



ELSEVIER

Journal of Systems Architecture 46 (2000) 483–511

JOURNAL OF
SYSTEMS
ARCHITECTURE

www.elsevier.com/locate/sysarc

On the design of IP routers Part 1: Router architectures

James Aweya *

Systems Design Engineer, Nortel Networks, P.O. Box 3511, Station C, Ottawa, Canada K1Y 4H7

Received 12 February 1999; received in revised form 27 May 1999; accepted 19 July 1999

Abstract

Internet Protocol (IP) networks are currently undergoing transitions that mandate greater bandwidths and the need to prepare the network infrastructures for converged traffic (voice, video, and data). Thus, in the emerging environment of high performance IP networks, it is expected that local and campus area backbones, enterprise networks, and Internet Service Providers (ISPs) will use multigigabit and terabit networking technologies where IP routers will be used not only to interconnect backbone segments but also to act as points of attachments to high performance wide area links. Special attention must be given to new powerful architectures for routers in order to play that demanding role. In this paper, we describe the evolution of IP router architectures and highlight some of the performance issues affecting IP routers. We identify important trends in router design and outline some design issues facing the next generation of routers. It is also observed that the achievement of high throughput IP routers is possible if the critical tasks are identified and special purpose modules are properly tailored to perform them. © 2000 Elsevier Science B.V. All rights reserved.

Keywords: Internet Protocol networks; Routers; Performance; Design

1. Introduction

The rapid growth in the popularity of the Internet has caused the traffic on the Internet to grow drastically every year for the last several years. It has also spurred the emergence of many Internet Service Providers (ISPs). To sustain growth, ISPs need to provide new differentiated services, e.g., tiered service, support for multimedia applications, etc. The routers in the ISPs' networks play a critical role in providing these

services. Internet Protocol (IP) traffic on private enterprise networks has also been growing rapidly for some time. These networks face significant bandwidth challenges as new application types, especially desktop applications uniting voice, video, and data traffic need to be delivered on the network infrastructure. This growth in IP traffic is beginning to stress the traditional processor-based design of current-day routers and as a result has created new challenges for router design.

Routers have traditionally been implemented purely in software. Because of the software implementation, the performance of a router was limited by the performance of the processor executing the protocol code. To achieve wire-speed routing, high-performance processors together

* Corresponding author. Tel.: +1-613-763-6491; fax: +1-613-763-5692.

E-mail address: aweyaj@nortelnetworks.com (J. Aweya).

with large memories were required. This translated into higher cost. Thus, while software-based wire-speed routing was possible at low-speeds, for example, with 10 megabits per second (Mbps) ports, or with a relatively smaller number of 100 Mbps ports, the processing costs and architectural implications make it difficult to achieve wire-speed routing at higher speeds using software-based processing.

Fortunately, many changes in technology (both networking and silicon) have changed the landscape for implementing high-speed routers. Silicon capability has improved to the point where highly complex systems can be built on a single integrated circuit (IC). The use of 0.35 μm and smaller silicon geometries enables application specific integrated circuit (ASIC) implementations of millions of gate-equivalents. Embedded memory (SRAM, DRAM) and microprocessors are available in addition to high-density logic. This makes it possible to build single-chip, low-cost routing solutions that incorporate both hardware and software as needed for best overall performance.

In this paper we investigate the evolution of IP router designs and highlight the major performance issues affecting IP routers. The need to build fast IP routers is being addressed in a variety of ways. We discuss these in various sections of the paper. In particular, we examine the architectural constraints imposed by the various router design alternatives. The scope of the discussion presented here does not cover more recent *label switching* routing techniques such as IP Switching [1], the Cell Switching Router (CSR) architecture [2], Tag Switching [3], and Multiprotocol Label Switching (MPLS), which is a standardization effort underway at the Internet Engineering Task Force (IETF). The discussion is limited to routing techniques as described in RFC 1812 [4].

In Section 2, we briefly review the basic functionalities in IP routers. The IETF's *Requirements for IP Version 4 Routers* [4] describes in great detail the set of protocol standards that Internet Protocol version 4 (IPv4) routers need to conform to. Section 3 presents the design issues and trends that arise in IP routers. In this section, we present an overview of IP router designs and point out their major innovations and limitations. The intercon-

nection unit (or switch fabric) is one of the main components in a router. The switch fabric can be implemented using many different techniques which are described in detail in a number of publications (e.g., [5–7]). Section 4 presents an overview of these switch fabric technologies. The concluding remarks of the paper are given in Section 5.

2. Basic functionalities in IP routers

Generally, routers consist of the following basic components: several network interfaces to the attached networks, processing module(s), buffering module(s), and an internal interconnection unit (or switch fabric). Typically, packets are received at an inbound network interface, processed by the processing module and, possibly, stored in the buffering module. Then, they are forwarded through the internal interconnection unit to the outbound interface that transmits them on the next hop on the journey to their final destination. The aggregate packet rate of all attached network interfaces needs to be processed, buffered and relayed. Therefore, the processing and memory modules may be replicated either fully or partially on the network interfaces to allow for concurrent operations.

A generic architecture of an IP router is given in Fig. 1. Fig. 1(a) shows the basic architecture of a typical router: the controller card (which holds the CPU), the router backplane, and interface cards. The CPU in the router typically performs such functions as path computations, routing table maintenance, and reachability propagation. It runs which ever routing protocol is needed in the router. The interface cards consist of adapters that perform inbound and outbound packet forwarding (and may even cache routing table entries or have extensive packet processing capabilities). The router backplane is responsible for transferring packets between the cards. The basic functionalities in an IP router can be categorized as: route processing, packet forwarding, and router special services. The two key functionalities for packet routing are route processing (i.e., path computation, routing table maintenance, and reachability

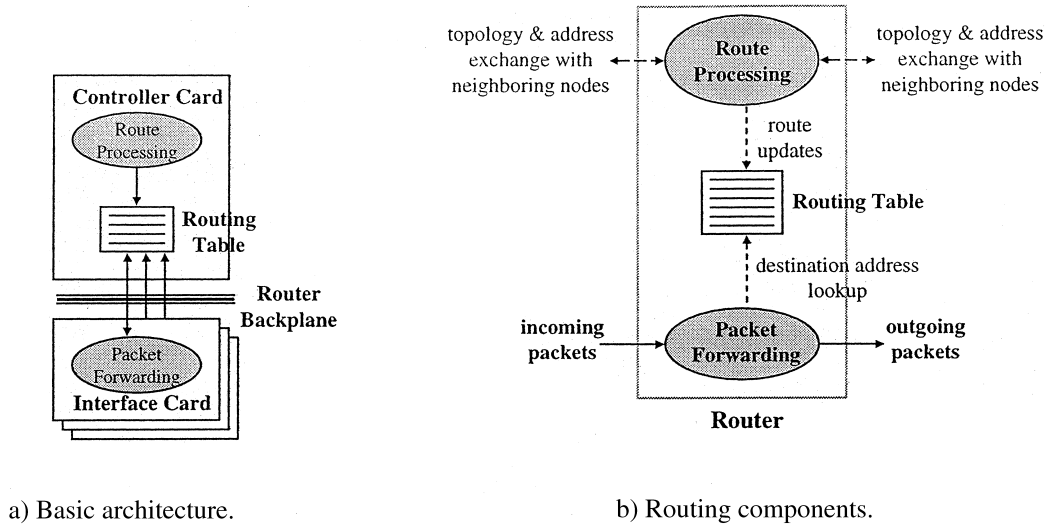


Fig. 1. Generic architecture of a router.

propagation) and packet forwarding (see Fig. 1(b)). We discuss the three functionalities in more detail below.

2.1. Route processing

Routing protocols are the means by which routers gain information about the network. Routing protocols map network topology and store their view of that topology in the routing table. Thus, route processing includes routing table construction and maintenance using routing protocols such as the Routing Information Protocol (RIP) and Open Shortest Path First (OSPF) [8–10]. The routing table consists of routing entries that specify the destination and the next-hop router through which the datagram should be forwarded to reach the destination. Route calculation consists of determining a route to the destination: network, subnet, network prefix, or host.

In static routing, the routing table entries are created by default when an interface is configured (for directly connected interfaces), added by, for example, the route command (normally from a system bootstrap file), or created by an Internet Control Message Protocol (ICMP) redirect (usually when the wrong default is used) [11]. Once configured, the network paths will not change.

With static routing, a router may issue an alarm when it recognizes that a link has gone down, but will not automatically reconfigure the routing table to reroute the traffic around the disabled link. Static routing, used in LANs over limited distances, requires basically the network manager to configure the routing table. Thus, static routing is fine if the network is small, there is a single connection point to other networks, and there are no redundant routes (where a backup route can be used if a primary route fails). Dynamic routing is normally used if any of these three conditions do not hold true.

Dynamic routing, used in internetworking across wide area networks, automatically reconfigures the routing table and recalculates the least expensive path. In this case, routers broadcast advertisement packets (signifying their presence) to all network nodes and communicate with other routers about their network connections, the cost of connections, and their load levels. Convergence, or reconfiguration of the routing tables, must occur quickly before routers with incorrect information misroute data packets into dead ends. Some dynamic routers can also rebalance the traffic load.

The use of dynamic routing does not change the way an IP forwarding engine performs routing at

the IP layer. What changes is the information placed in the routing table – instead of coming from the route commands in bootstrap files, the routes are added and deleted dynamically by a routing protocol, as routes change over time. The routing protocol adds a routing policy to the system, choosing which routes to place in the routing table. If the protocol finds multiple routes to a destination, the protocol chooses which route is the best, and which one to insert in the table. If the protocol finds that a link has gone down, it can delete the affected routes or add alternate routes that bypass the problem.

A network (including several networks administered as a whole) can be defined as an autonomous system. A network owned by a corporation, an ISP, or a university campus often defines an autonomous system. There are two principal routing protocol types [8–11]: those that operate within an autonomous system, or the Interior Gateway Protocols (IGPs), and those that operate between autonomous systems, or Exterior Gateway Protocols (EGPs). Within an autonomous system, any protocol may be used for route discovery, propagating, and validating routes. Each autonomous system can be independently administered and must make routing information available to other autonomous systems. The major IGPs include RIP, OSPF, and IS-IS (Intermediate System to Intermediate System). Some EGPs include EGP and Border Gateway Protocol (BGP). Refs. [8–10] present detail discussions on routing protocols.

2.2. Packet forwarding

2.2.1. Forwarding process

In this section, we briefly review the forwarding process in IPv4 routers. More details of the forwarding requirements are given in Ref. [4]. A router receives an IP packet on one of its interfaces and then forwards the packet out of another of its interfaces (or possibly more than one, if the packet is a multicast packet), based on the contents of the IP header. As the packet is forwarded hop by hop, the packet's (original) network layer header (IP header) remains relatively unchanged, containing the complete set of instructions on how

to forward the packet (IP tunneling may call for prepending the packet with other IP headers in the network). However, the data-link headers and physical-transmission schemes may change radically at each hop in order to match the changing media types.

Suppose that the router receives a packet from one of its attached network segments. The router verifies the contents of the IP header by checking the protocol version, header length, packet length, and header checksum fields. The protocol version must be equal to 4 for IPv4 which we assume in this paper. The header length must be greater than or equal to the minimum IP header size (20 bytes). The length of the IP packet, expressed in bytes, must also be larger than the minimum header size. In addition, the router checks that the entire packet has been received, by checking the IP packet length against the size of the received Ethernet packet, for example, in the case where the interface is attached to an Ethernet network. To verify that none of the fields of the header have been corrupted, the 16-bit ones-complement checksum of the entire IP header is calculated and verified to be equal to 0xffff. If any of these basic checks fail, the packet is deemed to be malformed and is discarded without sending an error indication back to the packet's originator.

Next, the router verifies that the time-to-live (TTL) field is greater than 1. The purpose of the TTL field is to make sure that packets do not circulate forever when there are routing loops. The host sets the packet's TTL field to be greater than or equal to the maximum number of router hops expected on the way to the destination. Each router decrements the TTL field by 1 when forwarding; when the TTL field is decremented to 0, the packet is discarded, and an Internet Control Message Protocol (ICMP) TTL Exceeded message is sent back to the host. On decrementing the TTL, the router must update the packet's header checksum. RFC 1141 [92] contains implementation techniques for computing the IP checksum. Since a router often changes only the TTL field (decrementing it by 1), a router can incrementally update the checksum when it forwards a received packet, instead of calculating the checksum over the entire IP header again.

The router then looks at the destination IP address. The address indicates a single destination host (unicast), a group of destination hosts (multicast), or all hosts on a given network segment (broadcast). Unicast packets are discarded if they were received as data-link broadcasts or as multicasts; otherwise, multiple routers may attempt to forward the packet, possibly contributing to a broadcast storm. In packet forwarding, the destination IP address is used as a key for the routing table lookup. The best-matching routing table entry is returned, indicating whether to forward the packet and, if so, the interface to forward the packet out of and the IP address of the next IP router (if any) in the packet's path. The next-hop IP address is used at the output interface to determine the link address of the packet, in case the link is shared by multiple parties (such as an Ethernet, Token Ring or Fiber Distributed Data Interface (FDDI) network), and is consequently not needed if the output connects to a point-to-point link.

In addition to making forwarding decisions, the forwarding process is responsible for making packet classifications for quality of service (QoS) control and access filtering. Flows can be identified based on source IP address, destination IP address, TCP/UDP port numbers as well as IP type of service (TOS) field. Classification can even be based on higher layer packet attributes.

If the packet is too large to be sent out of the outgoing interface in one piece (that is, the packet length is greater than the outgoing interface's Maximum Transmission Unit (MTU)), the router attempts to split the packet into smaller fragments. Fragmentation, however, can affect performance adversely [12]. Host may instead wish to prevent fragmentation by setting the Don't Fragment (DF) bit in the Fragmentation field. In this case, the router does not fragment but instead drops the packet and sends an ICMP Destination Unreachable (subtype Fragmentation Needed and DF Set) message back to the host. The host uses this message to calculate the minimum MTU along the packet's path [13], which is in turn used to size future packets.

The router then prepends the appropriate data-link header for the outgoing interface. The IP

address of the next hop is converted to a data-link address, usually using the Address Resolution Protocol (ARP) [14] or a variant of ARP, such as Inverse ARP [15] for Frame Relay subnets. The router then sends the packet to the next hop, where the process is repeated.

An application can also modify the handling of its packets by extending the IP headers of its packets with one or more IP options. IP options are used infrequently for regular data packets, because most internet routers are heavily optimized for forwarding packets having no options. Most IP options (such as the record-route and timestamp options) are used to aid in statistics collection but do not affect a packet's path. However, the strict-source route and the loose-source route options can be used by an application to control the path its packets take. The strict-source route option is used to specify the exact path that the packet will take, router by router. The utility of strict-source route is limited by the maximum size of the IP header (60 bytes), which limits to 9 the number of hops specified by the strict-source route option. The loose-source route is used to specify a set of intermediate routers (again, up to 9) through which the packet must go on the way to its destination. Loose-source routing is used mainly for diagnostic purposes, such as an aid to debugging internet routing problems.

2.2.2. Route lookup

For a long time, the major performance bottleneck in IP routers has been the time it takes to look up a route in the routing table. The problem is defined as that of searching through a database (routing table) of destination prefixes and locating the longest prefix that matches the destination IP address of a given packet. A routing table lookup returns the best-matching routing table entry, which tells the router which interface to forward the packet out of and the IP address of the packet's next hop. Prefix matching was introduced in the early 1990s, when it was foreseen that the number of end users as well as the amount of routing information on the Internet would grow enormously. The address classes A, B, and C (allowing sites to have 24, 16, and 8 bits respectively for addressing) proved too inflexible and wasteful of the

address space. To make better use of this scarce resource, especially the class B addresses, bundles of class C networks were given out instead of class B addresses. This resulted in massive growth of routing table entries. Thus, Classless Inter-Domain Routing (CIDR) [16] introduced a fundamental concept in route aggregation which allows for contiguous blocks of address space to be aggregated and advertised as singular routes instead of individual entries in a routing table. This helps keep the size of the routing tables smaller than if each address were advertised individually.

Many routers have a *default route* in their routing table. This is a routing table entry for the zero-length prefix 0/0 (or 0.0.0.0). The default route matches every destination, although it is overridden by all more specific prefixes. For example, suppose that an organization's intranet has one router attaching itself to the public Internet. All of the organization's other routers can then use a default route pointing to the Internet connection rather than knowing about all of the public Internet's routes.

The first approaches for longest prefix matching used radix trees or modified Patricia trees [17,18] combined with hash tables, (Patricia stands for "Practical Algorithm to Retrieve Information Coded in Alphanumeric" [19]). These trees are binary trees, whereby the tree traversal depends on a sequence of single-bit comparisons in the key, the destination IP address. These lookup algorithms have complexity proportional to the number of address bits which, for IPv4 is only 32. In the worst case it takes time proportional to the length of the destination address to find the longest prefix match. Worse, the commonly used Patricia algorithm may need to backtrack to find the longest match, leading to poor worst-case performance. The performance of Patricia is somewhat data dependent. With a particularly unfortunate collection of prefixes in the routing table, the lookup of certain addresses can take as many as 32 bit comparisons, one for each bit of the destination IP address.

Early implementations of routers, however, could not afford such expensive computations. Thus, one way to speed up the routing table lookup is to try to avoid it entirely. The routing

table lookup provides the next hop for a given IP destination. Some routers cache this *IP destination-to-next-hop association* in a separate database that is consulted (as the front end to the routing table) before the routing table lookup as shown in Fig. 2. Finding a particular destination in this database of IP destination-to-next-hop association is easier because an exact match is done instead of the more expensive best-match operation of the routing table. So, most routers relied on route caches [20,21]. The route caching techniques rely on there being enough locality in the traffic so that the cache hit rate is sufficiently high and the cost of a routing lookup is amortized over several packets.

This front-end database might be organized as a *hash table* [22]. After the router has forwarded several packets, if the router sees any of these destinations again (a *cache hit*), their lookups will be very quick. Packets to new destinations will be slower, however, because the cost of a failed hash lookup will have to be added to the normal routing table lookup. Front-end caches to the routing table can work well at the edge of the Internet or within organizations. However, cache schemes do not seem to work well in the Internet's core. The large number of packet destinations seen by the core routers can cause caches to overflow or for their lookup to become slower than the routing table lookup itself. Cache schemes are not really effective when the hash bucket size (the number of destinations that hash to the same value) starts getting large. Also, the frequent routing changes seen in the core routers can force them to invali-

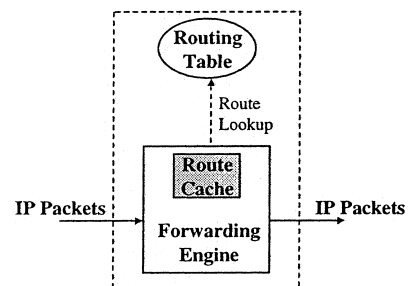


Fig. 2. A route cache used as a front-end database to the routing table.

date their caches frequently, leading to a small number of cache hits.

Typically, two types of packets arrive at a router: packets to be forwarded to another network or packets destined to the router itself. Whether a packet to a router causes a routing table reference depends on the router implementation. Some implementations may speed up routing table lookups. One possibility is for the router to explicitly check each incoming packet against a table of all of the router's addresses to see if there is a match. This explicit check means that the routing table is never consulted about packets destined to the router. Another possibility is to use the routing table for all packets. Before the packet is sent, the router checks if the packet is to its own address on the appropriate network. If the packet is for the router, then it is never transmitted. The explicit check after the routing table lookup requires checking a smaller number of router addresses at the increased cost of a routing table lookup. New routing table lookup algorithms are still being developed in attempts to build even faster routers. Recent examples are found in Refs. [23–28].

The basic idea of one of the recent algorithms [23] is to create a small and compressed data structure that represents large lookup tables using a small amount of memory. The technique exploits the sparseness of actual entries in the space of all possible routing table entries. This results in a lower number of memory accesses (to a fast memory) and hence faster lookups. The proposal reduces the routing table to very efficient representation of a binary tree such that the majority of the table can reside in the primary cache of the processor, allowing route lookups at gigabit speeds. In addition, the algorithm does not have to calculate expensive perfect hash functions, although updates to the routing table are still not easy. Ref. [24] proposes another approach for implementing the compression and minimizing the complexity of updates.

The recent work of Waldvogel et al. [25] presents an alternative approach which reduces the number of memory references rather than compact the routing table. The main idea is to first create a perfect hash table of prefixes for each prefix length.

A binary search among all prefix lengths, using the hash tables for searches amongst prefixes of a particular length, can find the longest prefix match for an N bit address in $O(\log N)$ steps. Although the algorithm has very fast execution times, calculating perfect hashes can be slow and can slow down updates.

Hardware based techniques for route lookup are also actively being investigated both in research and commercial designs (e.g., [29–32]). Other designs of forwarding engine have concentrated on IP packet header processing in hardware, to remove the dependence upon caching, and to avoid the cost of high-speed processor. Designs based upon content-addressable memory have been investigated [29], but such memory is too far expensive to be applied to a large routing table. Hardware route lookup and forwarding is also under active investigation in both research and commercial designs [30,31]. The argument for a software-based implementation stresses flexibility. Hardware implementations can generally achieve a higher performance at lower cost but are less flexible.

2.3. Router special services

Besides dynamically finding the paths for packets to take towards their destinations, routers also implement other functions. Anything beyond core routing functions falls into this category, for example, authentication and access services such as packet filtering for security/firewall purposes. Companies often put a router between their company network and the Internet and then configure the router to prevent unauthorized access to the company's resources from the Internet. This configuration may consist of certain patterns (for example, source and destination address and TCP port) whose matching packets should not be forwarded or of more complex rules to deal with protocols that vary their port numbers over time, such as the File Transfer Protocol (FTP). Such routers are called firewalls. Similarly, internet service providers (ISPs) often configure their routers to verify the source address in all packets received from the ISP's customers. This foils certain security attacks and makes other attacks easier to

trace back to their source. Similarly, ISPs providing dial-in access to their routers typically use Remote Authentication Dial-In User Service (RADIUS) [33] to verify the identity of the person dialing in.

Often other functions, less directly related to packet forwarding, also get incorporated into IP routers. Examples of these non-forwarding functions include network management components such as Simple Network Management Protocol (SNMP) and Management Information Bases (MIBs). Routers also play an important role in TCP/IP congestion control algorithms. When an IP network is congested, routers cannot forward all the packets they receive. By simply discarding some of their received packets, routers provide feedback to TCP congestion control algorithms, such as the TCP slow-start algorithm [34,35]. Early Internet routers simply discarded excess packets instead of queueing them onto already full transmit queues; these routers are termed drop-tail gateways. However, this discard behavior was found to be unfair, favoring applications that send larger and more bursty data streams. Modern Internet routers employ more sophisticated, and fairer, drop algorithms, such as Random Early Detection (RED) [36].

Algorithms also have been developed that allow routers to organize their transmit queues so as to give resource guarantees to certain classes of traffic or to specific applications. These queueing or link scheduling algorithms include Weighted Fair Queueing (WFQ) [37] and Class Based Queueing (CBQ) [38]. A protocol called RSVP [39] has been developed that allows hosts to dynamically signal to routers which applications should get special queueing treatment. However, RSVP has not yet been deployed, with some people arguing that queueing preference could more simply be indicated by using the Type of Service (TOS) bits in the IP header [40,41].

Some vendors allow the collection of traffic statistics on their routers, for example, how many packets and bytes are forwarded per receiving and transmitting interface on the router. These statistics are used for future capacity planning. They can also be used by ISPs to implement usage-based charging schemes for their customers.

3. IP router architectures

In this section, we review the main architectures proposed for IP routers. Specifically, we examine important trends in IP router design and outline some design issues facing the next generation of routers. The aim is not to provide an exhaustive review of existing architectures, but instead to give the reader a perspective on the range of options available and the associated trade-off between performance, functionality, and complexity.

3.1. Bus-based router architectures with single processor

The first generation of IP router was based on software implementations on a single general-purpose central processing unit (CPU). These routers consist of a general-purpose processor and multiple interface cards interconnected through a shared bus as depicted in Fig. 3.

Packets arriving at the interfaces are forwarded to the CPU which determines the next hop address and sends them back to the appropriate outgoing interface(s). Data are usually buffered in a centralized data memory [42,43], which leads to the disadvantage of having the data cross the bus twice, making it the major system bottleneck. Packet processing and node management software (including routing protocol operation, routing table maintenance, routing table lookups, and other control and management protocols such as ICMP, SNMP) are also implemented on the central pro-

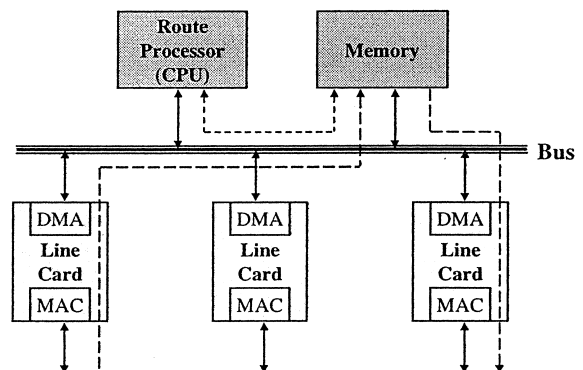


Fig. 3. Traditional bus-based router architecture.

cessor. Unfortunately, this simple architecture yields low performance for the following reasons:

- The central processor has to process all packets flowing through the router (as well as those destined to it). This represents a serious processing bottleneck.
- Some major packet processing tasks in a router involve memory intensive operations (e.g., table lookups) which limits the effectiveness of processor power upgrades in boosting the router packet processing throughput. Routing table lookups and data movements are the major consumer of overall packet processing cycles. Packet processing time does not decrease linearly if faster processors are used because of the sometimes dominating effect of memory access rate.
- Moving data from one interface to the other (either through main memory or not) is a time consuming operation that often exceeds the packet header processing time. In many cases, the computer input/output (I/O) bus quickly becomes a severe limiting factor to overall router throughput.

Since routing table lookup is a time-consuming process of packet forwarding, some traditional software-based routers cache the IP destination-to-next-hop association in a separate database that is consulted as the front end to the routing table before the routing table lookup. The justification for route caching is that packet arrivals are temporally correlated, so that if a packet belonging to a new flow arrives then more packets belonging to the same flow can be expected to arrive in the near future. Route caching of IP destination/next-hop address pairs will decrease the average processing time per packet if locality exists for packet addresses [20,21]. Still, the performance of the traditional bus-based router depends heavily on the throughput of the shared bus and on the forwarding speed of the central processor. This architecture cannot scale to meet the increasing throughput requirements of multigigabit network interface cards.

3.2. Bus-based router architectures with multiple processors

For the second generation IP routers, improvement in the shared-bus router architecture

was introduced by distributing the packet forwarding operations. In some architectures, distributing fast processors and route caches, in addition to receive and transmit buffers, over the network interface cards reduces the load on the system bus. Other second generation routers remedy this problem by employing multiple forwarding engines (dedicated solely for packet forwarding operation) in parallel since a single CPU cannot keep up with requests from high-speed input ports. An advantage of having multiple forwarding engines serving as one pool is the ease of balancing loads from the ports when they have different speeds and utilization levels. We review, in this section, these second generation router architectures.

3.2.1. Architectures with route caching

This architecture reduces the number of bus copies and speeds up packet forwarding by using a route cache of frequently seen addresses in the network interface as shown in Fig. 4. Packets are therefore transmitted only once over the shared bus. Thus, this architecture allows the network interface cards to process packets locally some of the time.

In this architecture, a router keeps a central master routing table and the satellite processors in the network interfaces each keep only a modest cache of recently used routes. If a route is not in a network interface processor's cache, it would request the relevant route from the central table. The route cache entries are traffic-driven in that the first packet to a new destination is routed by the main CPU (or route processor) via the central routing table information and as part of that forwarding operation, a route cache entry for that destination is then added in the network interface. This allows subsequent packet flows to the same destination network to be switched based on an efficient route cache match. These entries are periodically aged out to keep the route cache current and can be immediately invalidated if the network topology changes. At high-speeds, the central routing table can easily become a bottleneck because the cost of retrieving a route from the central table is many times more expensive than actually

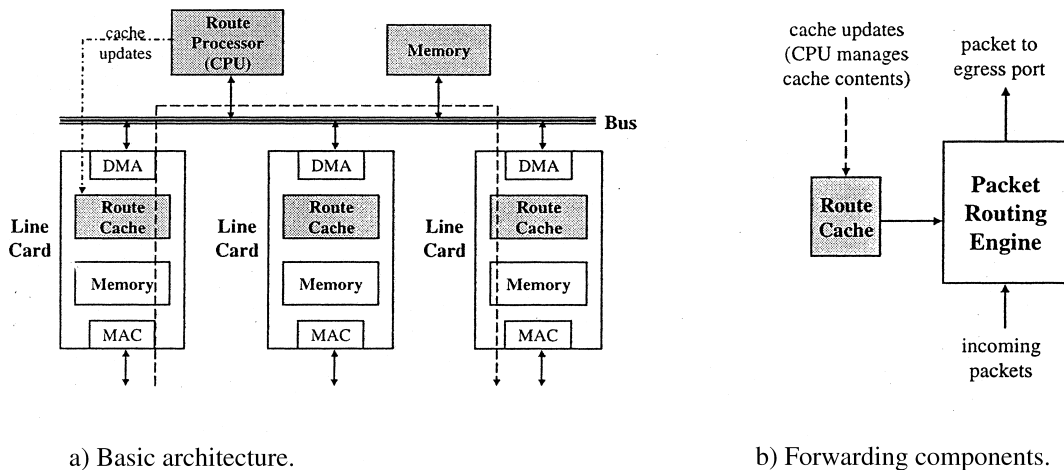


Fig. 4. Reducing the number of bus copies using a route cache in the network interface.

processing the packet local in the network interface.

A major limitation of this architecture is that it has a traffic dependent throughput and also the shared bus is still a bottleneck. The performance of this architecture can be improved by enhancing each of the distributed network interface cards with larger memories and complete forwarding tables. The decreasing cost of high bandwidth memories makes this possible. However, the shared bus and the general purpose CPU can neither scale to high capacity links nor provide traffic pattern independent throughput.

3.2.2. Architectures with multiple parallel forwarding engines

Another bus-based multiple processor router architecture is described in [44]. Multiple forwarding engines are connected in parallel to achieve high packet processing rates as shown in Fig. 5. The network interface modules transmit and receive data from the links at the required rates. As a packet comes into a network interface, the IP header is stripped by a control circuitry, augmented with an identifying tag, and sent to a forwarding engine for validation and routing. While the forwarding engine is performing the routing function, the remainder of the packet is deposited in an input buffer (in the network in-

terface) in parallel. The forwarding engine determines which outgoing link the packet should be transmitted on, and sends the updated header fields to the appropriate destination interface module along with the tag information. The packet is then moved from the buffer in the source interface module to a buffer in the destination interface module and eventually transmitted on the outgoing link.

The forwarding engines can each work on different headers in parallel. The circuitry in the interface modules peels the header off of each packet and assigns the headers to the forwarding engines in a round-robin fashion. Since in some (real time) applications packet order maintenance is an issue, the output control circuitry also goes round-robin, guaranteeing that packets will then be sent out in the same order as they were received. Better load-balancing may be achieved by having a more intelligent input interface which assigns each header to the lightest loaded forwarding engine [44]. The output control circuitry would then have to select the next forwarding engine to obtain a processed header from by following the demultiplexing order followed at the input, so that order preservation of packets is ensured. The forwarding engine returns a new header (or multiple headers, if the packet is to be fragmented), along with routing information (i.e., the immediate destination of the packet). A route processor runs the routing protocols and

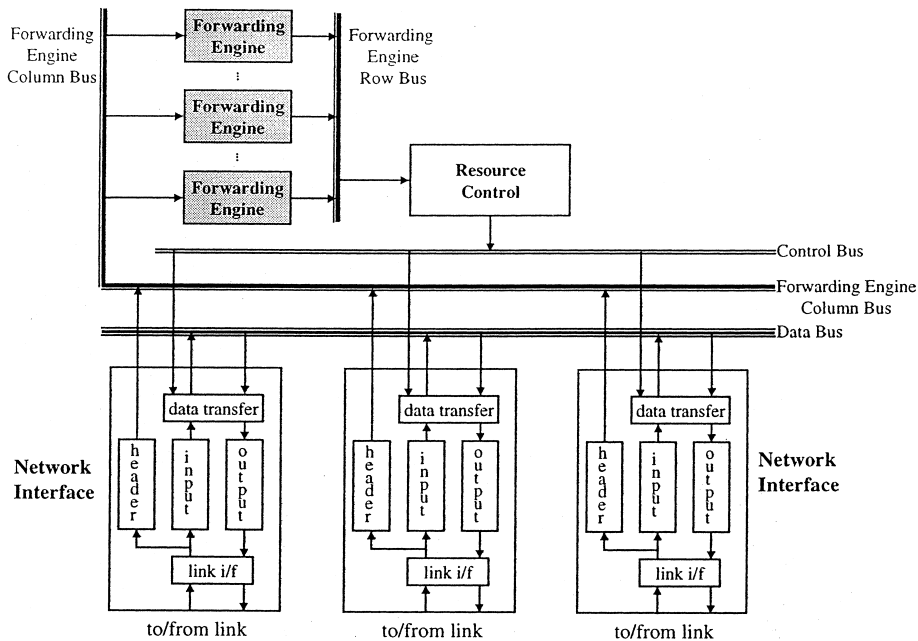


Fig. 5. Bus-based router architecture with multiple parallel forwarding engines.

creates a forwarding table that is used by the forwarding engines.

The choice of this architecture was premised on the observation that it is highly unlikely that all interfaces will be bottlenecked at the same time. Hence sharing of the forwarding engines can increase the port density of the router. The forwarding engines are only responsible for resolving next-hop addresses. Forwarding only IP headers to the forwarding engines eliminates an unnecessary packet payload transfer over the bus. Packet payloads are always directly transferred between the interface modules and they never go to either the forwarding engines or the route processor unless they are specifically destined to them.

3.3. Switch-based router architectures with multiple processors

To alleviate the bottlenecks of the second generation of IP routers, the third generation of routers were designed with the shared bus replaced by a switch fabric. This provides sufficient bandwidth for transmitting packets between interface cards

and allows throughput to be increased by several orders of magnitude. With the interconnection unit between interface cards not the bottleneck, the new bottleneck is packet processing.

The multigigabit router (MGR) is an example of this architecture [45]. The design has dedicated IP packet forwarding engines with route caches in them. The MGR consists of multiple line cards (each supporting one or more network interfaces) and forwarding engine cards, all connected to a high-speed (crossbar) switch as shown in Fig. 6. The design places forwarding engines on boards distinct from line cards. When a packet arrives at a line card, its header is removed and passed through the switch to a forwarding engine. The remainder of the packet remains on the inbound line card. The forwarding engine reads the header to determine how to forward the packet and then updates the header and sends the updated header and its forwarding instructions back to the inbound line card. The inbound line card integrates the new header with the rest of the packet and sends the entire packet to the outbound line card for transmission. The MGR, like most routers,

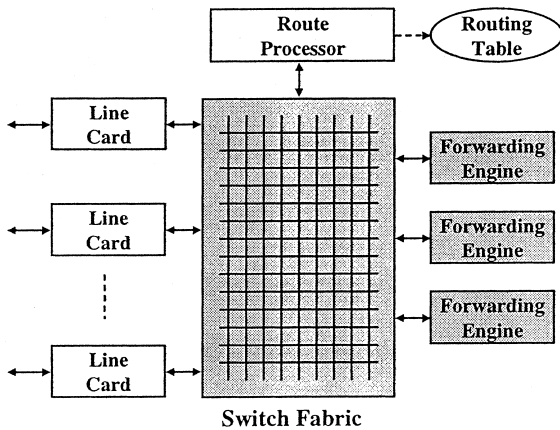


Fig. 6. Switch-based router architecture with multiple forwarding engines.

also has a control (and route) processor that provides basic management functions such as generation of routing tables for the forwarding engines and link (up/down) management. Each forwarding engine has a set of the forwarding tables (which are a summary of the routing table data).

Similar to other cache-driven architectures, the forwarding engine checks to see if the cached route matches the destination of the packet (a cache hit). If not, the forwarding engine carries out an extended lookup of the forwarding table associated with it. In this case, the engine searches the routing table for the correct route, and generates a version of the route for the route cache. Since the forwarding table contains prefix routes and the route cache is a cache of routes for particular destination, the processor has to convert the forwarding table entry into an appropriate destination-specific cache entry.

3.4. Limitation of IP packet forwarding based on route caching

Regardless of the type of interconnection unit used (bus, shared memory, crossbar, etc.), a route cache can be used in conjunction with a (centralized) processing unit for IP packet forwarding [20,21,45]. In this section, we examine the limitations of route caching techniques before we proceed with the discussion on other router architectures.

The route cache model creates the potential for cache misses which occur with “demand-caching” schemes as described above. That is, if a route is not found in the forwarding cache, the first packet(s) then looks to the routing table maintained by the CPU to determine the outbound interface and then a cache entry is added for that destination. This means when addresses are not found in the cache, the packet forwarding defaults to a classical software-based route lookup (sometimes described as a “slow-path”). Since the cache information is derived from the routing table, routing changes cause existing cache entries to be invalidated and reestablished to reflect any topology changes. In networking environments which frequently experience significant routing activity (such as in the Internet) this can cause traffic to be forwarded via the main CPU (the slow path), as opposed to via the route cache (the fast path).

In enterprise backbones or public networks, the combination of highly random traffic patterns and frequent topology changes tends to eliminate any benefits from the route cache, and performance is bounded by the speed of the software slow path which can be many orders of magnitude lower than the caching fast path.

This demand-caching scheme, maintaining a very fast access subnet of the routing topology information, is optimized for scenarios whereby the majority of the traffic flows are associated with a subnet of destinations. However, given that traffic profiles at the core of the Internet (and potentially within some large enterprise networks) do not follow closely this model, a new forwarding paradigm is required that would eliminate the increasing cache maintenance resulting from growing numbers of topologically dispersed destinations and dynamic network changes.

The performance of a product using the route cache technique is influenced by the following factors:

- how big the cache is,
- how the cache is maintained (the three most popular cache maintenance strategies are random replacement, first-in-first-out (FIFO), and least recently used (LRU)), and

- what the performance of the slow path is, since at least some percentage of the traffic will take the slow path in any application.

The main argument in favor of cache-based schemes is that a cache hit is at least less expensive than a full route lookup (so a cache is valuable provided it achieves a modest hit rate). Even with an increasing number of flows, it appears that packet bursts and temporal correlation in the packet arrivals will continue to ensure that there is a strong chance that two datagrams arriving close together will be headed for the same destination.

In current backbone routers, the number of flows that are active at a given interface can be extremely high. Studies have shown that an OC-3 interface might have an average of 256,000 flows active concurrently [46]. It is observed in [47] that, for this many flows, use of hardware caches is extremely difficult, especially if we consider the fact that a fully associative hardware cache is required. So caches of such size are most likely to be implemented as hash tables since only hash tables can be scaled to these sizes. However, the $O(1)$ lookup time of a hash table is an average case result, and the worst-case performance of a hash table can be poor since multiple headers might hash into the same location. Due to the large number of flows that are simultaneously active in a router and due to the fact that hash tables generally cannot guarantee good hashing under all arrival patterns, the performance of cache-based schemes is heavily traffic dependent. If a large number of new flows arrive at the same time, the slow path of the router can be overloaded, and it is possible that packet loss can occur due to the (slow path) processing overload and not due to output link congestion.

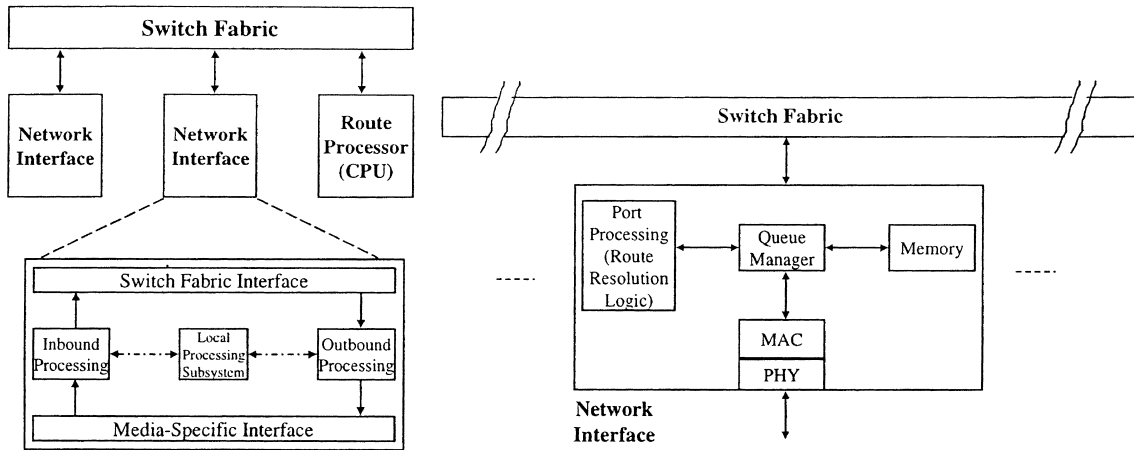
Some architectures have been proposed that avoid the potential overload of continuous cache churn (which results in a performance bottleneck). These designs use instead a forwarding database in each network interface which mirrors the entire content of the IP routing table maintained by the CPU (route processor). That is, there is a one-to-one correspondence between the forwarding database entries and routing table prefixes; therefore no need to maintain a route cache [48–50]. By eliminating the route cache, the architecture fully

eliminates the slow path. This offers significant benefits in terms of performance, scalability, network resilience and functionality, particularly in large complex networks with dynamic flows. These architectures can best accommodate the changing network dynamics and traffic characteristics resulting from increasing numbers of short flows typically associated with Web-based applications and interactive type sessions.

3.5. Switch-based router architectures with fully distributed processors

From the discussion in the preceding sections, we find that the three main bottlenecks in a router are processing power, memory bandwidth, and internal bus bandwidth. These three bottlenecks can be avoided by using a distributed switch-based architecture with properly designed network interfaces. Since routers are mostly dedicated systems not running any specific application tasks, off-loading processing to the network interfaces reflects a proper approach to increase the overall router performance. A successful step towards building high performance routers is to add some processing power to each network interface in order to reduce the processing and memory bottlenecks. General-purpose processors and/or dedicated very large scale integrated (VLSI) components can be applied. The third bottleneck (internal bus bandwidth) can be solved by using special mechanisms where the internal bus is in effect a switch (e.g., shared memory, crossbar, etc.) thus allowing simultaneous packet transfers between different pairs of network interfaces. This arrangement must also allow for efficient multicast capabilities.

We review in this section, decentralized router architectures where each network interface is equipped with appropriate processing power and buffer space. Our focus is on high performance packet processing subsystems that form an integral part of multiport routers for high-speed networks. The routers have to cope with extremely high aggregate packet rates flowing through the system and, thus, require efficient processing and memory components. A generic modular switch-based router architecture is shown in Fig. 7.



a) Functional diagram [48-50].

b) Generic architecture.

Fig. 7. A generic switch-based distributed router architecture.

Each network interface provides the processing power and the buffer space needed for packet processing tasks related to all the packets flowing through it. Functional components (inbound, outbound, and local processing elements) process the inbound, outbound traffic and time-critical port processing tasks. They perform the processing of all protocol functions (in addition to quality of service (QoS) processing functions) that lie in the critical path of data flow. In order to provide QoS guarantees, a port may need to classify packets into predefined service classes. Depending on router implementation, a port may also need to run data-link level protocols such as Serial Line Internet Protocol (SLIP), Point-to-Point Protocol (PPP) and IEEE 802.1Q VLAN (Virtual LAN), or network-level protocols such as Point-to-Point Tunneling Protocol (PPTP). The exact features of the processing components depend on the functional partitioning and implementation details. Concurrent operation among these components can be provided. The network interfaces are interconnected via a high performance switch that enables them to exchange data and control messages. In addition, a CPU is used to perform some centralized tasks. As a result, the overall processing and buffering capacity is distributed over the available interfaces and the CPU.

The Media-Specific Interface (MSI) performs all the functions of the physical layer and the Media Access Control (MAC) sublayer (in the case of the IEEE 802 protocol model). The Switch Fabric Interface (SFI) is responsible for preparing the packet for its journey across the switch fabric. The SFI may prepend an internal routing tag containing port of exit, the QoS priority, and drop priority, onto the packet.

In order to decrease latency and increase throughput, more concurrency has to be achieved among the various packet handling operations (mainly header processing and data movement). This can be achieved by using parallel or multiprocessor platforms in each high performance network interface. Obviously, there is a price for that (cost and space) which has to be considered in the router design.

To analyze the processing capabilities and to determine potential performance bottlenecks, the functions and components of a router and, especially, of all its processing subsystems have to be identified. Therefore, all protocols related to the task of a router need to be considered. In an IP router, the IP protocol itself as well as additional protocols, such as ICMP, ARP, RARP (Reverse ARP), BGP, etc., are required.

A distinction can be made between the processing tasks directly related to packets being for-

warded through the router and those related to packets destined to the router, such as maintenance, management or error protocol data. Best performance can be achieved when packets are handled by multiple heterogeneous processing elements, where each element specializes in a specific operation. In such a configuration, special purpose modules perform the time critical tasks in order to achieve high throughput and low latency. Time critical tasks are the ones related to the regular data flow. The non-time critical tasks are performed in general purpose processors (CPU). A number of commercial routers follow this design approach (e.g., [47,51–62]).

3.5.1. Critical data path processing (fast path)

The processing tasks directly related to packets being forwarded through the router can be referred to as the *time critical processing tasks*. They form the *critical path* (sometimes called the *fast path*) through a router and need to be highly optimized in order to achieve multigigabit rates. These processing tasks comprise all protocols involved in the critical path (e.g., Logical Link Control, (LLC) Subnetwork Access Protocol (SNAP) and IP) as well as ARP which can be processed in the network interface because it needs direct access to the network, even though it is not time critical. The time critical tasks mainly consist of header checking, and forwarding (and may include segmentation) functions. These protocols directly affect the performance of an IP router in terms of the number of packets that can be processed per second.

Most high-speed routers implement the fast path functions in hardware. Generally, the fast path of IP routing requires the following functions: IP packet validation, destination address parsing and table lookup, packet lifetime control (TTL update), and checksum calculation. The fast path may also be responsible for making packet classifications for QoS control and access filtering. The vast majority of packets flowing through an IP router need to have only these operations performed on them. While they are not trivial, it is possible to implement them in hardware, thereby providing performance suitable for high-speed routing.

3.5.2. Non-critical data path processing (slow path)

Packets destined to a router, such as maintenance, management or error protocol data are usually not time critical. However, they have to be integrated in an efficient way that does not interfere with the packet processing and, thus, does not slow down the time-critical path. Typical examples of these *non-time critical processing* tasks are error protocols (e.g., ICMP), routing protocols (e.g., RIP, OSPF, BGP), and network management protocols (e.g., SNMP). These processing tasks need to be centralized in a router node and typically reside above the network or transport protocols.

As shown in Fig. 8, only protocols in the forwarding path of a packet through the IP router are implemented on the network interface itself. Other protocols such as routing and network management protocols are implemented on the CPU. This way, the CPU does not adversely affect performance because it is located out of the data path, where it maintains route tables and determines the policies and resources used by the network interfaces. As an example, the CPU subsystem can be attached to the switch fabric in the same way as a regular network interface. In this case, the CPU subsystem is viewed by the switch fabric as a regular network interface. It has, however, a com-

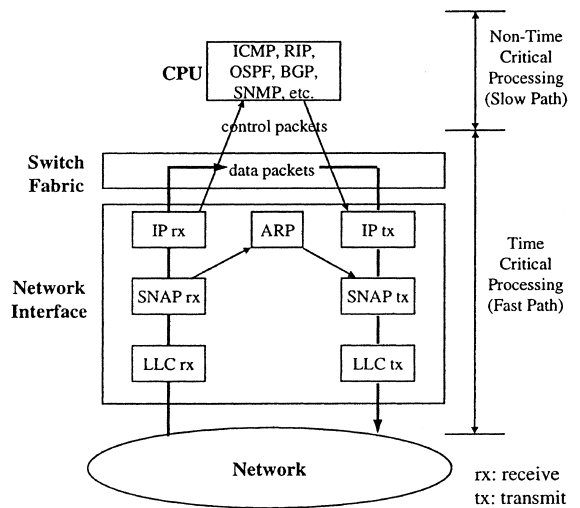


Fig. 8. Example IEEE 802 protocol entities in an IP router [Adapted from 48].

pletely different internal architecture and function. This subsystem receives all non-critical protocol data units and requests to process certain related protocol entities (ICMP, SNMP, RIP, OSPF, BGP, etc.). Any protocol data unit that needs to be sent on any network by these protocol entities is sent to the proper network interface, as if it was just another IP datagram relayed from another network.

The CPU subsystem can communicate with all other network interfaces through the exchange of coded messages across the switch fabric (or on a separate control bus [44,47]). IP datagrams generated by the CPU protocol entities are also coded in the same format. They carry the IP address of the next hop. For that, the CPU needs to access its individual routing table. This routing table can be the master table of the entire router. All other routing tables in the network interfaces will be exact replicas (or summaries in compressed table format) of the master table. Routing table updates in the network interfaces can then be done by broadcasting (if the switch fabric is capable of that) or any other suitable data push technique. Any update information (e.g., QoS policies, access control policies, packet drop policies, etc.) originated in the CPU has to be broadcast to all network interfaces. Such special data segments are received by the network interfaces which take care of the actual write operation in their forwarding tables. Updates to the routing table in the CPU are done either by the various routing protocol entities or by management action. This centralization is reasonable since routing changes are assumed to happen infrequently and not particularly time critical. The CPU can also be configured to handle any packet whose destination address cannot be found in the forwarding table in the network interface card.

3.5.3. Fast path or slow path?

It is not always obvious which router functions are to be implemented in the fast path or slow path. Some router designers may choose to include the ARP processing in the fast path instead of in the slow path of a router for performance reasons, and because ARP needs direct access to the physical network. Other may argue for ARP im-

plementation in the slow path instead of the fast path. For example, in the ARP used for Ethernet [14], if a router gets a datagram to an IP address whose Ethernet address it does not know, it is supposed to send an ARP message and hold the datagram until it gets an ARP reply with the necessary Ethernet address. When the ARP is implemented in the slow path, datagrams for which the destination link layer address is unknown are passed to the CPU, which does the ARP and, once it gets the ARP reply, forwards the datagram and incorporates the link-layer address into future forwarding tables in the network interfaces.

There are other functions which router designers may argue to be not critical and are more appropriate to be implemented in the slow path. IP packet fragmentation and reassembly, source routing option, route recording option, timestamp option, and ICMP message generation are examples of such functions. It can be argued that packets requiring these functions are rare and can be handled in the slow path: a practical product does not need to be able to perform “wire-speed” routing when infrequently used options are present in the packet. Since such packets having such options comprise a small fraction of the total traffic, they can be handled as exception conditions. As a result, such packets can be handled by the CPU, i.e., the slow path. For IP packet headers with error, generally, the CPU can instruct the inbound network interface to discard the errored datagram [48]. In some cases, the CPU will generate an ICMP message. Alternatively, in the cache-based scheme [45], templates of some common ICMP messages such as the TimeExceeded message are kept in the forwarding engine and these can be combined with the IP header to generate a valid ICMP message.

An IP packet can be fragmented by a router, that is, a single packet can be segmented into multiple, smaller packets and transmitted onto the output ports. This capability allows a router to forward packets between ports where the output is incapable of carrying a packet of the desired length; that is, the MTU of the output port is less than that of the input port. Fragmentation is good in the sense that it allows communication between

end systems connected through links with dissimilar MTUs. It is bad in that it imposes a significant processing burden on the router, which must perform more work to generate the resulting multiple output datagrams from the single input IP datagram. It is also bad from a data throughput point of view because, when one fragment is lost, the entire IP datagram must be retransmitted. The main arguments for implementing fragmentation in the slow path is that IP packet fragmentation can be considered an “exception condition”, outside of the fast path. Now that IP MTU discovery [13] is prevalent, fragmentation should be rare.

Reassembly of fragments may be necessary for packets destined for entities within the router itself. These fragments may have been generated either by other routers in the path between the sender and the router in question or by the original sending end system itself. Although fragment reassembly can be a resource-intensive process (both in CPU cycles and memory), the number of packets sent to the router is normally quite low relative to the number of packets being routed through. The number of fragmented packets destined for the router is a small percentage of the total router traffic. Thus, the performance of the router for packet reassembly is not critical and can be implemented in the slow path.

The fast path and slow path functions are summarized in Fig. 9. Fig. 9 further categorizes the slow path router functions into two: those performed on a packet-by-packet basis (that is, optional or exception conditions) and those performed as background tasks.

3.5.4. Protocol entities and IP processing in the distributed router architecture

The IP protocol is the most extensive entity in the packet processing path of an IP router and, thus, IP processing typically determines the achievable performance of a router. Therefore, a decomposition of IP that enables efficient multiprocessing is needed in the distributed router architecture. An example of a typical functional partitioning in the distributed router architecture is shown in Fig. 10.

This distributed multiprocessing architecture, means that the various processing elements can work in parallel on their own tasks with little dependence on the other processors in the system. This architecture decouples the tasks associated with determining routes through the network from the time-critical tasks associated with IP processing. The results of this is an architecture with high levels of aggregate system performance and the ability to scale to increasingly higher performance levels.

Network interface cards built with general-purpose processors and complex communication protocols tend to be more expensive than those built using ASICs and simple communication protocols. Choosing between ASICs and general-purpose processors for an interface card is not straightforward. General-purpose processors tend to be more expensive, but allow extensive port functionality. They are also available off-the-shelf, and their price/performance ratio improves yearly [63]. ASICs are not only cheaper, but can also provide operations that are specific to routing, such as traversing a Patricia tree.

Typical Router Fast-Path Functions	Typical Router Slow-Path Functions	
	Packet-by-Packet Operations	Background Tasks
<ul style="list-style-type: none"> - IP header validation - Routing table lookup & packet classification - Time-to-live (TTL) update - Checksum update 	<ul style="list-style-type: none"> - Fragmentation and reassembly - Source routing option - Route recording option - Timestamp option - ICMP message generation 	<ul style="list-style-type: none"> - Routing protocols (RIP, OSPF, BGP, etc.) - Network management (SNMP) - Router configuration (BOOTP, DHCP, etc.)

Fig. 9. Typical IP router fast-path and slow-path functions.

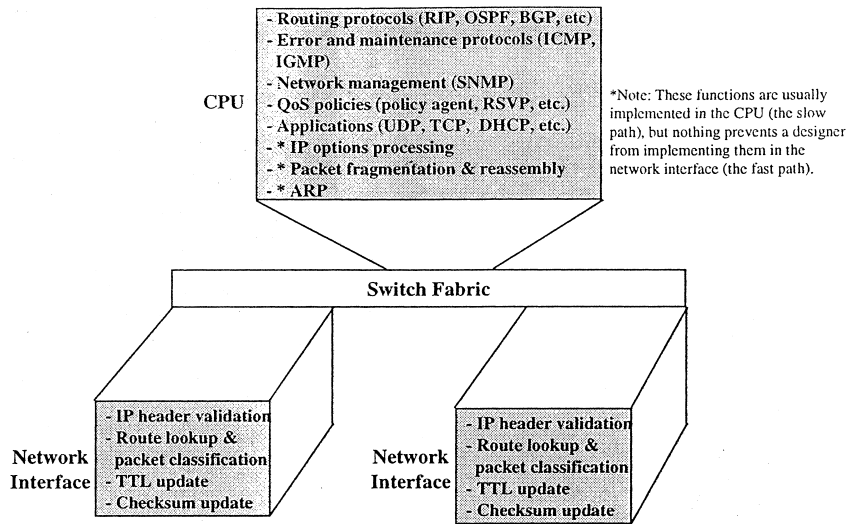


Fig. 10. An example functional partitioning in the distributed router architecture.

Moreover, the lack of flexibility with ASICs can be overcome by implementing functionality in the route processor (e.g., ARP, fragmentation and reassembly, IP options, etc.).

Some router designers often observe that the IPv4 specification is very stable and say that it would be more cost effective to implement the forwarding engine in an ASIC. It is argued that ASIC can reduce the complexity on each system board by combining a number of functions into individual chips that are designed to perform at high speeds. It may also be argued that the Internet is constantly evolving in a subtle way that require programmability and as such a fast processor is appropriate for the forwarding engine.

The forwarding database in a network interface consists of several cross-linked tables as illustrated in Fig. 11. This database can include IP routes (unicast and multicast), ARP tables, and packet filtering information for QoS and security/access control.

Now, let us take a generic shared memory router architecture and then trace the path of an IP packet as it goes through an ingress port and out of an egress port. The IP packet processing steps are shown in Fig. 12.

The IP packet processing steps are as follows with the step numbers corresponding to those in Fig. 12:

1. *IP header validation.* As a packet enters an ingress port, the forwarding logic verifies all Layer 3 information (header length, packet length, protocol version, checksum, etc.).
2. *Route lookup and header processing.* The router then performs an IP address lookup using the packet's destination address to determine the egress (or outbound) port, and performs all IP forwarding operations (TTL decrement, header checksum, etc.).
3. *Packet classification.* In addition to examining the Layer 3 information, the forwarding engine examines Layer 4 and higher layer packet attributes relative to any QoS and access control policies.
4. With the Layer 3 and higher layer attributes in hand, the forwarding engine performs one or more parallel functions:
 - associates the packet with the appropriate priority and the appropriate egress port(s) (an internal routing tag provides the switch fabric with the appropriate egress port information, the QoS priority queue the packet is to be stored in, and the drop priority for congestion control),
 - redirects the packet to a different (overridden) destination (ICMP redirect),
 - drops the packet according to a congestion control policy (e.g., RED [36]), or a security policy, and

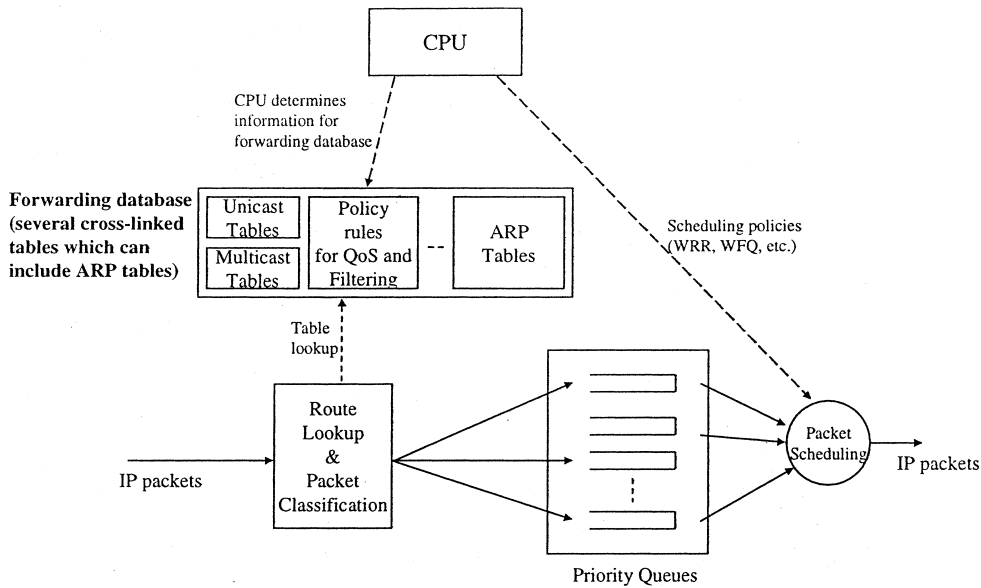


Fig. 11. Forwarding database consisting of several cross-linked tables.

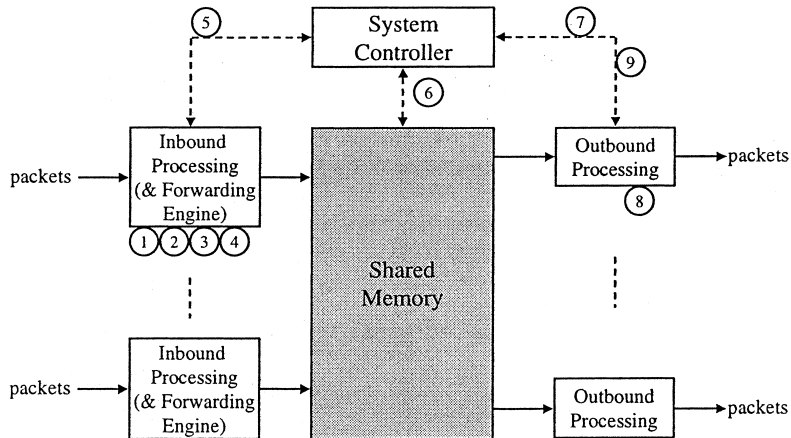


Fig. 12. IP packet processing in a shared memory router architecture.

- performs the appropriate accounting functions (statistics collection, etc.).
- 5. The forwarding engine notifies the system controller that a packet has arrived.
- 6. The system controller reserves a memory location for the arriving packet.
- 7. Once the packet has been passed to the shared memory, the system controller signals the appropriate egress port(s). For

- multicast traffic, multiple egress ports are signalled.
- 8. The egress port(s) extracts the packet from the known shared memory location using any of a number of algorithms: Weighted Fair Queueing (WFQ), Weighted Round-Robin (WRR), Strict Priority (SP), etc.
- 9. When the destination egress port(s) has retrieved the packet, it notifies the system control-

ler, and the memory location is made available for new traffic.

Section 3 does not attempt an exhaustive review of all possible router architectures. Instead, its aim is to identify the basic approaches that have been proposed and classify them according to the design trade-offs they represent. In particular, it is important to understand how different architectures fare in terms of performance, implementation complexity, and scalability to higher link speeds. In Section 4, we present an overview of the most common switch fabrics used in routers.

4. Typical switch fabrics of routers

Switch fabric design is a very well studied area, especially in the context of Asynchronous Transfer Mode (ATM) switches [5–7], so in this section, we examine briefly the most common fabrics used in router design. The switch fabric in a router is responsible for transferring packets between the other functional blocks. In particular, it routes user packets from the input modules to the appropriate output modules. The design of the switch fabric is complicated by other requirements such as multicasting, fault tolerance, and loss and delay priorities. When these requirements are considered, it becomes apparent that the switch fabric should have additional functions, e.g., concentration, packet duplication for multicasting if required, packet scheduling, packet discarding, and congestion monitoring and control.

Virtually all IP router designs are based on variations or combinations of the following basic approaches: shared memory; shared medium; distributed output buffered; space division (e.g., crossbar). Some important considerations for the switch fabric design are: throughput, packet loss, packet delays, amount of buffering, and complexity of implementation. For given input traffic, the switch fabric designs aim to maximize throughput and minimize packet delays and losses. In addition, the total amount of buffering should be minimal (to sustain the desired throughput without incurring excessive delays) and implementation should be simple.

4.1. Shared medium switch fabric

In a router, packets may be routed by means of a shared medium e.g., bus, ring, or dual bus. The simplest switch fabric is the bus. Bus-based routers implement a monolithic backplane comprising a single medium over which all inter-module traffic must flow. Data are transmitted across the bus using Time Division Multiplexing (TDM), in which each module is allocated a time slot in a continuously repeating transmission. However, a bus is limited in capacity and by the arbitration overhead for sharing this critical resource. The challenge is that it is almost impossible to build a bus arbitration scheme fast enough to provide non-blocking performance at multigigabit speeds.

An example of a fabric using a time-division multiplexed (TDM) bus is shown in Fig. 13. Incoming packets are sequentially broadcast on the bus (in a round-robin fashion). At each output, address filters examine the internal routing tag on each packet to determine if the packet is destined for that output. The address filters passes the appropriate packets through to the output buffers.

It is apparent that the bus must be capable of handling the total throughput. For discussion, we assume a router with N input ports and N output ports, with all port speeds equal to S (fixed size) packets per second. In this case, a packet time is defined as the time required to receive or transmit an entire packet at the port speed, i.e., $1/S$ s. If the bus operates at a sufficiently high speed, at least NS packets/s, then there are no conflicts for bandwidth and all queueing occurs at the outputs.

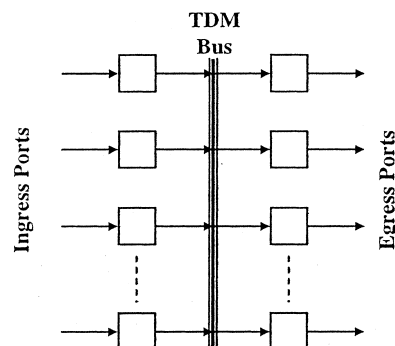


Fig. 13. Shared medium switch fabric: a TDM bus.

Naturally, if the bus speed is less than NS packets/s, some input queueing will probably be necessary.

In this architecture, the outputs are modular from each other, which has advantages in implementation and reliability. The address filters and output buffers are straightforward to implement. Also, the broadcast-and-select nature of this approach makes multicasting and broadcasting natural. For these reasons, the bus type switch fabric has found a lot of implementation in routers. However, the address filters and output buffers must operate at the speed of the shared medium, which could be up to N times faster than the port speed. There is a physical limit to the speed of the bus, address filters, and output buffers; these limit the scalability of this approach to large sizes and high speeds. Either the size N or speed S can be large, but there is a physical limitation on the product NS . As with the shared memory approach (to be discussed next), this approach involves output queueing, which is capable of the optimal throughput (compared to simple first-in first-out (FIFO) input queueing). However, the output buffers are not shared, and hence this approach requires more total amount of buffers than the shared memory fabric for the same packet loss rate.

4.2. Shared memory switch fabric

A typical architecture of a shared memory fabric is shown in Fig. 14. Incoming packets are typically converted from a serial to parallel form

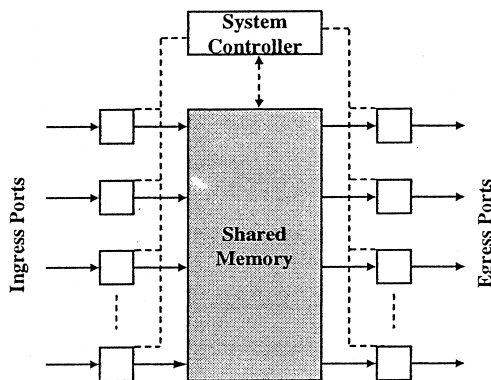


Fig. 14. A shared memory switch fabric.

which are then written sequentially into a (dual port) random access memory. Their packet headers with internal routing tags are typically delivered to a memory controller, which decides the order in which packets are read out of the memory. The outgoing packets are demultiplexed to the outputs, where they are converted from parallel to serial form. Functionally, this is an output queueing approach, where the output buffers all physically belong to a common buffer pool. The output buffered approach is attractive because it can achieve a normalized throughput of one under a full load [64,65]. Sharing a common buffer pool has the advantage of minimizing the amounts of buffers required to achieve a specified packet loss rate. The main idea is that a central buffer is most capable of taking advantage of statistical sharing. If the rate of traffic to one output port is high, it can draw upon more buffer space until the common buffer pool is (partially or) completely filled. For these reasons it is a popular approach for router design (e.g., [51–53,56,58–60]).

Unfortunately, the approach has its disadvantages. As the packets must be written into and read out from the memory one at a time, the shared memory must operate at the total throughput rate. It must be capable of reading and writing a packet (assuming fixed size packets) in every $1/NS$ s, that is, N times faster than the port speed. As the access time of random access memories is physically limited, this speed-up factor N limits the ability of this approach to scale up to large sizes and fast speeds. Either the size N or speed S can be large, but the memory access time imposes a limit on the product NS , which is the total throughput. Moreover, the (centralized) memory controller must process (the routing tags of) packets at the same rate as the memory. This might be difficult if, for instance, the controller must handle multiple priority classes and complicated packet scheduling. Multicasting and broadcasting in this approach will also increase the complexity of the controller. Multicasting is not natural to the shared memory approach but can be implemented with additional control circuitry. A multicast packet may be duplicated before the memory or read multiple times from the memory. The first approach obviously requires more memory be-

cause multiple copies of the same packet are maintained in the memory. In the second approach, a packet is read multiple times from the same memory location [78–80]. The control circuitry must keep the packet in memory until it has been read to all the output ports in the multicast group.

A single point of failure is invariably introduced in the shared memory-based design because adding a redundant switch fabric to this design is so complex and expensive. As a result, shared memory switch fabrics are best suited for small capacity systems.

4.3. Distributed output buffered switch fabric

The distributed output buffered approach is shown in Fig. 15. Independent paths exist between all N^2 possible pairs of inputs and outputs. In this design, arriving packets are broadcast on separate buses to all outputs. Address filters at each output determine if the packets are destined for that output. Appropriate packets are passed through the address filters to the output queues.

This approach offers many attractive features. Naturally there is no conflict among the N^2 independent paths between inputs and outputs, and hence all queuing occurs at the outputs. As stated earlier, output queuing achieves the optimal normalized throughput compared to simple FIFO input queuing [64,65]. Like the shared medium approach, it is also broadcast-and-select in nature

and, therefore, multicasting is natural. For multicasting, an address filter can recognize a set of multicast addresses as well as output port addresses. The address filters and output buffers are simple to implement. Unlike the shared medium approach, the address filters and buffers need to operate only at the port speed. All of the hardware can operate at the same speed. There is no speed-up factor to limit scalability in this approach. For these reasons, this approach has been taken in some commercial router designs (e.g., [57]).

Unfortunately, the quadratic N^2 growth of buffers means that the size N must be limited for practical reasons. However, in principle, there is no severe limitation on S . The port speed S can be increased to the physical limits of the address filters and output buffers. Hence, this approach might realize a high total throughput NS packets per second by scaling up the port speed S . The Knockout switch was an early prototype that suggested a trade-off to reduce the amount of buffers at the cost of higher packet loss [66]. Instead of N buffers at each output, it was proposed to use only a fixed number L buffers at each output (for a total of NL buffers which is linear in N), based on the observation that the simultaneous arrival of more than L packets (cells) to any output was improbable. It was argued that $L=8$ is sufficient under uniform random traffic conditions to achieve a cell loss rate of 10^{-6} for large N .

4.4. Space division switch fabric: the crossbar switch

Optimal throughput and delay performance is obtained using output buffered switches. Moreover, since upon arrival, the packets are immediately placed in the output buffers, it is possible to better control the latency of the packet. This helps in providing QoS guarantees. While this architecture appears to be especially convenient for providing QoS guarantees, it has serious limitations: the output buffered switch memory speed must be equal to at least the aggregate input speed across the switch. To achieve this, the switch fabric must operate at a rate at least equal to the aggregate of all the input links connected to the switch. However, increasing line rate (S) and increasing switch size (N) make it extremely difficult to significantly

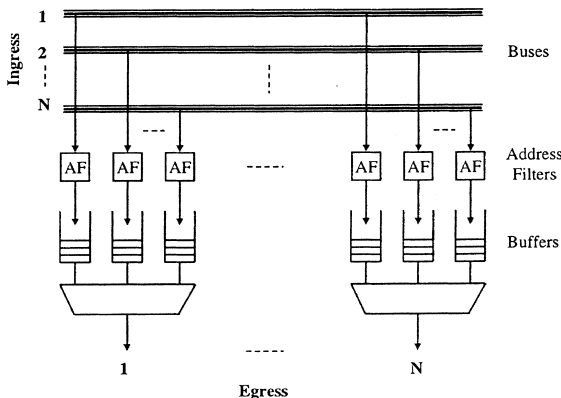


Fig. 15. A distributed output buffered switch fabric.

speedup the switch fabric, and also build memories with a bandwidth of order $O(NS)$.

At multigigabit and terabit speeds it becomes difficult to build output buffered switches. As a result some high-speed implementations are based on the input buffered switch architecture. One of the most popular interconnection networks used for building input buffered switches is the crossbar because of its (i) low cost, (ii) good scalability and (iii) non-blocking properties. Crossbar switches have an architecture that, depending on the implementation, can scale to very high bandwidths. Considerations of cost and complexity are the primary constraints on the capacity of switches of this type. The crossbar switch (see Fig. 16) is a simple example of a space division fabric which can physically connect any of the N inputs to any of the N outputs. An input buffered crossbar switch has the crossbar fabric running at the link rate. In this architecture buffering occurs at the inputs, and the speed of the memory does not need to exceed the speed of a single port. Given the current state of technology, this architecture is widely considered to be substantially more scalable than output buffered or shared memory switches. However, the crossbar architecture presents many technical challenges that need to be overcome in order to provide bandwidth and delay guarantees. Examples of commercial routers that use crossbar switch fabrics are [54,55,61].

We start with the issue of providing bandwidth guarantees in the crossbar architecture. For the case where there is a single FIFO queue at each

input, it has long been known that a serious problem referred to as head-of-line (HOL) blocking [64] can substantially reduce achievable throughput. In particular, the well-known results of [64] is that for uniform random distribution of input traffic, the achievable throughput is only 58.6%. Moreover, Li [67] has shown that the maximum throughput of the switch decreases monotonically with increasing burst size. Considerable amount of work has been done in recent years to build input buffered switches that match the performance of an output buffered switch. One way of reducing the effect of HOL blocking is to increase the speed of the input/output channel (i.e., the speedup of the switch fabric). Speedup is defined as the ratio of the switch fabric bandwidth and the bandwidth of the input links. There have been a number of studies such as [68,69] which showed that an input buffered crossbar switch with a single FIFO at the input can achieve about 99% throughput under certain assumptions on the input traffic statistics for speedup in the range of 4–5. A more recent simulation study [70] suggested that speedup as low as 2 may be sufficient to obtain performance comparable to that of output buffered switches.

Another way of eliminating the HOL blocking is by changing the queueing structure at the input. Instead of maintaining a single FIFO at the input, a separate queue per each output can be maintained at each input (virtual output queues (VOQs)). However, since there could be contention at the inputs and outputs, there is a necessity for an arbitration algorithm to schedule packets between various inputs and outputs (equivalent to the matching problem for bipartite graphs). It has been shown that an input buffered switch with VOQs can provide asymptotic 100% throughput using a maximum matching algorithm [71]. However, the complexity of the best known maximum match algorithm is too high for high-speed implementation. Moreover, under certain traffic conditions, maximum matching can lead to starvation. Over the years, a number of maximal matching algorithms have been proposed [72–75].

As stated above, increasing the speedup of the switch fabric can improve the performance

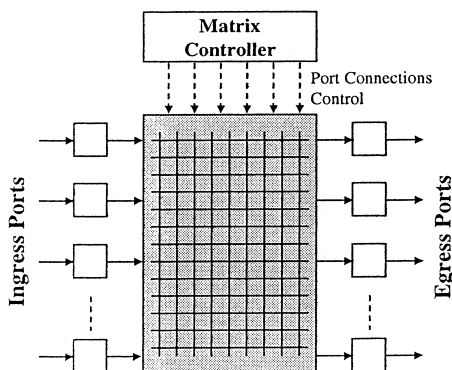


Fig. 16. A crossbar switch.

of an input buffered switch. However, when the switch fabric has a higher bandwidth than the links, buffering is required at the outputs too. Thus a combination of input buffered and output buffered switch is required, i.e., Combined Input and Output Buffered (CIOB). The goal of most designs then is to find the minimum speedup required to match the performance of an output buffered switch using a CIOB and VOQs. McKeown et al. [76] shown that a CIOB switch with VOQs is always work conserving if speedup is greater than $N/2$. In a recent work, Prabhakar et al. [77] showed that a speed of 4 is sufficient to emulate an output buffered switch (with an output FIFO) using a CIOB switch with VOQs.

Multicast in the space division fabrics is simple to implement but has some consequences. For example, a crossbar switch (with input buffering) is naturally capable of broadcasting one incoming packet to multiple outputs. However, this would aggravate the HOL blocking at the input buffers. Approaches to alleviate the HOL blocking effect would increase the complexity of buffer control. Other inefficient approaches in crossbar switches require an input port to write out multiple copies of the packet to different output ports one at a time. This does not support the one-to-many transfers required for multicasting as in the shared bus architecture and the fully distributed output buffered architectures. The usual concern about making multiple copies is that it reduces effective switch throughput. Several approaches for handling multicasting in crossbar switches have been proposed, e.g., [81]. Generally, multicasting increases the complexity of space division fabrics.

4.5. Other issues in router switch fabric design

We have described above four typical design approaches for router switch fabrics. Needless to say, endless variations of these designs can be imagined but the above are the most common fabrics found in routers. There are other issues applicable to understanding the trade-offs involved in any new design. We discuss some of these issues next.

4.5.1. Construction of large router switch fabrics

With regard to the construction of large switch fabrics, most of the four basic switch fabric design approaches are capable of realizing routers of limited throughput. The shared memory and shared medium approaches can achieve a throughput limited by memory access time. The space division approach has no special constraints on throughput or size, only physical factors do limit the maximum size in practice [82]. There are physical limits to the circuit density and number of input/output (I/O) pins. Interconnection complexity and power dissipation become more difficult issues with fabric size [83]. In addition, reliability and repairability become difficult with size. Modifications to maximize the throughput of space division fabrics to address HOL blocking increases the implementation complexity.

It is generally accepted that large router switch fabrics of 1 terabits per second (Tbps) throughput or more cannot be realized simply by scaling up a fabric design in size and speed. Instead, large fabrics must be constructed by interconnection of switch modules of limited throughput. The small modules may be designed following any approach, and there are various ways to interconnect them [84–86].

4.5.2. Fault tolerance and reliability

With the rapid growth of the Internet and the emergence of growing competition between Internet Service Providers (ISPs), reliability has become an important issue for IP routers. In addition, multigigabit routers will be deployed in the core of enterprise networks and the Internet. Traffic from thousands of individual flows pass through the switch fabric at any given time [46]. Thus, the robustness and overall availability of the switch fabric becomes a critically important design parameter. As in any communication system, fault tolerance is achieved by adding redundancy to the crucial components. In a router, one of the most crucial components is the packet routing and buffering fabric. Redundancy may be added in two ways: by duplicating copies of the switch fabric or by adding redundancy within the fabric [87–90]. In addition to redundancy, other considerations

include detection of faults and isolation and recovery.

4.5.3. Buffer management and quality of service

The prioritization of mission critical applications and the support of IP telephony and video conferencing create the requirement for supporting QoS enforcement with the switch fabric. These applications are sensitive to both absolute latency and latency variations.

Beyond best-effort service, routers are beginning to offer a number of QoS or priority classes. Priorities are used to indicate the preferential treatment of one traffic class over another. The switch fabrics must handle these classes of traffic differently according to their QoS requirements. In the output buffered switch fabric, for example, typically the fabric will have multiple buffers at each output port and one buffer for each QoS traffic class. The buffers may be physically separate or a physical buffer may be divided logically into separate buffers.

Buffer management here refers to the discarding policy for the input of packets into the buffers (e.g., Random Early Detection (RED)), and the scheduling policy for the output of packets from the buffers (e.g., strict priority, weighted round-robin (WRR), weighted fair queueing (WFQ), etc.). Buffer management in the IP router involves both dimensions of time (packet scheduling) and buffer space (packet discarding). The IP traffic classes are distinguished in the time and space dimensions by their packet delay and packet loss priorities. We therefore see that buffer management and QoS support is an integral part of the switch fabric design.

5. Conclusions and open problems

IP provides a high degree of flexibility in building large and arbitrary complex networks. Interworking routers capable of forwarding aggregate data rates in the multigigabit and terabit per second range are required in emerging high performance networking environments. This paper has presented an evaluation of typical approaches proposed for designing high-speed routers. We

have focused primarily on the architectural overview and the design of the components that have the highest effect on performance.

First, we have observed that high-speed routers need to have enough internal bandwidth to move packets between its interfaces at multigigabit and terabit rates. The router design should use a switched backplane. Until very recently, the standard router used a shared bus rather than a switched backplane. While bus-based routers may have satisfied the early needs of IP networks, emerging demands for high bandwidth, QoS delivery, multicast, and high availability place the bus architecture at a significant disadvantage. For high speeds, one really needs the parallelism of a switch with superior QoS, multicast, scalability, and robustness properties. Second, routers need enough packet processing power to forward several million packets per second (Mpps). Routing table lookups and data movements are the major consumers of processing cycles. The processing time of these tasks does not decrease linearly if faster processors are used. This is because of the sometimes dominating effect of memory access rate.

It is observed that while an IP router must, in general, perform a myriad of functions, in practice the vast majority of packets need only a few operations performed in real-time. Thus, the performance critical functions can be implemented in hardware (the fast path) and the remaining (necessary, but less time-critical) functions in software (the slow path). IP contains many features and functions that are either rarely used or that can be performed in the background of high-speed data forwarding (for example, routing protocol operation and network management). The router architecture should be optimized for those functions that must be performed in real-time, on a packet-by-packet basis, for the majority of the packets. This creates an optimized routing solution that route packets at high speed at a reasonable cost.

It has been observed in [63] that the cost of a router port depends on, (1) the amount and kind of memory it uses, (2) its processing power, and (3) the complexity of the protocol used for communication between the port and the route processor.

This means the design of a router involve trade-offs between performance, complexity, and cost.

Router ports built with general-purpose processors and complex communication protocols tend to be more expensive than those built using ASICs and simple communication protocols. Choosing between ASICs and general-purpose processors for an interface card is not straightforward. General-purpose processors are costlier, but allow extensive port functionality, are available off-the-shelf, and their price/performance ratio improves rapidly with time [63]. ASICs are not only cheaper, but can also provide operations that are specific to routing, such as traversing a Patricia tree. Moreover, the lack of flexibility with ASICs can be overcome by implementing functionality in the route processor.

The cost of a router port is also proportional to the type and size of memory on the port. SRAMs offer faster access times, but are more expensive than DRAMs. Buffer memory is another parameter that is difficult to size. In general, the rule of thumb is that a port should have enough buffers to support at least one bandwidth-delay product worth of packets, where the delay is the mean end-to-end delay and the bandwidth is the largest bandwidth available to TCP connections traversing that router. This sizing allows TCP to increase their transmission windows without excessive losses.

The cost of a router port is also determined by the complexity of the internal connections between the control paths and the data paths in the port card. In some designs, a centralized controller sends commands to each port through the switch fabric and the port's internal buffers. Careful engineering of the control protocol is necessary to reduce the cost of the port control circuitry and also the loss of command packets which will certainly need retransmission.

The design principles of router fabrics are well known, at least for some common switch fabrics. In many cases it is possible to leverage many of the design principles from ATM switches into next generation multigigabit routers. Although the design principles for most fabrics are well established by now, the design is complicated by considerations for fault tolerance, multicasting, and

queuing priorities. Fault tolerance implies the necessity of redundant resources in the fabric. The straightforward approach of duplicating parallel planes may also improve throughput. Alternatively, various ways exist to add redundancy within a single fabric plane. Multicasting depends on the fabric design because some design approaches inherently broadcast every packet. For space division fabrics, including large fabrics constructed by interconnection of small modules, multicasting can be complex.

It will also be necessary to manage multiple buffers to satisfy the QoS requirements of the different traffic classes. Buffer management consist of packet scheduling and discarding policies, which determine the manner in which packets are input into and output from the buffers. The packet scheduling policy attempts to satisfy the delay requirements as indicated by delay priorities. The packet discarding policy attempts to satisfy the loss requirements as indicated by packet loss priorities.

Significant advances have been made in router designs to address the most demanding customer issues regarding high-speed packet forwarding (e.g., route lookup algorithms, high-speed switching cores and forwarding engines), low per-port cost, flexibility and programmability, reliability, and ease of configuration. While these advances have been made in the design of IP routers, some important open issues still remain to be resolved. These include packet classification and resource provisioning, improved price/performance router designs, "active networking" [91] and ease of configuration, reliability and fault tolerance designs, and Internet billing/pricing. Extensive work is being carried out both in the research community and industry to address these problems.

References

- [1] P. Newman, T. Lyon, G. Minshall, Flow labelled IP: a connectionless approach to ATM, in: Proceedings of the IEEE Infocom '96, San Francisco, CA, March 1996, pp. 1251–1260.
- [2] Y. Katsube, K. Nagami, H. Esaki, Toshiba's router architecture extensions for ATM: overview, in: IETF RFC 2098, April 1997.

- [3] Y. Rekhter, B. Davie, D. Katz, E. Rosen, G. Swallow, Cisco systems tag switching architecture overview, in: IETF RFC 2105, February 1997.
- [4] F. Baker, Requirements for IP Version 4 routers, in: IETF RFC 1812, June 1995.
- [5] H. Ahmadi, W. Denzel, A survey of modern high performance switching techniques, *IEEE J. Selected Areas Commun.* 7 (1989) 1091–1103.
- [6] F. Tobagi, Fast packet switch architectures for broadband integrated services digital networks, *Proc. IEEE* 78 (1990) 133–178.
- [7] R.Y. Awdeh, H.T. Mouftah, Survey of ATM switch architectures, *Computer Net. and ISDN Syst.* 27 (1995) 1567–1613.
- [8] R. Perlman, *Interconnections: Bridges and Routers*, Addison-Wesley, Reading, MA, 1992.
- [9] C. Huitema, *Routing in the Internet*, Prentice Hall, Englewood Cliffs, NJ, 1996.
- [10] J. Moy, *OSPF: Anatomy of an Internet Routing Protocol*, 1998.
- [11] W.R. Stevens, *TCP/IP Illustrated*, vol. 1, The Protocols, Addison-Wesley, Reading, MA, 1994.
- [12] C.A. Kent, J.C. Mogul, Fragmentation considered harmful, *Computer Commun. Rev.* 17 (5) (1987) 390–401.
- [13] J. Mogul, S. Deering, Path MTU Discovery IETF RFC 1191, April 1990.
- [14] D. Plummer, Ethernet address resolution protocol: or converting network protocol addresses to 48-bit ethernet addresses for transmission on ethernet hardware, in: IETF RFC 826, November 1982.
- [15] T. Bradley, C. Brown, Inverse address resolution protocol, in: IETF RFC 1293, January 1992.
- [16] V. Fuller et al., Classless inter-domain routing, in: IETF RFC 1519, June 1993.
- [17] K. Sklower, A tree-based packet routing table for Berkeley Unix, in: USENIX, Winter '91, Dallas, TX, 1991.
- [18] W. Doeringer, G. Karjoth, M. Nassehi, Routing on longest-matching prefixes, *IEEE/ACM Trans. Networking* 4 (1) (1996) 86–97.
- [19] D.R. Morrison, Patricia – practical algorithm to retrieve information coded in alphanumeric, *J. ACM* 15 (4) (1968) 515–534.
- [20] D.C. Feldmeier, Improving gateway performance with a routing-table cache, in: Proceedings of the IEEE Infocom '88, New Orleans, LA, March 1988.
- [21] C. Partridge, Locality and route caches, in: NSF Workshop on Internet Statistics Measurement and Analysis, San Diego, CA, February 1996.
- [22] D. Knuth, *The Art of Computer Programming*, vol. 3, Sorting and Searching, Addison-Wesley, Reading, MA, 1973.
- [23] M. Degermark, Small forwarding tables for fast routing lookups, in: Proceedings of the ACM SIGCOMM '97, Cannes, France, September 1997.
- [24] H.Y. Tzeng, Longest prefix search using compressed trees, in: Proceedings of the Globecom '98, Sydney, Australia, November 1998.
- [25] M. Waldvogel, G. Varghese, J. Turner, B. Plattner, Scalable high speed IP routing lookup, in: Proceedings of the ACM SIGCOMM '97, Cannes, France, September 1997.
- [26] V. Srinivasan, G. Varghese, Faster IP lookups using controlled prefix expansion, in: Proceedings of the ACM SIGMETRICS, May 1998.
- [27] S. Nilsson, G. Karlsson, Fast address look-up for internet routers, in: Proceedings of the IEEE Broadband Communications '98, April 1998.
- [28] E. Filippi, V. Innocenti, V. Vercellone, Address lookup solutions for gigabit switch/router, in: Proceedings of the Globecom '98, Sydney, Australia, November 1998.
- [29] A.J. McAuley, P. Francis, Fast routing table lookup using CAMs, in: Proceedings of the IEEE Infocom '93, San Francisco, CA, March 1993, pp. 1382–1391.
- [30] T.B. Pei, C. Zukowski, Putting routing tables in silicon, *IEEE Network* 6 (1992) 42–50.
- [31] M. Zitterbart et al., HearT: high performance routing table lookup, in: Fourth IEEE Workshop on Architecture and Implementation of High Performance Communications Subsystems, Thessaloniki, Greece, June 1997.
- [32] P. Gupta, S. Lin, N. McKeown, Routing lookups in hardware at memory access speeds, in: Proceedings of the IEEE Infocom '98, March 1998.
- [33] C. Rigney, A. Rubens, W. Simpson, S. Willens, Remote Authentication Dial-In User Service (RADIUS), in: IETF RFC 2138, April 1997.
- [34] V. Jacobson, Congestion avoidance and control, *Comput. Commun. Rev.* 18 (4) (1988).
- [35] W. Stevens, TCP slow start, congestion avoidance, fast retransmit and fast recovery algorithms, in: IETF RFC 2001, January 1997.
- [36] S. Floyd, V. Jacobson, Random early detection gateways for congestion avoidance, *IEEE/ACM Trans. Networking* 1 (4) (1993).
- [37] A. Demers, S. Keshav, S. Shenker, Analysis and simulation of a fair queueing algorithm, *Journal of Internetworking: Research and Experience* 1 (1990).
- [38] S. Floyd, V. Jacobson, Link sharing and resource management models for packet networks, *IEEE/ACM Trans. Networking* 3 (4) (1995).
- [39] L. Zhang, S. Deering, D. Estrin, S. Shenker, D. Zappala, RSVP: A new resource reservation protocol, in: IEEE Network, September 1993.
- [40] K. Nichols, S. Blake, Differentiated services operational model and definitions, in: IETF Work in progress.
- [41] S. Blake, An architecture for differentiated services, in: IETF Work in progress.
- [42] S.F. Bryant, D.L.A. Brash, The DECNIS 500/600 multi-protocol bridge router and gateway, *Digital Technical Journal* 5 (1) (1993).
- [43] P. Marimuthu, I. Viniotis, T.L. Sheu, A parallel router architecture for high speed LAN internetworking, in: Proceedings of the 17th IEEE Conference on Local Computer Networks, Minneapolis, Minnesota, September 1992.

- [44] S. Asthana, C. Delph, H.V. Jagadish, P. Krzyzanowski, Towards a gigabit IP router, *Journal of High Speed Networks* 1 (4) (1992).
- [45] C. Partridge, A 50Gb/s IP router, *IEEE/ACM Trans. Networking* 6 (3) (1998) 237–248.
- [46] K. Thomson, G.J. Miller, R. Wilder, Wide-area traffic patterns and characteristics, in: *IEEE Network*, December 1997.
- [47] V.P. Kumar, T.V. Lakshman, D. Stiliadis, Beyond best effort: router architectures for the differentiated services of tomorrow's internet, in: *IEEE Commun. Mag.*, May 1998, pp. 152–164.
- [48] A. Tantawy, O. Koufopavlou, M. Zitterbart, J. Abler, On the design of a multigigabit IP router, *Journal of High Speed Networks* 3 (1994) 209–232.
- [49] O. Koufopavlou, A. Tantawy, M. Zitterbart, IP-Routing among gigabit networks, in: S. Rao (Ed.), *Interoperability in Broadband Networks*, IOS Press, 1994, pp. 282–289.
- [50] O. Koufopavlou, A. Tantawy, M. Zitterbart, A comparison of gigabit router architectures, in: E. Fdida (Ed.), *High Performance Networking*, North-Holland, Amsterdam, 1994.
- [51] Implementing the Routing Switch: How to Route at Switch Speeds and Switch Costs, White Paper, Bay Networks, 1997.
- [52] Catalyst 8510 Architecture, White Paper, Cisco Systems, 1998.
- [53] Catalyst 8500 Campus Switch Router Series, White Paper, Cisco Systems, 1998.
- [54] Cisco 12000 Gigabit Switch Router, White Paper, Cisco Systems, 1997.
- [55] Performance Optimized Ethernet Switching, Cajun White Paper #1, Lucent Technologies.
- [56] Internet Backbone Routers and Evolving Internet Design, White Paper, Juniper Networks, September 1998.
- [57] The Integrated Network Services Switch Architecture and Technology, White Paper, Berkeley Networks, 1997.
- [58] Torrent IP9000 Gigabit Router, White Paper, Torrent Networking Technologies, 1997.
- [59] Wire-Speed IP Routing, White Paper, Extreme Networks, 1997.
- [60] PE-4884 Gigabit Routing Switch, White Paper, Packet Engines, 1997.
- [61] GRF 400 White Paper: A Practical IP Switch for Next-Generation Networks, White Paper, Ascend Communications, 1998.
- [62] Rule Your Networks: An Overview of StreamProcessor Applications, White Paper, NEO Networks, 1997.
- [63] S. Keshav, R. Sharma, Issues and Trends in Router Design, *IEEE Commun. Mag.*, May 1998, pp. 144–151.
- [64] M. Karol, M. Hluchyj, S. Morgan, Input versus output queueing on a space-division packet switch, *IEEE Trans. Commun.* 35 (1987) 1337–1356.
- [65] M. Hluchyj, M. Karol, Queueing in high-performance packet switching, *IEEE J. Selected Areas Commun.* 6 (1988) 1587–1597.
- [66] Y.S. Yeh, M. Hluchyj, A.S. Acampora, The knockout switch: a simple modular architecture for high-performance packet switching, *IEEE J. Selected Areas Commun.* 5 (8) (1987) 1274–1282.
- [67] S.Q. Li, Performance of a non-blocking space-division packet switch with correlated input traffic, in: *Proceedings of the IEEE Globecom '89*, 1989, pp. 1754–1763.
- [68] C.Y. Chang, A.J. Paulraj, T. Kailath, A broadband packet switch architecture with input and output queueing, in: *Proceedings of the Globecom '94*, 1994.
- [69] I. Iliadis, W. Denzel, Performance of packet switches with input and output queueing, in: *Proc. ICC '90*, 1990.
- [70] R. Guerin, K.N. Sivarajan, Delay and throughput performance of speed-up input-queueing packet switches, in: *IBM Research Report RC 20892*, June 1997.
- [71] N. McKeown, V. Anantharam, J. Walrand, Achieving 100% Throughput in an input-queued switch, in: *Proceedings of the IEEE Infocom '96*, 1996, pp. 296–302.
- [72] T.E. Anderson, S.S. Owicki, J.B. Saxe, C.P. Thacker, High speed switch scheduling for local area networks, *ACM Trans. Computer Sys.* 11 (4) (1993) 319–352.
- [73] D. Stiliadis, A. Verma, Providing bandwidth guarantees in an input-buffered crossbar switch, *Proc. IEEE Infocom '95* (1995) 960–968.
- [74] C. Lund, S. Phillips, N. Reingold, Fair prioritized scheduling in an input-buffered switch, in: *Proceedings of the Broadband Communications*, 1996.
- [75] A. Mekittikul, N. McKeown, A practical scheduling algorithm to achieve 100% throughput in input-queued switches, in: *Proceedings of the IEEE Infocom '98*, March 1998.
- [76] N. McKeown, B. Prabhakar, M. Zhu, Matching output queueing with combined input and output queueing, in: *Proceedings of the 35th Annual Allerton Conference on Communications, Control and Computing*, October 1997.
- [77] B. Prabhakar, N. McKeown, On the speedup required for combined input and output queueing switching, *Computer Systems Lab, Technical Report CSL-TR-97-738*, Stanford University.
- [78] N. Endo et al., Shared buffer memory switch for an ATM exchange, *IEEE Trans. Commun.* 41 (1993) 237–245.
- [79] H. Kuwahara, A shared buffer memory switch for an ATM exchange, in: *Proceedings of the ICC '89*, Boston, 1989, pp. 118–122.
- [80] Y. Sakurai, ATM switching systems for B-ISDN, *Hitachi Rev.* 40 (1990) 193–198.
- [81] N. McKeown, Fast Switched Backplane for a Gigabit Switched Router, Technical Report, Department of Electrical Engineering, Stanford University.
- [82] T. Lee, A modular architecture for very large packet switches, *IEEE Trans. Commun.* 38 (1990) 1097–1106.
- [83] T. Bandwell, Physical design issues for very large scale atm switching systems, *IEEE J. Selected Areas Commun.* 9 (1991) 1227–1238.
- [84] K. Murano, Technologies towards broadband ISDN, *IEEE Commun. Mag.* 28 (1990) 66–70.

- [85] A. Itoh et al., Practical implementation and packaging technologies for a large-scale ATM switching system, *IEEE J. Selected Areas Commun* 9 (1991) 1280–1288.
- [86] T. Kozaki et al., 32×32 Shared buffer type ATM switch VLSI's for B-ISDN's, *IEEE J. Selected Areas Commun.* 9 (1991) 1239–1247.
- [87] V. Kumar, S. Wang, Reliability enhancement by time and space redundancy in multistage interconnection networks, *IEEE Trans. Reliability* 40 (1991) 461–473.
- [88] G. Adams et al., A survey and comparison of fault-tolerant multistage interconnection networks, *IEEE Comput.* 20 (1987) 14–27.
- [89] D. Agrawal, Testing and fault-tolerance of multistage interconnection networks, *IEEE Comput.* 15 (1982) 41–53.
- [90] A. Itoh, A fault-tolerant switching network for BISDN, *IEEE J. Selected Areas Commun.* 9 (1991) 1218–1226.
- [91] D. Tennenhouse et al., A Survey of Active Network Research, in: *IEEE Commun. Mag.*, January 1997.
- [92] A. Rijssinghani, Computation of the Internet Checksum via Incremental Update, in: *IETF RFC 1624*, May 1994.



James Aweya received his B.Sc. (Hon) degree from the University of Science and Technology, Kumasi, Ghana, and his M.Sc. degree from the University of Saskatchewan, Saskatoon, Canada, all in electrical engineering. He recently completed his Ph.D. studies in electrical engineering at the University of Ottawa. In 1986 he joined the Electricity Corporation of Ghana, where he worked on the design and operation of electric distribution systems. From 1987 to 1990, he was with the Volta River Authority where he was engaged in the instrumentation and control of electric power systems. Since 1996 he has been with Nortel Networks where he is a Systems Design Engineer. He is currently involved in the design of resource management functions for IP and ATM networks, architectural issues and design of IP routers, and network architectures and protocols. He has published more than 25 journal and conference papers and has 5 patents pending. In addition to computer networks and protocols, he has other interest in fuzzy logic control, neural networks, and application of artificial intelligence to computer networking.