also needs to support quality of service (QOS) selection. We believe that both can be met by the application of switching technology to the design of an IP router. Our objective is simply to make IP go fast and offer QOS support by integrating state-of-the-art switching technology with IP routing and forwarding. We aim to combine the flexibility of IP with the speed of switching.

2.2 Why ATM?

Switching offers the ability to support high bandwidth links at wire speed and switches of very high aggregate throughput may be implemented. A switch achieves this high capacity by implementing the data path completely in hardware. Thus there is a tradeoff between the flexibility of the router, which can make an independent forwarding decision for every packet, and the speed of the switch, which requires state to be pre-established in its forwarding tables.

ATM offers scalability of both link bandwidth and switch capacity. It is also well suited to the application of VLSI implementation and the widespread, almost frenzied, interest in ATM promises the rapid decrease in cost that comes from volume production. In short, we choose ATM because the hardware is now standardized and available, it is fast, and the price tag is falling.

2.3 Why connectionless?

The considerable success of the Internet is in large part due to the connectionless nature of IP. IP is built on a very low level building block — datagram forwarding. No assumptions are made regarding the services provided by the underlying network beyond the ability to forward a datagram in the direction of the destination. This has permitted IP to operate over a very wide range of underlying network technologies. Multiple types of communications service are offered by enhancement of the basic forwarding service. Datagram forwarding requires no state to be maintained for individual connections. This has proven extremely robust in the presence of failures.

In a connection oriented network, possibly as much as 90% of the signalling code is there to handle error conditions. This code is impossible to thoroughly test and is almost never totally correct. A soft state approach in which state in the network is periodically refreshed covers

a large spectrum of possible error conditions with a very simple recovery mechanism.

3. Connectionless service implementation

Since ATM is itself connection oriented, the heart of the problem is to make use of the speed and capacity of the switching hardware without sacrificing the scalability and flexibility that come from the connectionless nature of IP. A number of approaches to the implementation of IP over ATM have been proposed in the literature and in the standards bodies [4] which we now review before presenting our own approach.

3.1 LAN emulation

In order to deploy ATM technology within the local area as quickly as possible it was proposed that ATM LANs emulate the architecture of the “legacy” shared medium LANs [5–7]. ATM emulates the physical shared medium by establishing an ATM multicast group between all of the end stations that belong to the ATM LAN segment. An address resolution server translates the 48-bit MAC address into an ATM address. Once the ATM address of the destination has been discovered, a point-to-point ATM virtual connection may be established to the destination using the ATM signalling protocol.

There are certain disadvantages to this solution. LAN emulation is a bridging solution and it is well known that bridging does not scale well to large networks. The incompatibility between Ethernet/IEEE 802.3 and Token Ring is faithfully reproduced by LAN emulation. The multicast service and the address resolution service are both implemented by servers. This represents a single point of failure until a design for redundancy and the distribution of the server database is developed. The solution is complex. ATM requires at least 40,000 lines of signalling code on the host, and the LAN emulation client and server each require 10,000–20,000 lines of code, all to emulate a service typically implemented in an Ethernet device driver in about 2,000 lines. Also, network management is complicated by the need to separately manage both the legacy components of the architecture and the new ATM components.

3.2 Classical IP over ATM

While LAN emulation uses ATM to emulate the properties of a shared medium LAN, classical IP over ATM [8, 9] replaces the shared medium LAN with an ATM subnet. However, the LAN-based paradigm of “classical” IP is preserved. Thus a group of ATM connected hosts is gathered together and treated as a Logical IP Subnetwork (LIS) in much the same way as a group of hosts connected to a shared medium LAN is

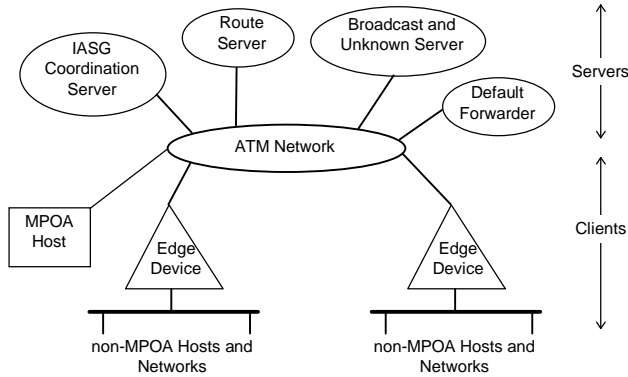


Figure 1: Multiprotocol Over ATM Architecture

viewed as a subnet. Logical IP subnets must be interconnected via routers regardless of whether direct connection via ATM is possible.

Classical IP over ATM offers no broadcast service within the logical subnet and the multicast service is still under discussion [10]. Many applications developed within the context of the classical IP model assume the ability to broadcast to all hosts on the local subnet. This broadcast ability is often used for auto-configuration and its absence in an ATM subnet results in increased configuration requirements. Both LAN emulation and classical IP over ATM are simple and require no modifications to IP but they do not scale well to large networks because all communication between ATM LAN segments or logical subnets must proceed via routers. These routers become a significant throughput bottleneck.

3.3 Routing over large clouds

Routing over large clouds addresses the issue of communication between different logical subnets within a large homogenous network, called a cloud or a non-broadcast multi-access (NBMA) network, such as ATM. The problem consists of locating the exit point on the cloud “nearest” to any given destination and to obtain the ATM address for that exit point. The signalling protocol may then be used to establish an ATM connection across the ATM cloud to the exit point.

The NBMA address resolution protocol (NARP) [11] offers a server based solution which bears some similarity to the classical IP over ATM address resolution server. The NBMA next hop resolution protocol (NHRP) [12] offers a superset of the NARP functionality where the servers are themselves routers designed to offer an address resolution service rather than a packet forwarding service.

3.4 Multiprotocol over ATM

Multiprotocol over ATM (MPOA) is all of the above. It is a combination of the basic concepts of LAN

emulation, classical IP over ATM, and NHRP, integrated into a single protocol specification. MPOA also adds protocols to replicate servers and distribute the database for reasons of capacity and availability [13].

The logical components of the MPOA architecture are shown in Fig. 1. An edge device connects legacy subnets to the ATM network. It is typically a simple bridge/router with a cache of routes and address translations that it obtains from the MPOA servers. Hosts and edge devices are grouped together into IASGs[§] which are the MPOA equivalent of an IP subnet or an emulated LAN. The IASG coordination server offers address resolution and MAC layer bridging for the hosts and edge devices in the IASGs it supports. A frame whose MAC address is unknown is passed to the broadcast/unknown server to be forwarded to all ports on all edge devices within the IASG. For network layer packet forwarding an edge device will query a route server to obtain the ATM address of the exit point nearest the destination and will cache the information for future use. The route servers run multiple internetwork layer routing protocols and the service operates in a similar manner to NHRP. Network layer multicast is handled by the broadcast/unknown server.

4. In search of the connectionless connection

4.1 Obscured by clouds

Both routing over large clouds and multiprotocol over ATM seek to reap the performance benefits of direct interconnection of the data path at the link layer (i.e. switching) across a large network, without the need for packet forwarding at the internetwork layer. This stretches the architectural model of classical IP. Classical IP assumes that subnets are interconnected at the network layer by IP routers and that anything within a subnet is reachable at the link layer, typically with a shared medium LAN. By seeking to achieve direct connectivity across subnets we are no longer conforming to this model. Both approaches obscure the real topology of the underlying network from the internetwork layer routing protocol. To IP the physical network becomes a large opaque cloud which results in some significant problems.

First there is a duplication of functionality. Both IP and ATM each require routing protocols. Not only does this imply duplication of the routing protocols but it also leads to duplication of the management functions. In addition, management functions are required to handle the

[§] See [13] for the derivation of the acronyms. As a general guide append “FG” to everything as a reminder that components are logical “functional groups” rather than physical devices.

interaction between the two. This makes it much more difficult to locate problems. When connectivity is lost it is much more difficult to determine where the fault lies if two separate routing protocols are involved. It is also possible for undetected routing loops to be formed in certain situations [9].

For efficient multicast capability IP requires knowledge of the underlying network topology. Without this information, a multicast packet arriving at a router for transmission to leaf nodes attached to the cloud must be replicated at the router and each copy must be transmitted separately across the cloud. This is clearly inefficient as multiple copies of the same packet will be transmitted unnecessarily across data links within the cloud. The replication function would be far better implemented at the appropriate forks in the multicast tree within the cloud. For low bandwidth datagram traffic this is simply inefficient but it could prove very disruptive for higher bandwidth real-time traffic with quality of service requirements, such as IP voice and video multimedia applications.

All routers and route servers connected to a large cloud may be considered to be logically one hop away from each other. Conventional routing protocols have N^2 scaling difficulties in the case where each router has many neighbors, arising from the routing table size, the amount of routing update processing, and the amount of routing update traffic generated. But if the routers are not logically fully meshed the reliability of the network is reduced because it is possible that two hosts that are physically connected have lost logical connectivity. Introducing logical connectivity on top of a physical network can only reduce the reliability since more systems need to be functioning to achieve connectivity.

In a physical network it is a simple matter for a router to discover its neighbors and then fire up a routing protocol to connect to the network. It is also simple for a host to discover what services are available on the network. In a cloud environment, however, the cloud can be of any arbitrary size and topology, so a router cannot discover its neighbors, it must be assigned them. It also makes it much more difficult for a host to automatically discover network services. This lack of support for auto-configuration leads to greatly increased management and manual configuration requirements. Also, in a cloud model, routers are required to interconnect multiple logical subnets. This requires the configuration of multiple logical interfaces on a single physical interface — the “one-armed router” — which also increases the configuration requirements.

4.2 Flows and soft state

The concept of a flow has emerged within the IP community over the past few years. A flow is a sequence of packets sent from a particular source to a particular (unicast or multicast) destination that are related in terms of their routing and any local handling policy they may require [14]. It performs a similar function in a connectionless network to the role the connection plays in a connection oriented network. In IP version 6 the inclusion of a flow label in the packet header allows the forwarding process to be enhanced by caching routing decisions. Also network resources may be reserved on behalf of a flow to offer quality of service guarantees.

IP is connectionless but many applications above IP employ a connection oriented transport protocol. The most efficient mapping of IP onto ATM must consider the characteristics of the application, or at least the transport protocol, in deciding whether to establish an end-to-end ATM connection on behalf of any specific flow [9, 15]. Flows carrying real-time traffic, flows with quality of service requirements, or flows likely to have a long holding time, will be handled most efficiently by mapping them into an individual ATM connection. Short duration flows, and database queries would best be handled by connectionless hop-by-hop packet forwarding between IP routers using shared, pre-established ATM connections between the routers. This is particularly true for exchanges such as DNS lookups that consist of a single packet in each direction. Establishing an end-to-end ATM connection for every IP packet flow would impose a heavy load on the ATM signalling protocol and impose unnecessary delay on query-response traffic.

One of the reasons that IP scales well to large networks is due to its connectionless nature. If a router or a link fails in a moderately well connected network, IP simply routes around the failure. If we establish end-to-end connections across an ATM cloud the failure of a link or router will invalidate all associated connections. This will exert a heavy load on the signalling protocol to re-establish all of the ATM connection state. Also, many connections that are not directly associated with the failed component will become sub-optimal, perhaps highly sub-optimal when the topology changes. It is also possible that routing loops can form after a topology change until the old routing information is purged from the address resolution servers and route servers [16].

It is clear that in order to take advantage of the efficiency of switching at the data link layer, and to offer quality of service guarantees, state information must be maintained within the switches. However, the simplicity and robustness of IP is much more likely to be preserved if the state is maintained locally rather than on an end-to-end basis and if the state is “soft” rather than “hard”. Soft

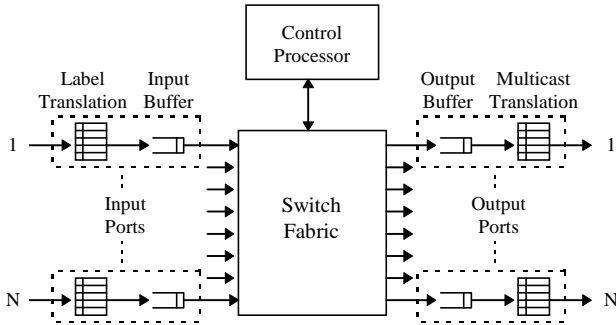


Figure 2: General Structure of an ATM Switch

state is state information that is installed within a network for reasons of performance enhancement but is not crucial to the correct operation of the network [2, 17]. It is typically designed to be refreshed periodically such that many possible error conditions may be corrected by simply timing out old state. This leads us to consider the possibility of a connectionless implementation of ATM.

4.3 Connectionless ATM

The general structure of an ATM switch is illustrated in Fig. 2 [18]. Each input port has a hardware lookup table indexed by the ATM label field (VPI/VCI) in the incoming cell. For unicast traffic each entry in the table contains the new label to be used on the next hop, the output port on which to transmit the cell, and information related to the quality of service the cell should receive. Entries in the table are established by the control processor in response to a user signalling connection request. The control processor needs to run the signalling protocol and a routing protocol in order to establish connections and the network signalling protocol itself must be a connectionless datagram transfer protocol in which the datagrams are signalling protocol elements.

However, the general structure of a high performance router is very similar to that of an ATM switch. Each input port has a hardware lookup table indexed by the header field of incoming packets. Each entry in the table specifies the output interface on which the packet should be transmitted and may also contain other information related to the processing of the packet. Entries in the table are established by the control processor. The control processor needs to run a routing protocol in order to determine on which output port to forward each packet.

The major difference between the router and the ATM switch is that the input port tables of the router are treated as a cache. If there is no entry in the input port table for an incoming packet a request is made to the control processor for forwarding information. Also, entries in the table are removed after a certain timeout interval in order to refresh the cached information in case the route has changed. “Soft” state information is maintained in the input port

tables of the router. The cache is implemented in the input port to speed up packet processing, it does not affect the underlying connectionless nature of the router.

Given that the general structure of the router and the ATM switch are fundamentally the same one might consider whether the input port lookup table of the ATM switch may be treated as a cache — whether the table entries may be treated as a local cache of “soft” state information. Our motivation for using ATM switching at the data link layer is: its ability to support high bandwidth links and very high aggregate capacity switches; the independence of the switching mechanism from the bandwidth of the link; and its ability to support quality of service guarantees. All of these benefits are a direct result of implementing the data path in hardware. Provided the correct information gets loaded into the input port lookup tables it does not make any difference to the data path whether a “hard” state, end-to-end connection oriented approach is taken, or a “soft” state, connectionless cached approach. However, it makes a substantial difference to the properties of the resulting network.

5. Flow labelled IP

Our goal is to achieve the most efficient implementation of IP on top of fast switching hardware. We now consider using standard ATM hardware but completely changing the control software in order to operate each switch in a connectionless manner. The result will be a router with attached switching hardware that has the ability to cache routing decisions in the switching hardware. We call this an *IP switch* since it allows packet flows to be switched, bypassing the router, once the routing information has been cached in the switch.

5.1 A switch by any other name ...

Consider the general structure of an ATM switch illustrated in Fig. 2. To construct an IP switch we take the hardware as it stands, without any modification, but completely remove the software resident in the control processor above AAL-5. Thus we remove the signalling, any existing routing protocol, and any LAN emulation server or address resolution servers, etc. In place of the ATM software we load a standard IP router software package. We will need to define a simple flow management protocol (IFMP) to associate IP flows with ATM virtual channels, a flow classifier to decide whether to switch each flow, and a driver to control the switch hardware.

At system startup a default forwarding ATM virtual channel is established between the IP routing software running on the control processor and that of each of its

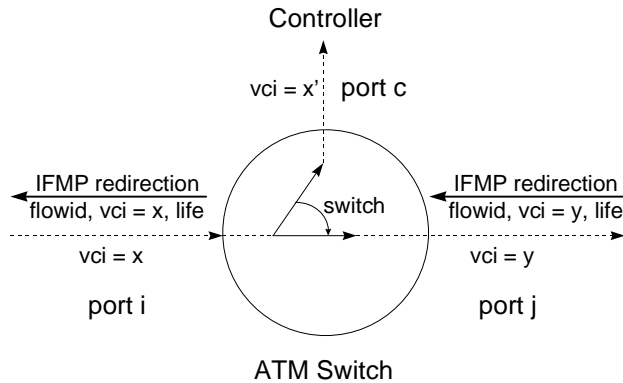


Figure 3: Establishing a switched flow

neighbors. The default forwarding channel is used for the hop-by-hop connectionless forwarding of IP datagrams. Thus we now have the ability to forward IP packets but to gain the benefit of the switching hardware we need a mechanism to associate an IP flow with a specific ATM label (virtual path and virtual channel identifier VPI/VCI).

5.2 Flow classification

We characterize an IP flow according to the fields in the IP/TCP/UDP header that determine the routing decision such as: type of service, protocol, source address, destination address, source port, destination port, etc. Two packets belong to the same flow if the values of these fields are identical. Several different flow types may be defined each characterized by a different set of header fields (though the set of flow types must be ordered so that one can perform a most specific match operation.)

When a packet is received across a default forwarding channel it is reassembled and submitted to the control processor for forwarding. The processor forwards the packet in the normal manner but it also performs a flow classification on the packet to determine whether future packets belonging to the same flow should be switched directly in the ATM hardware or continue to be forwarded hop-by-hop by the router software.

Flow classification is a local policy decision. The flow classifier inspects the contents of the fields that characterize the flow and makes its decision based upon a local policy expressed in a table. For example, by looking for well known source or destination port numbers one can identify the application. Flows belonging to FTP data connections may be configured to be switched, but DNS queries could be forwarded as datagrams. (The performance of such a classification is investigated in the following section.)

5.3 Flow management protocol (IFMP)

If the processor decides that the flow should be switched it selects a free label (label x) from the label

space of the input port (port i) on which the packet was received, fig. 3. We make the assumption that virtual channels are unidirectional so the ATM label space (VPI/VCI range) of the incoming direction of each link is owned by the input port to which it is connected. The processor also selects a free label (label x') on its control port (port c). (The control port is the port, either real or virtual, by which the control processor is connected to the switch.) The switch driver is then instructed to map label x on input port i to label x' on the control port c .

After making the entry in the translation table of the switch input port the processor sends an IFMP [19] Redirection message upstream to the previous hop from which the packet came. The redirection message contains the label x , a flow identifier, and a lifetime. The flow identifier contains the set of header fields that characterize the flow. The redirection message requests the upstream host or router to transmit all further packets with header fields that match those specified in the flow identifier on the ATM virtual channel specified by the label. The lifetime field specifies the length of time for which this redirection is valid. Unless the flow state is refreshed, when the lifetime expires, this binding of flow and label should be deleted and further packets belonging to the flow will be transmitted on the default forwarding channel.

From this point packets belonging to the flow will arrive at the switch controller, port c , with the ATM VPI/VCI label x' . The packets will still be reassembled and forwarded by the IP forwarding software but the process is speeded up because the previous routing decision for this flow was cached in the router software and can be indexed by the label x' .

The real benefit of switching comes when the downstream router or host also runs the same redirection algorithm. When the router receives a redirection message from its downstream neighbor on port j , redirecting the flow to label y , it can switch all further traffic belonging to that flow directly within the ATM hardware. The router does this by instructing the switch to map label x on port i to label y on port j . Thus the traffic is no longer sent to the control processor but is switched directly to the required output port.

When an IP switch accepts a redirection message it also changes the encapsulation it uses for the redirected flow. The encapsulation used for IP packets on the default forwarding channel is the standard LLC/SNAP encapsulation over AAL-5. The encapsulation used for each IP packet on a flow redirected to a specific virtual channel removes all of the header fields that characterize the flow from the header of each packet [20]. The IP packet with the resulting compressed header is then encapsulated in AAL-5 and transmitted on the specified

virtual channel. The fields that are removed are stored by the router that issued the redirection and are associated with the specified ATM virtual channel. The complete packet, including TTL and checksum fields, may be reconstructed using the incoming label to access the stored header fields. This approach is taken for security reasons. It allows an IP switch to act as a security firewall without having to inspect the contents of each packet. It prevents a user from establishing a switched flow to a permitted destination or service behind a firewall and then changing the IP packet header to gain access to a prohibited destination.

Conceptually, each IP switch maintains a background refresh timer. When the background refresh timer expires, the state of every flow is examined. If a flow has received traffic since the last refresh period its state is refreshed. Flow state is refreshed by sending a redirect message upstream with the same label and flow identifier as the original and a new lifetime. If a flow has received no traffic since the last refresh period its cached state is removed. This will involve issuing an IFMP Reclaim message upstream to reclaim the label for reuse. The flow state is not deleted until an IFMP Reclaim Ack message is received to acknowledge release of the requested label. (Reclaim messages may also be used to release labels in use if the free label space is close to exhaustion.) For flows that are labelled, but not switched, the control processor can examine its own state to see whether the flow has received any traffic in the previous refresh period. For flows that are switched the control processor must query the switch hardware to discover whether a specific channel has recently been active.

The flow management protocol is advisory in nature. The decision to accept a redirection request is local and redirection messages may be ignored. Redirection messages are not acknowledged since the first packet arriving on the new virtual channel will indicate acceptance of the request. The protocol is also symmetric in that no distinction is made between a user interface (UNI) and a network interface (NNI). This leads to a very simple implementation.

5.4 Point-to-point

LAN emulation, and classical IP over ATM, etc. seek to establish a logical shared medium network model on top of ATM. However, we propose a point-to-point network model — a much more natural model for ATM. All routing protocols deal well with point-to-point links. Point-to-point links existed in IP before the advent of Ethernet multi-access broadcast links and there may be as many point-to-point links (SLIP and PPP) in the Internet today as there are shared medium links.

We have adopted a point-to-point network model rather than a cloud model. The Internet is proof that IP can scale to very large networks without requiring the concept of a data-link cloud. Also we have been careful to separate the act of labelling a flow from that of switching a flow. Choosing to switch a labelled flow is a purely local decision. From outside an IP switch, one cannot determine whether a particular flow has been switched or forwarded other than its increased performance. This separation of labelling and switching, and the local nature of the switching decision, ensures scalability to large networks. The labelling or switching decision for any particular link has no effect on the rest of the network.

5.5 Multicast

An IP switch can support IP multicast without any modification to the Internet Group Management Protocol (IGMP). Flow redirection proceeds in exactly the same manner as for unicast traffic. At an IP switch, where an incoming multicast flow is replicated into a number of branches, each branch may be individually redirected by its downstream neighbor. If the incoming multicast flow is labelled, the multicast capability of the ATM switch may be used directly on those outgoing branches that have redirected the flow onto a specific virtual channel. The switch can also send a copy of the multicast flow to the control processor so that branches that have not decided to redirect the flow may receive their copies of the traffic over the default forwarding channel.

IP multicast offers a multipoint-to-multipoint service. Any sender can transmit traffic to the multicast group. Individual flows, however, are point-to-multipoint since each flow is specific to a single source. ATM hardware only offers a point-to-multipoint multicast service. The result of the flow redirection process for a multicast group will be to establish a point-to-multipoint virtual channel from every sender that has recently transmitted traffic to the group.

5.6 Quality of service

Each IP switch can make its own quality of service decisions according to local policy. Each flow is classified as part of the forwarding operation and quality of service information may be included in the flow classification decision based upon the application, the type of service field in the IP header, the protocol, etc. Each IP switch may also base its quality of service decision on the capabilities of the underlying ATM switch hardware. For current generation switches, separating the traffic into real-time and best-effort flows may be all that can be supported. Future switch designs are likely to be able to offer sophisticated scheduling capabilities [1, 21, 22].

Individual quality of service requests for each flow may also be supported using the resource reservation protocol (RSVP) [23, 24]. RSVP performs a similar function to the ATM signalling protocol in that it can reserve network resources for a particular flow. The major differences are that RSVP is receiver initiated and that it uses a “soft” state approach rather than the “hard” state of ATM signalling. RSVP allows an application to specify the traffic characteristics of a flow using a *flowspec*, similar in nature to the traffic descriptor of ATM traffic management. A reservation request may be accepted or denied by each IP switch in the path using an admission control policy. Resources are reserved by configuring the queueing and scheduling hardware within the ATM switch, and the flow may be policed by configuring the policing (UPC) hardware in the ATM switch according to the *flowspec*. An IP switch should be configured to redirect all flows requesting bandwidth reservation in order to switch them. This will allow the traffic management capabilities of the ATM hardware to be employed to guarantee the requested quality of service.

5.7 Latency

Transmission of the first packet on a new flow across a network of IP switches has the effect of leaving an ATM connection in its wake if all of the IP switches are configured to switch that type of packet. For a connection oriented transport protocol such as TCP the first packet on a new flow will be the SYN packet in the forward direction and the SYN ACK packet in the return direction. This exchange is used in the three-way handshake that establishes the transport connection. In the typical case, for a network of IP switches, by the time the SYN ACK has returned to the source a specific ATM virtual channel will have been established from source to destination in the forward direction. Thus as soon as the first data packet is sent on the new transport connection it will be carried on an ATM virtual connection. If the virtual connection is still in process of being established on any link, data is forwarded by the routing software on the default forwarding channel.

6. Simulation results

The performance gain that results from the integration of an ATM switch with an IP router is somewhat dependent upon the characteristics of the incident traffic. If all of the traffic consists of single packet queries and single packet responses between a very large population of sources and destinations, the ATM switch will add very little. However, even in this situation, our approach will offer better performance than connection based IP over ATM since it can offer connectionless datagram

forwarding. Many applications, such as file transfer and real time audio or video, transmit a significant quantity of information after establishing a connection. Also some applications such as remote login, while not transferring large quantities of information, have long holding times and tend to transmit a large number of small packets. These will benefit from the establishment of a connection within the switch both by removing the per packet processing overhead and by reducing the transfer delay by allocating a higher quality of service to such applications.

To investigate the benefit of our approach we obtained a traffic trace from the Internet backbone^{**}. The trace contains five minutes of traffic taken on Sep 25, 1995. It was taken by monitoring an FDDI ring that connects traffic from the San Francisco Bay Area to and from the Internet backbone. The trace includes a timestamp, IP source and destination addresses, the packet length, and source and destination port numbers for each packet. Backbone traffic is most likely to present a worst case for flow switching because of the large number of independent conversations that are multiplexed together. If we can demonstrate a performance enhancement with Internet backbone traffic, there will be a far greater enhancement possible for campus and corporate backbone traffic.

We first investigated the flow characteristics of each of the protocols present in the trace. The protocol is defined as either the IP protocol number of the packet, or if that indicates TCP or UDP, the well-known port number from either the source port or the destination port numbers. For each packet in the trace we check to see if a suitable flow is available. If so, the statistics of the flow are updated, else a new flow is created. For this purpose a flow is defined as suitable if it was established between the same source and destination IP addresses and for the same protocol as the packet being processed. A flow is deleted after it has remained idle for 60 seconds although the duration of the flow is recorded as being the time from when it was created until the time of the last packet transmitted on the flow. The traffic flow analysis presented in [25] suggests that a timeout value of the order of 60 seconds is a reasonable compromise between the size of the flow table and the probability of deleting flows that will shortly become active again.

The results are presented in table 1 for all protocols with a recognizable protocol or port number that contributed more than 0.05% of the total packets in the trace. (The table accounts for about 82% of the total number of packets.) For each protocol the table gives the

^{**} We are grateful to K. Claffy and Hans-Werner Braun, Applied Network Research, San Diego Supercomputer Center, for making the trace available to us. The trace is available at <ftp://www.nlanr.net/Traces>.

Protocol	port		%flows	%pkts	%bytes	flows/s	pkts/s	duration	pkts/flow	bytes/pkt
IP in IP		✓	0.04	2.73	2.57	0.09	456	173.1	2307	253
TCP ftp-data	20	✓	0.76	12.09	15.18	2.17	2018	118.2	525	338
TCP ftp-cntrl	21		1.55	0.74	0.23	6.50	124	38.6	16	83
TCP telnet	23	✓	1.39	4.81	1.61	4.24	803	114.3	114	90
TCP smtp	25		10.26	4.80	2.82	49.49	802	18.2	15	158
UDP dns	53		45.30	5.57	3.04	216.56	929	15.4	4	147
TCP gopher	70	✓	0.45	0.54	0.55	1.87	91	43.3	40	275
TCP http	80	✓	17.94	40.21	41.53	72.98	6717	56.5	74	278
TCP pop-v3	110		0.08	0.05	0.03	0.41	9	27.0	21	148
TCP authent	113		2.12	0.19	0.05	10.54	32	9.0	3	64
TCP nntp	119	✓	0.35	6.56	6.59	0.68	1096	176.7	627	270
UDP ntp	123		5.01	0.20	0.06	25.02	33	1.37	1.3	83
TCP netbios	139	✓	0.03	0.08	0.15	0.11	14	69.8	82	501
UDP snmp	161		1.35	0.26	0.11	6.14	43	17.9	6	115
TCP login	513	✓	0.09	0.24	0.14	0.31	41	88.1	92	156
TCP cmd	514	✓	0.01	0.13	0.07	0.06	21	49.1	316	149
TCP audio	1397	✓	0.00	2.20	2.62	0.01	367	167.9	15653	321
TCP AOL	5190	✓	0.18	0.46	0.38	0.51	77	129.8	84	223
TCP X-11		✓	0.08	0.66	0.53	0.18	111	160.6	276	217

Table 1: Flow Statistics per Protocol

percentage of the total number of flows, packets, and bytes contributed by that protocol. It gives the mean number of flow setups/s after the initial startup phase, the mean number of packets/s, the average duration of each flow, the average number of packets transmitted across each flow, and the mean number of bytes per packet. Protocols with characteristics for which it appears worth establishing a flow are marked “✓”. These are protocols with an average flow duration in excess of about 20 seconds and which transmit an average of more than about 40 packets per flow (an arbitrary limit). If we assume that these flow characteristics are a property of the application behind the protocol, and the manner in which people use the applications, then for any individual protocol the results should remain relatively independent of the position in the network that the measurement is made. The characteristics of each protocol should also change only slowly over time. Thus we can use this information in deciding whether to establish a flow for any given packet.

It is interesting to note that http (web traffic) shows an average of 74 packets per flow, much higher than the value typically quoted (about 15–20 packets per flow). This is because we are looking at host pair flows, which allow multiple TCP connections between the same two IP addresses to share a single flow, rather than assuming a separate flow for each TCP connection.

In a second experiment each packet is first classified according to its protocol. If it belongs to a protocol marked “✓” in table 1 it is suitable for switching. If neither the source nor destination port numbers are well

known (less than 1024 or a recognized registered number) we assume the packet is suitable for switching if belongs to TCP but not if it is UDP. (This is the best guess we can make for packets that do not have a recognizable port number.) For those packets classified for switching we check to see if a suitable flow exists. If the search fails a flow is created. In this experiment a flow is suitable if it has the same source and destination IP address as the packet being processed. Flows are deleted after a timeout of 60 seconds. In this experiment 84% of the packets and 91% of the bytes in the trace are recognized as suitable for switching. A mean of 92 flows per second are established after an initial startup phase of 60 seconds, with a 95th percentile of 116 flows per second. The average number of established flows in the flow table is 15,500.

The experiment was repeated with all packets classified for switching. Thus, a suitable flow must exist, or be established, for every packet. This simulates the purely connection-oriented approaches to IP over ATM. In this case a mean of 422 flows/s must be established with an average of 42,000 entries in the flow table. This clearly demonstrates the advantage of connectionless forwarding for short lived flows.

The trace contains an average of 16,700 incoming packets per second of which, in the second experiment, about 14,100 are recognized as suitable for switching and the remaining 2,600 must be forwarded by the processor. We may assume that in an efficient implementation the work required to establish a flow is approximately equivalent to the number of packets that must be sent to

set it up and tear it down. This we may assume to be about 9 packets plus the original packet that must be forwarded before the flow is established. Thus, the amount of work required to establish 92 flows is approximately equivalent to the work required to forward 920 packets. So, if flows are established for all of the packets marked as suitable for switching, 16,700 packets per second may be handled with an amount of work equivalent to that required to forward 3,520 packets per second. By adding the ATM switch we are able to handle approximately 4.5 times more traffic. While one cannot dwell too heavily on the results from a single traffic trace this is a very encouraging result for Internet backbone traffic.

7. Conclusion

Current proposals for supporting IP over ATM networks hide the real network topology from the IP layer by treating the data link layer as a large opaque cloud. This leads to complexity, inefficiency and duplication of functionality. The cloud approach is untested and its scaling properties unproven, yet the Internet is proof that IP can scale to very large networks without requiring the concept of a cloud at the data link layer.

We have proposed a connectionless approach to integrate IP with fast ATM switching hardware. The IP routing decision is cached as soft state in the ATM switch such that future packets belonging to the same flow may be switched in hardware rather than forwarded by software. We believe that this approach combines the simplicity and robustness of IP with the speed and capacity of ATM.

Simulation results using a traffic trace from the Internet backbone indicate that for this trace 84% of the packets and 91% of the bytes were recognized as suitable for switching. Assuming an efficient implementation this suggests that for this trace the addition of an ATM switch could increase the traffic capacity of the routing software by up to 4.5 times. For campus and corporate backbone traffic the increase in capacity is likely to be much higher.

Flow labelled IP is the marriage of the fast but dumb to the not-so-fast but flexible in the firm conviction that the result will be both fast and flexible.

Acknowledgement

We would like to thank Bob Hinden, Fong-Ching Liaw, and Eric Hoffman, for their contribution to the architecture presented here, for many hours of patient discussion, and many more lines of code.

References

- [1] D. D. Clark, S. Shenker, and L. Zhang, "Supporting real-time applications in an integrated services packet network," Proc. ACM SIGCOMM, Comp. Commun. Review 22(4), Sep. 1992, 14–26.
- [2] R. Braden, D. Clark, and S. Shenker, "Integrated services in the Internet architecture: An overview," IETF RFC 1633, Jly. 1994.
- [3] A. S. Thyagarajan, S. L. Casner, and S. E. Deering, "Making the Mbone real," Proc. INET, Honolulu, Jun. 1995, 465–473
- [4] A. Alles, "ATM Internetworking," May 1995, <<http://www.cisco.com/warp/public/614/12.html>>.
- [5] J-Y. LeBoudec, E. Port, and H. L. Truong, "Flight of the FALCON," IEEE Commun. Mag., Feb. 1993, 50–56.
- [6] P. Newman, "ATM local area networks," IEEE Commun. Mag., Mar. 1994, 86–98.
- [7] N. Kavak, "Data communication in ATM networks," IEEE Network Mag., May 1995, 28–37.
- [8] M. Laubach, "Classical IP and ARP over ATM," IETF RFC 1577, Jan. 1994.
- [9] R. G. Cole, D. H. Shur, and C. Villamizar, "IP over ATM: A framework document," IETF Internet Draft, draft-ietf-ipatm-framework-doc-03.ps, Jun. 1995.
- [10] G. Armitage, "Support for multicast over UNI 3.1 based ATM networks," IETF Internet Draft, draft-ietf-ipatm-ipmc-05.txt, May 1995.
- [11] J. Heinanen and R. Govindan, "NBMA address resolution protocol (NARP)," IETF RFC 1735, Dec. 1994.
- [12] D. Katz and D. Piscitello, "NBMA next hop resolution protocol (NHRP)," IETF Internet Draft, draft-ietf-rolc-nhrp-04.txt, May 1995.
- [13] C. Brown, "Baseline text for MPOA," ATM Forum/95-0824r4, Nov. 1995.
- [14] S. Deering, R. Hinden, "Internet protocol, version 6 (IPv6)," IETF Internet Draft, draft-ietf-ipngwg-ipv6-spec-02.txt, Jun. 1995.
- [15] Y. Rekhter and D. Kandlur, "IP architecture extensions for ATM," IETF Internet Draft, draft-rekhter-ip-atm-architecture.txt, Jan. 1995.
- [16] B. Braden, J. Postel, and Y. Rekhter, "Internet extensions for shared media," IETF RFC 1620, May. 1994.
- [17] D. D. Clark, "The design philosophy of the DARPA Internet protocols," Proc. ACM SIGCOMM, Comp. Commun. Review 18(4), Aug. 1988, 106–114.
- [18] P. Newman, "ATM technology for corporate networks," IEEE Commun. Mag., Apr. 1992, 90–101.
- [19] "The IFMP protocol specification for IPv4," Ipsilon Networks, <<http://www.ipsilon.com>>.
- [20] "The transmission of flow labelled IPv4 on ATM data links," Ipsilon Networks, <<http://www.ipsilon.com>>.
- [21] S. Floyd and V. Jacobson, "Link-sharing and resource management models for packet networks," IEEE/ACM Trans. Networking, 3(4) Aug. 1995, 365–386.
- [22] I. Wakeman et al., "Implementing real time packet forwarding policies using streams," Usenix 1995 Technical Conference, New Orleans, Jan. 1995, 71–82.

- [23] R. Braden et al., “Resource ReSerVation Protocol — Version 1 functional specification,” IETF Internet Draft, draft-ietf-rsvp-spec-05.ps, Mar. 1995.
- [24] L. Zhang et al., “RSVP: A new resource ReSerVation Protocol,” IEEE Network Mag., Sep. 1993, 8–18.
- [25] K. C. Claffy, H.-W. Braun, and G. C. Polyzos “A parameterizable methodology for Internet traffic flow profiling,” IEEE J. Selected Areas in Commun. 13(8), Oct. 1995, 1481–1494.