(III)

90        10        30

# Abstract

Continuous rapid growth of the Internet in recent years makes it the most probable future integrated services network. However, current Internet architecture is inadequate in providing real-time applications. It cannot guarantee delay bound requirements of real-time applications. Moreover, non-real-time applications may be terminated if real-time traffic causes congestion. Internet 2 is thus proposed to meet future needs. In order to support the realization of future broadband Internet with guarantee of wide quality of service (QoS) requirements, this integrated project constantly improves the results in the previous year and further investigates the following key technologies:

A. High-capacity (Gigabit) routers: In this project, we design and develop high-capacity routers. We have developed a prototype router, which consists of a data path module, a queueing module, a classifier, and a scheduler in last two years. In this sub-project, we continue developing two key technologies including scheduling and switching. In the prototype, we implemented SCFQ mechanism due to the consideration of implementation complexity. However, the delay bound and fairness of SCFQ is not as good as WFQ. In new design, we implement the FFQ to improve the performance of the scheduler. In the queueing management of the prototype, we adopted the shared memory architecture. It has a disadvantage that its capacity is limited for switch with many ports because of heavy memory access. It is possible to develop a large-scale input queueing system, but there exists a head-of-line blocking problem. In this sub-project, we investigate the CIOQ architecture. We have developed a algorithm to emulate an output-queueing switch and evaluate the performance of CIOQ with finite buffers

B. Measurement-based Admission Control and Congestion Avoidance Schemes for Controlled-Load Service in Broadband Internet: Many real-time applications are capable of adapting their transmission to the network state and can as well tolerate occasional delay

bound violations in the presence of transient network congestion. For this type of applications, an absolutely reliable bound on packet delivery times is not required. Moreover, many non-real-time applications such as on-line transaction processing and distributed simulation would desire a congestion-free packet delivery service from the network. This new type of service is called "Congestion-Free Service" with the guarantee of a maximum packet loss rate. This work considers resource allocation in the support of "Congestion- Free Service" in Broadband Internet. First, we studied traffic characterization under different measurement models via analyzing traffic traces collected from National Taiwan University campus network. The preliminary results show that traffic load can be approximated by the Normal distribution. In the second part of the work, we proposed a dynamic bandwidth and queue management scheme to support "Congestion-Free Service." To better understand the characteristics of the broadband Internet, we have developed a tool called *AppMeasure* which provides three sets of service: *IPFlow*, *AnyTrace*, and *AppFlow*. IPFlow provides on-line tracking and reporting of individual IP flows. In AnyTrace, user can specify the *number of bytes* of a packet to trace. It serves as a measurement base to enable *all-layer* traffic accounting and analysis, and differentiated charging based on Quality of Service. Appflow tracks and reports any *type of traffic* specified by the user. It aims to integrate with policy-based QoS service

C. QoS routing: This sub-project describes the implementation of a class-based QoS routing algorithm, which is designed for low blocking probability and low overload. The class-based QoS routing algorithm is designed in the per-pair granularity. We design a software architecture and divide it into several modules. Then we describe what the modules do and how they work. At least the performance evaluation of the implementation is discussed. The results show that the costs, such as processing time of path computation and memory requirement of routing table, are expensive. The costs are what we have to pay for QoS routing supporting.

D. Integrated Service and DiffServ Technology: Network processors are emerging as a programmable alternative to the traditional ASIC-based solutions in scaling up the data-plane processing of network services. They serve as co-processors to offload data-plane traffic from the original general-purpose microprocessor. In this work, we illustrate the process and investigate performance issues in prototyping a DiffServ edge router with IXP1200, which consissts of one control-plane StrongARM core processor and six data-plane microengines, and stores classification and scheduling per-flow policy rules at SRAM and packets at SDRAM. The external benchmark shows that though the system can achieve aggregated wire-speed of 1.8Gbps in simple IP forwarding, the throughput drops to 200~300Mbps when performing DiffServ due to the double bottlenecks of SRAM and microengines. Through internal benchmarks, we found that performance bottlenecks may shift from one place to another given different network services and algorithms. For simple IP forwarding services, SDRAM is a nature bottleneck. However, it could shift to SRAM or microengines if heavy table access or computation is involved, respectively. We also identify the design pitfall of the hardware called the "MAC buffer overflow".

*Keywords:* Broadband Internet, Quality of service, QoS routing, Gigabit router, Signaling, Traffic measurement, Network Processor, DiffServ

A.                      Gigabit

SCFQ (Self-Clocked Fair Queueing)                        WFQ (Weighted

fair queueing)                     FFQ (frame-based fair queueing)

(input

queueing)                                    （ head-of-line blocking ）

CIOQ (combined input and output-queued)            buffer

B.

IETF

off-line

C.

D.                                                                                                                      ASIC

(co-processors)                                          (general-purpose  processor)

(DiffServ edge router)

IXP1200                                             IXP1200

StrongARM                    (core  processor)

(classification)          (scheduling)              SRAM                        SDRAM

(input  port)                (throughput)    50Mbps

(Per-Hop Behavior)    500              (flow)                    SRAM

(bottleneck)

(forwarding service)          SDRAM

SRAM        microengine

IXP1200                                              ，

6

' (MAC buffer overflow)

Gigabit

Circuit

switching                       Packet switching



Quality  of  Service

Fast  Ethernet

100Mbps                 Ethernet Switch










:

Internet                                                Quality

of Service                        QoS control




                                                    (Quality of Service,

QoS)




        GPS  WFQ                              delay bound

                            SCFQ                                        delay bound

                                                            FFQ

                    WFQ




    (input-queued)




                                                                        (1)

                    (FIFO)

    maximal                            100%                    (2)

(speedup)                          /            (input/output  link)

            (buffer)                      combined input and output-queued (CIOQ)

                    CIOQ

(QoS, quality of service)　　　　IETF

(Integrated Service)[19]

(Differentiated Service) [12]

(absolute)　　　　　(delay bound)　　　　　　　　　(queueing)

(Guaranteed Service)[14] [18]

*vat   nv   vic*

(admission control)　　　　(packet scheduling)　　　　(buffer management)

AppMeasure

Blocking probability and route granularity are important issues for QoS routing protocols. In general, finer granularity yields higher blocking probability but also higher computational complexity. A new algorithm proposed in [24] to show us a new thought of QoS routing after integrated the

11

researches mentioned above. This QoS routing protocol, with a novel idea of per-class forwarding with routing marks, is designed in order to achieve the high scalability and low blocking probability. The mechanism of the routing is based on per-pair granularity. It is shown by simulation that, with a small number of routing marks, the routing algorithm yields competitive blocking probability as compared to the routing algorithm that routes flows independently.

In this sub-project, we try to implement the class-based routing algorithm as an extension to a routing daemon, named Zebra. When a router receives a QoS request, the router may trigger following actions: path calculation, marking, and forwarding. In this project, we will only focus on the issue of routing, i.e., the path calculation task. The path calculation will find a new route with the least cost and abundant residual bandwidth based on Dijkstra's algorithm, which is also adopted by OSPF

The increasing link bandwidth demands even faster nodal processing especially for the data-plane traffic. The nodal data-plane processing may range from routing table lookup to various classifications for firewall, DiffServ and Web switching. The traditional general-purpose processor architecture is no longer scalable enough for wire-speed processing so that some ASIC components or co-processors are commonly used to *offload* the data-plane processing, while leaving only control-plane processing to the original processor.

Many ASIC-driven products have been announced in the market, such as the acceleration cards for encryption/decryption [29], VPN gateways [30], Layer 3 switches [31], DiffServ routers [32] and Web switches [33]. While these ASICs indeed speedup the data-plane packet processing with special hardware blocks, much wider memory buses, and faster execution process, they lack flexibility in *reprogrammability* and have a long development cycle which is usually   months or even years.

Network processors are emerging as an alternative solution to ASICs for providing scalability for data-plane packet processing while retaining reprogrammability. In this study,

we adopt Intel IXP1200 [34] network processor shown in Fig.1 which consists of one StrongARM core and six co-processors referred as microengines, so that developers can embed the control-plane and data-plane traffic management modules into the StrongARM core and microengines, respectively. Scalability concern in data-plane packet processing could be satisfied with the four zero context switching overhead hardware contexts in each of the six microengines and the instructions specialized for networking.
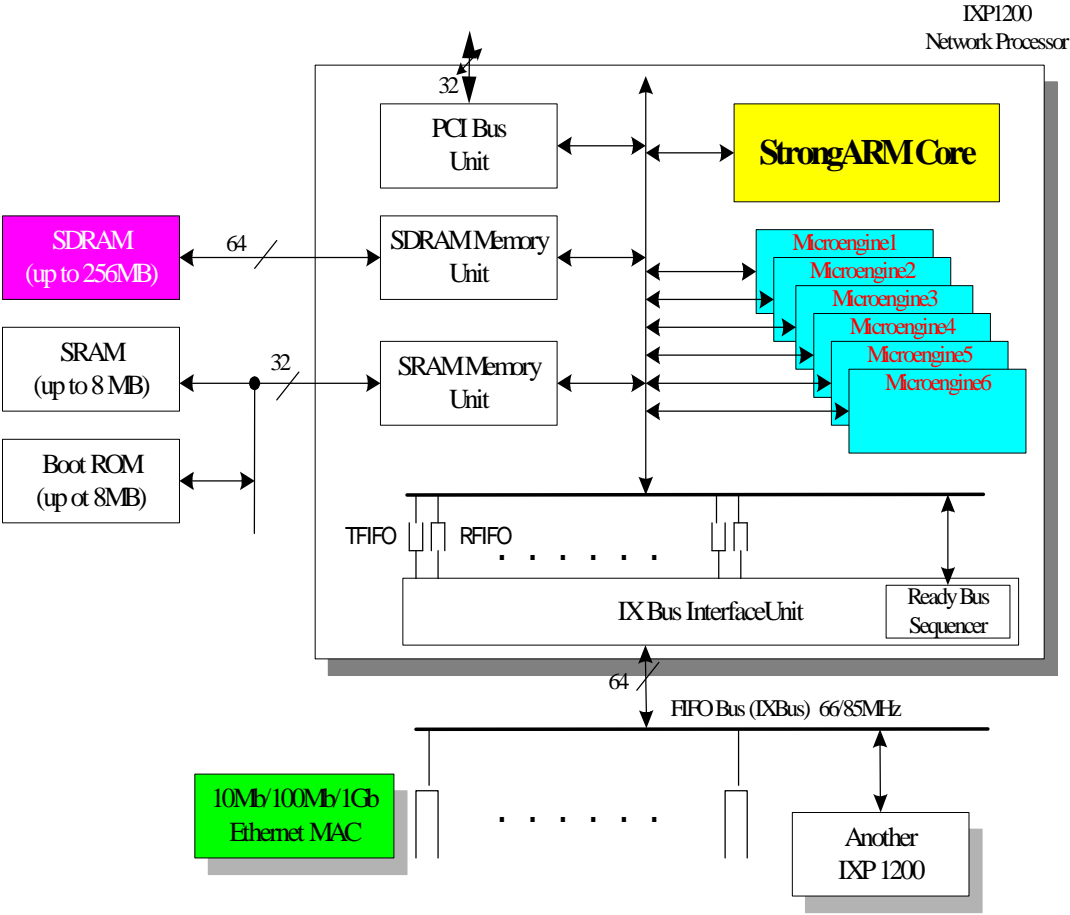


Fig. 1 Hardware architecture of IXP1200

# Gigabit

## A. FQ

FFQ

frame

frame

frame counter                                                             frame

frame

frame

state machine     I/Q

(finish time)

FFQ

(quantization)                                                                                    (time stamp)

                                                                                                            64

                                                                                        64

                                FFQ                                                FFQ

**FFQ**

## B. CIOQ

CIOQ                        corssbar                                                        HOL

            LCF/MUF                                output queueuing

                                CIOQ

                                (service scheduling)

                CIOQ                                        crossbar switch                        buffer

        phase    corssbar

                                                CIOQ

                                                                        Least Cushion

First / Most Urgent First    (LCF/MUF)

                        crossbar                        CIOQ

15

**CIOQ**

CIOQ         buffer           buffer

            uniform      16X16

   buffer    cell         correlated

burst length L    0.8     3L 4L buffer

## A.

NetFlow   Cisco            Cisco

              NetFlow

   AppMeasure       NetFlow

NetFlow   IP    IP    IP    IP    IP

  TOS         AppMeasure

  header      IP

content

Application flow

## B. AppMeasure

AppMeasure                    1    IPFlow              NetFlow

IP                        NetFlow                        Cisco NetFlow

2    AnyTrace

3    AppFlow

FTP    TELNET    WWW

application    AppMeasure

AppMeasure      AppCollect

AppCollect

AppMeasure

AppMeasure

Data Mining

AppMeasure                              Monitor

**. AppMeasure**

**. AppMeasure**

.

.

AppMeasure                                                  AppMeasure

AppCollect                              AppMeasure      AppCollect

user level

kernel level     AppCollect

Pcap                                              Pcap

AppCollect        AppCollect

(datagram)                  AppMeasure

AppMeasure                                              AppMeasure

overhead

AppMeasure

AppMeasure

AppMeasure

## A. Implementation Issue

### A.1 . Design Objectives and Scope

Our objective is to put the class-based QoS routing protocol into practice, and we try to limit the implementation complexity. There are some important assumptions which affect the design choices. These assumptions include:

}     Support for on-demand path recalculation.

}     Exchange QoS parameter using the metric field of LSAs.

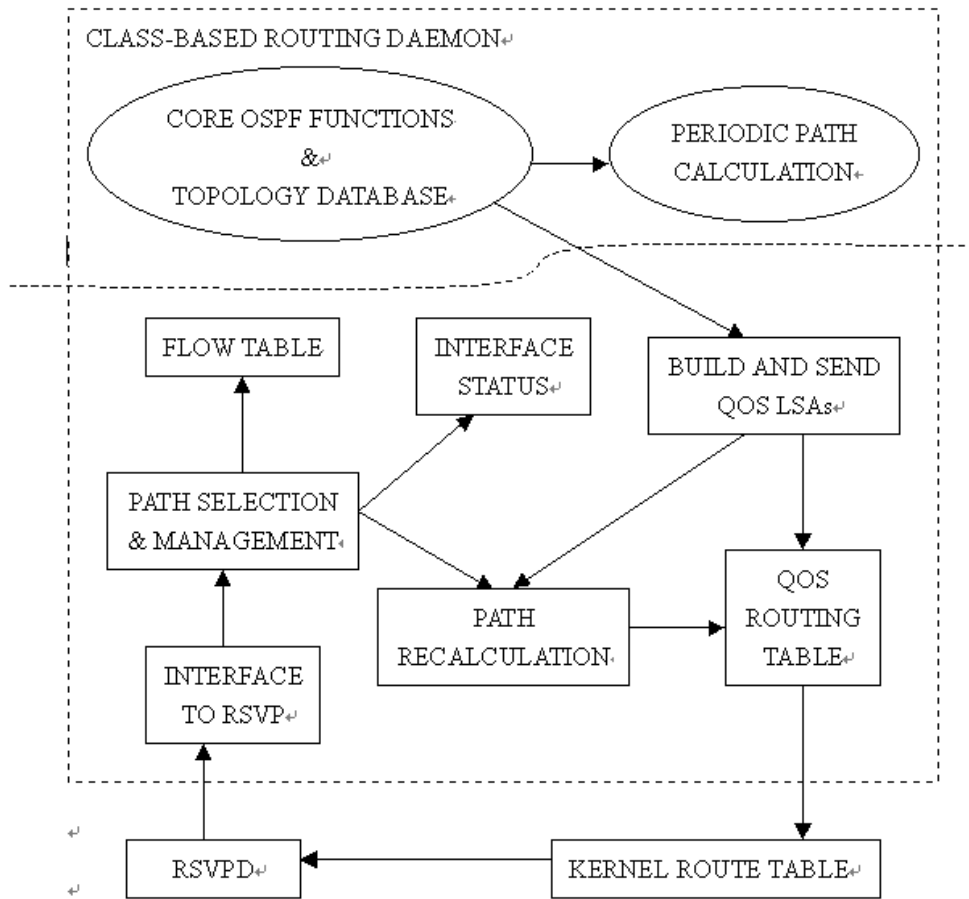}     Support per-class routing.

}     Interface to RSVP.

In addition, our implementation also relies on the following simplifying assumptions made in [23].

}    The scope of QoS route computation is limited in a single area.

}    All routers in the area run the class-based QoS route daemon.

}    All interfaces on a router are QoS capable.

}    Support hop-by-hop routing only.


### 3.2. Software Architecture

Figure 9 shows the architecture of the class-based QoS routing daemon. The modules of QoS routing consists of following modules:

}    **Interface status module** maintain parameters of network status from the original OSPF. The interface status records the residual bandwidth on the interface.

}    **QoS routing table module** records route cost and residual bandwidth along the path as QoS parameter in its route entry. the module stores the class-based routing table into kernel routing table using the data structure made for ECMP.

}    **QoS LSAs module** creates QoS LSAs by collecting status of all interfaces on the router and updates the QoS parameters, route cost and residual bandwidth, recorded in the routing table when a new QoS LSA arrives.

}    **Flow table (with traffic characteristics) module** records the QoS request of flows. The entry of the flow table contains flow informantion and traffic characteristics. The QoS requirements are part of RSVP messages.

}    **Path selection and management module** records the request of flows into the flow table and trigger path recalculation if necessary.

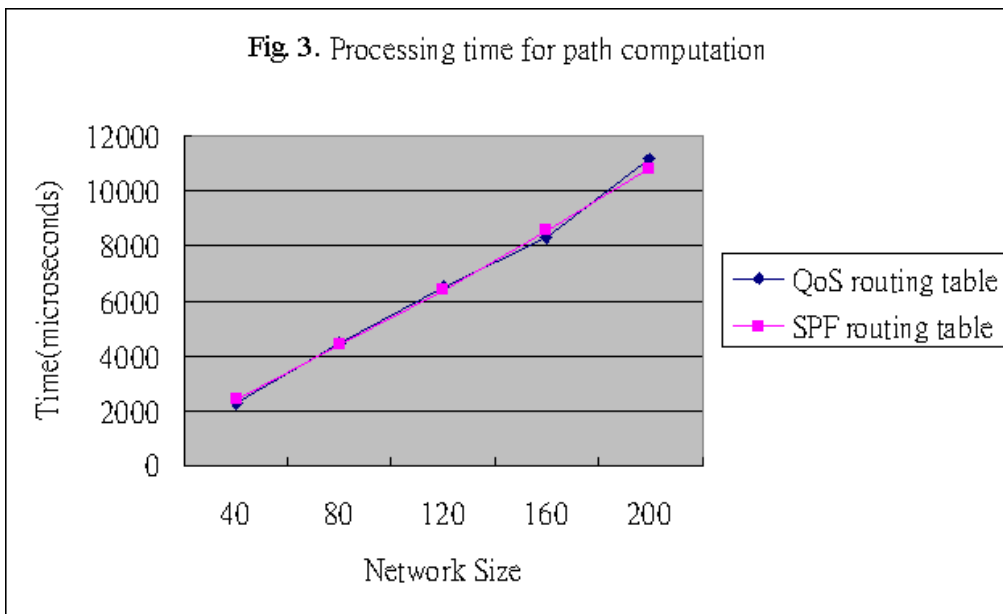}    **Path Recalculation module** finds the least cost route according to Dijkstra algorithm.

## B. Performance Evaluation

### B.1. Methodology

In this section, we will evaluate the performance of our implemetation. We explore three different dimensions in our comparisons: a) processing cost, b) memory requirement, and c) message generation and reception cost.

We construct a network topology to accomplish the experiment. The topology can be expanded by repeating a basic building block. The basic building block consists of 4 routers and 5 transit networks and is shown in Figure 10. The N x N mesh topology is constructed by repeating the basic block along two dimensions. We can control the network size by simply changing N. We set up the mesh topology by creating pseudo router LSAs in our implementation.

Routers    Transit networks



Fig. 3. Processing time for path computation



**B.2. Stand-Alone Cost**

It is important to measure the cost of our implementation. Among the router cost, processing time and routing table size can be measured on a single router. The cost measured on a single router is called stand-alone cost. To measure the traffic amount of LSAs exchange, we need at least two routers and monitor the packets of LSAs on the link connected to the two routers.
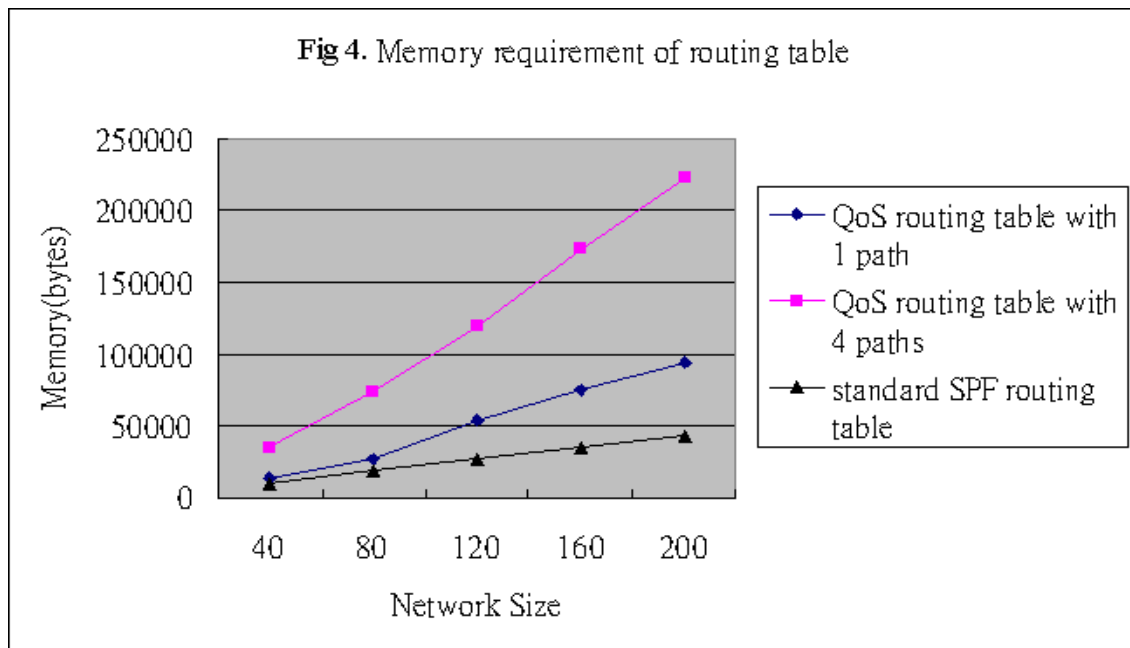
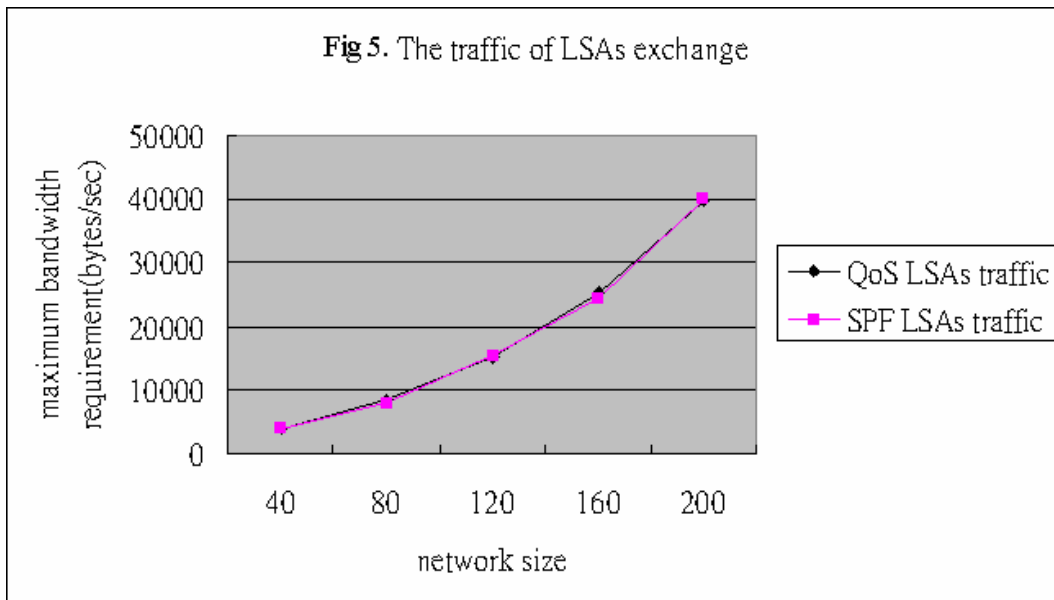B.2.1 Processing time of Routing Table Computation

Figure 11 shows the comparison of processing time used by QoS routing and standard SPF tree calculation. The standard SPF routing algorithm runs Dijkstra to find the SPF tree, and so does the

22

class-based QoS routing algorithm. The result provides that the sum of reciprocal bandwidth is good to be the cost function of the class-based QoS routing algorithm.

## B.2.2. Memory Requirements of the QoS Routing Table

Figure 12 shows the comparison of memory requirements of QoS routing table with full paths, QoS routing table within one path and standard SPF routing table. To provide accuracy of QoS routing, we maintain the lists of nexthops in each route entry. It results that the memory used to store the QoS routing table is much greater than the standard SPF routing table. It's necessary to record the nexthops to provide accuracy.



Fig 4. Memory requirement of routing table

Fig 5. The traffic of LSAs exchange

B.2.3. Link State Advertisements Generation and Receiption

As seen in figure 13, the curves of the QoS LSAs traffic and the standard SPF LSAs traffic are closed to each other. The major difference between the two kind of traffic is their frequency. We always monitor the QoS LSAs traffic on the link but the SPF LSAs traffic is seen for two or three times. Thus, the frequency of the QoS LSAs generation is much higher than the SPF LSAs. However, even the maximum LSAs traffic amount is still a little fragment of the whole bandwidth of a link. The influence of the LSAs traffic can be ignored.

In this work, we first explain the need of network processors for today's complex applications, and introduce the architecture and packet flow in IXP1200 shown in Fig. 14. Then we detail the mapping of DiffServ onto IXP1200, as shown in Fig. 15. There are two most important modules in DiffServ, classifier and scheduler, which are implemented with Multi-dimensional Range Matching [35] and Deficit Round Robin [36]. Finally we have external and internal benchmarks in order to find the bottlenecks in our implementation and possible design pitfalls of IXP1200.

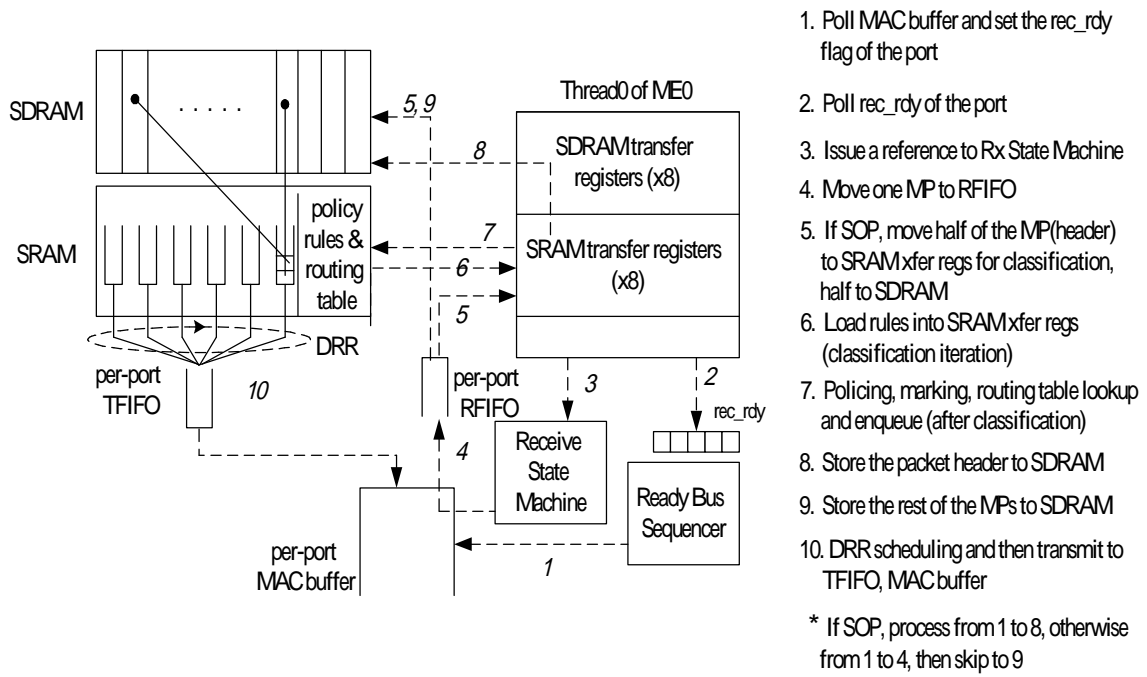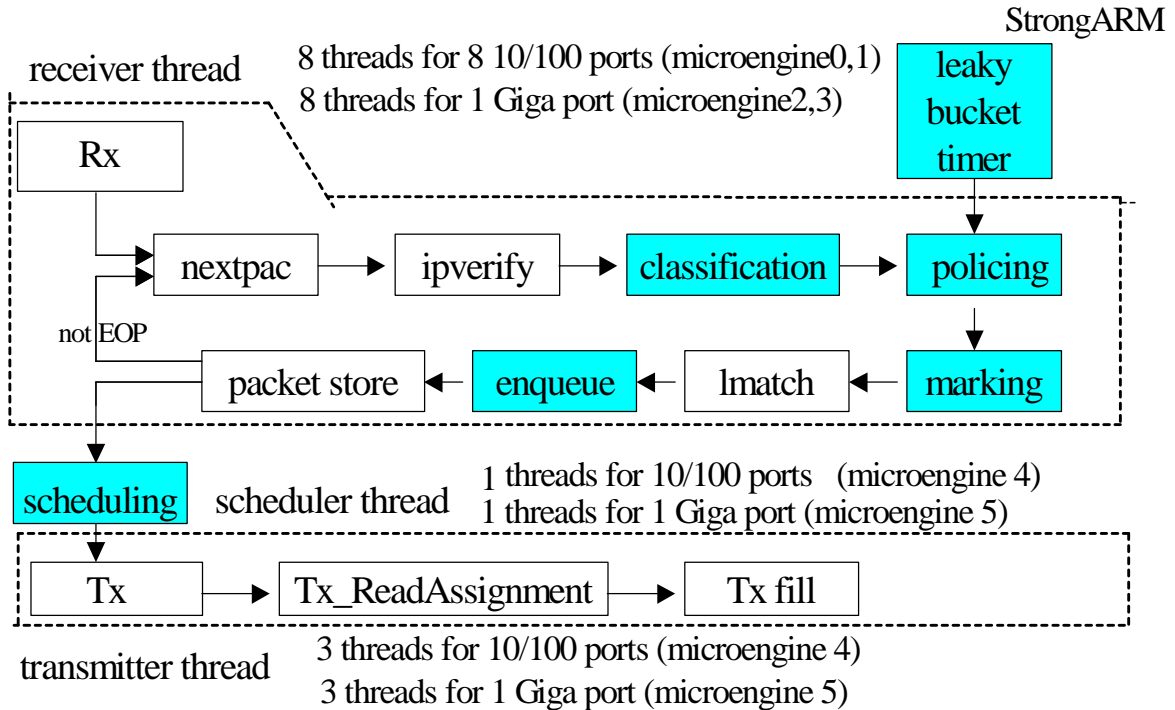24

Fig. 14. Detailed DiffServ packet flow in IXP1200.



Fig. 15. Data-plane architecture of DiffServ edge router over IXP1200.

The results of external benchmarks, as in Fig. 16, Fig. 17 and Fig. 18, have shown that our implementation can support well the PHBs in DiffServ at an aggregated throughput of 290Mbps. .We also identify the *MAC buffer overflow* which is described below. Fig. 19 shows

a diagram of packet reception. As we can see in Fig. 1, the rest of MPs, which are basic data units in IXP1200, are transferred from MAC buffer, RFIFO to SDRAM after the SOP (Start Of Packet) is classified. However, if SOP cannot be processed in time and the buffer is not large enough, the incoming MPs of the same packet could fill up the whole buffer and thus result in a packet drop, and then 100% packet loss.

Since both the *slow classification* and *small buffer* contribute to the MAC buffer overflow, we propose three solutions to avoid the two necessary conditions. They are (1). faster classification, (2). larger MAC buffer size, and (3). move the MPs into SDRAM before classification.
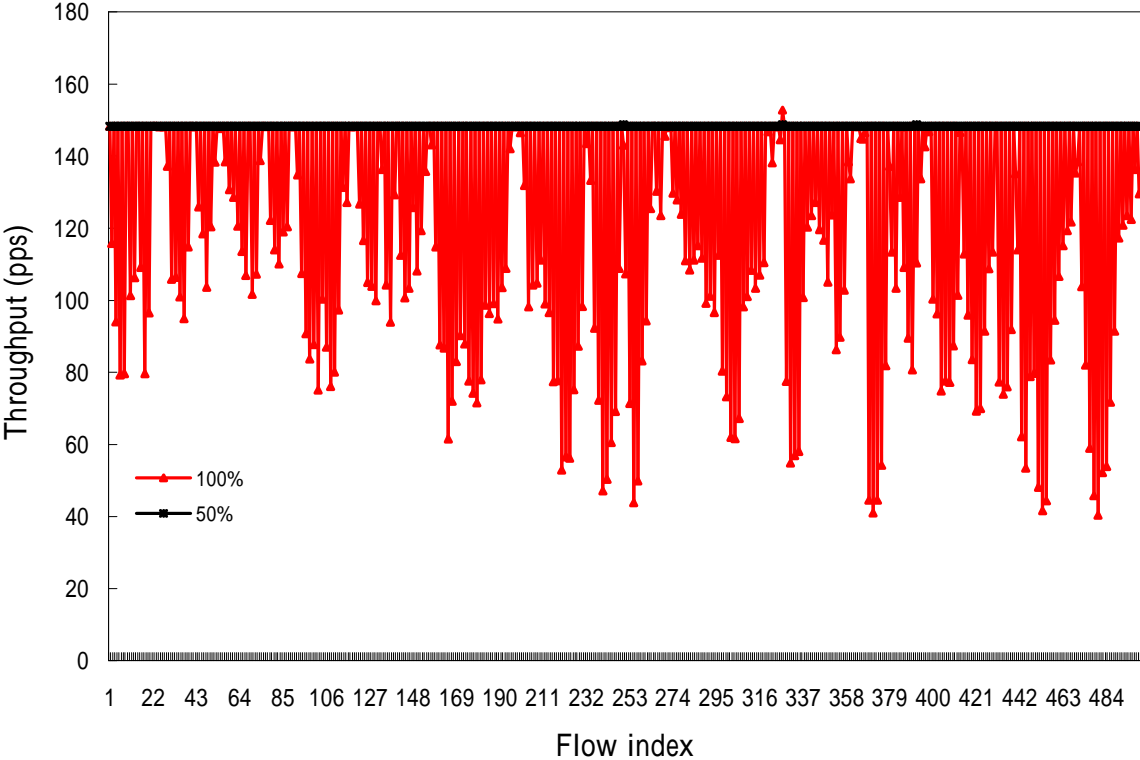


Fig. 16.    Flow fairness test (Len=64bytes, 500 flows, BW=74400/500=148pps, normal case).
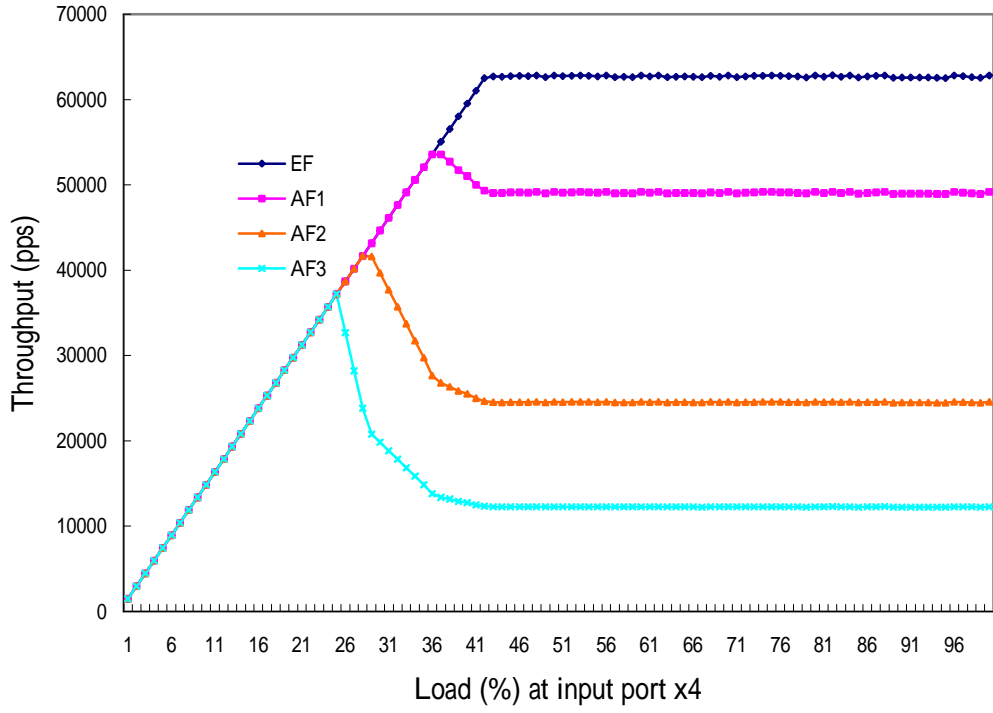
Fig. 17.    Priority and bandwidth control test (Len=64byte, EF=62500pps)
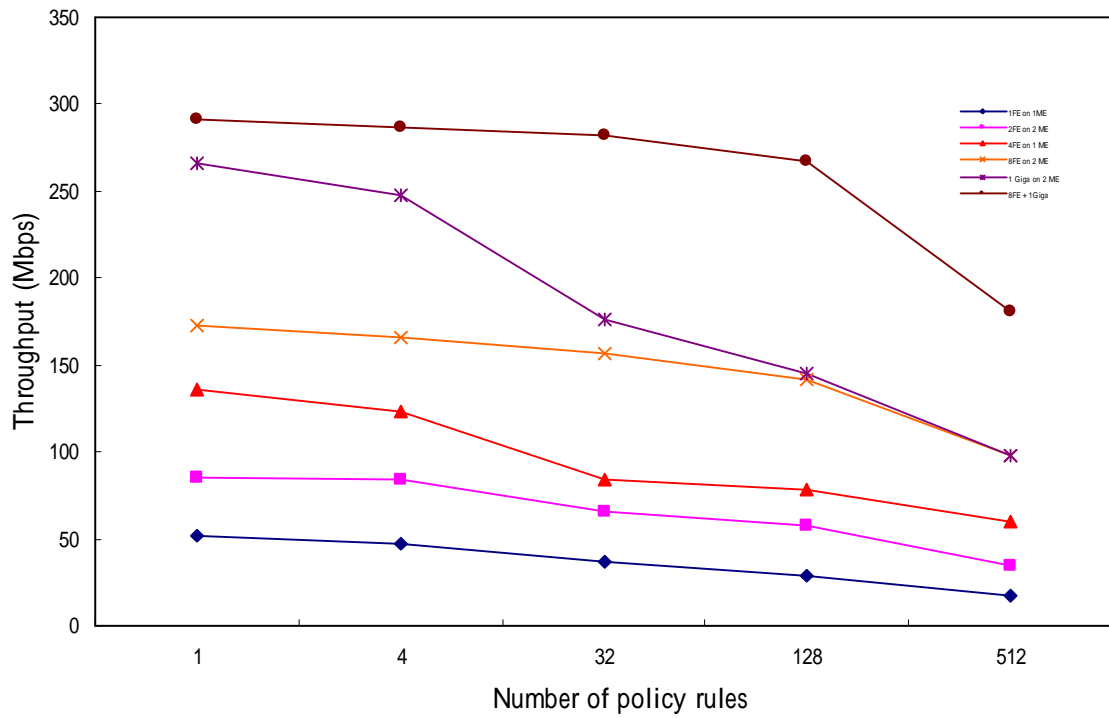


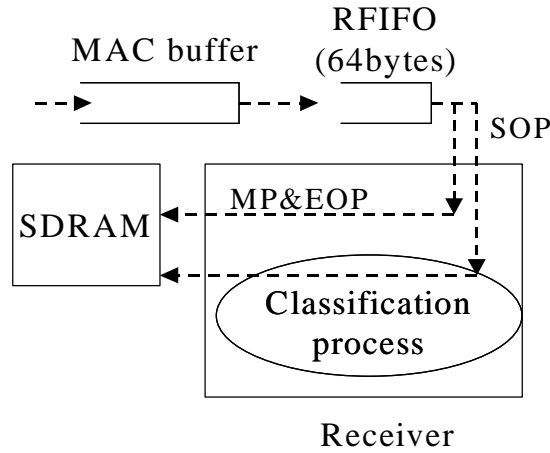Fig. 18    Aggregated throughput (Len=64bytes, worst case)

Fig. 19.   Receiving process of a packet

In both external and internal benchmarks, we identify the *double bottleneck* of both exclusive SRAM access and the lack of computing power in microengines inside the Range Matching DiffServ, as shown in Table 1. That is, the Range Matching DiffServ could still suffer from the other bottleneck after one of them is solved. Three methods are proposed to solve the bottleneck of SRAM accesses that leads to the low utilization of receiver microengines. First is to divide one large SRAM into many smaller *banks* at different interfaces. This could shorten the queuing delay of requests in the command queue if the requested addresses are in different memory banks. Second, we may adopt a new memory architecture, for example, RAMBUS DRAM (RDRAM) [37] in IQ2000 [38] that has a peak bandwidth of up to 1.6GBps which is two to three times of what SRAM supports. Third, an additional *cache* can be used to reduce the number of memory accesses because the traffic in the same time period usually shows locality in lookups of policy and routing tables.

While the SDRAM is the bottleneck in IP forwarding [39], we observe that the bottleneck may shift from one functional unit to another depending on the specific service, algorithm and the way input traffic is allocated to threads, as shown in Table 1. We also find that the SRAM bottleneck does not necessarily occur at 100% utilization, it could even occur at *55%* when the access is *bursty*.

Table 1. Bottlenecks in DiffServs of two algorithms

| Service or traffic allocation | Bottleneck |
|---|---|
| Linear search | SRAM |
| Range matching : | |
|   Single input port | SRAM |
|   8x100M input ports | ME |
|   1 gigabit port | ME |
|   8x100M and 1 gigabit | SRAM |

Gigabit

FFQ                                    SCFQ

                                        CIOQ                              LCF/MUF

                    CIOQ                    output queueing

buffer

                AppMeasure                        AppMeasure

                                                            GNU Zebra

            DiffServ                                    ASIC

LXP1200                                                1.8Gps

            DiffServ                200~300Mbps

[1.] Hui Ahang, "Service Discipline for Guaranteed Performance Service in Packet -Switching Networks," Proceedings of IEEE, vol. 83, no.10, Oct., 1995.

[2.] A.K. Parekh and R.G. Gallager, "A Generalized Processor Sharing Approach to Flow Control – The Single Node case," Proc. INFOCOM '92, vol. 2, May 1992, pp. 915-24.

[3.] J.C.R. Bennett and Hui Zhang, "WF$^2$Q: Worst-case Fair Weighted Fair Queueing," in Proc. IEEEINFOCOM '96, San Francisco, CA, Mar. 1996, pp. 120-128.

[4.] S.Jamaloddin Golestani, "A self-clocked Fair Queueing Scheme for Broadband Applications," in Proc. IEEE INFOCOM '94, Toronto, CA, June 1994, pp. 636-646.

[5.] Dimitrios Stiliadis and Anujan Varma, "Frame-based Fair Queueing: A New Traffic Scheduling algorithm for Packet-Switch Networks," Tech. Rep.UCSC-CRL-95-39, July 18, 1995.

[6.] Dimitrios Stiliadis and Anujan Varma, "Latency-rate servers: A general model for analysis of traffic scheduler algorithms," Tech. Rep. UCSC-CRL-95-38, U.C. Santa Cruz, Dept. of computer Engineering, July 1995.

[7.] David A. Patterson, John L. Hennessy, "Computer Organization & Design, The Hardware/Software Interface," 2nd Edition, Morgan Kaufmann, 1998.

[8.] A. Varma and D. Stiliadis, "Hardware Implementation of Fair Queuing Algorithms for Asynchronous Traffic Mode Networks," IEEE Communications Magazine, vol. 35, no. 12, Dec. 1997, pp. 54-68.

[9.] B. Prabhakar and N. McKeown, "On the Speedup Required for Combined Input and Output Queued Switching," Computer Systems Lab, Technical Report CSL-TR-97-738, Stanford University.

[10.]Karol, M. Hluchyj, and S. Morgan, "Input Versus Output Queueing on a Space Division Switch," IEEE Trans. Commun., vol. 35, pp. 1347-1763, Dec.

[11.]N. McKeown, V. Anantharam, and J. Walrand, "Achieving 100% Throughput in an Input-queued Switch," Proc. IEEE INFOCOM'96, pp. 296-302.

[12.]S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, W. Weiss, "An Architecture for Differentiated Services", IETF RFC 2475, Dec. 1998.

[13.]C.S. Chang and J.A. Thomas, "Effective Bandwidth in High-Speed Digital Networks," IEEE Journal on Selected Areas in Communications, Vol. 13, No. 6, pp. 1091-1100, August, 1995.

[14.] D. D. Clark, S. Shenker, L. Zhang, "Supporting Real-Time Applications in an Integrated Services Packet Network: Architecture and Mechanism", SIGCOMM'92, 1992.

[15.] S. Floyd, "Comments on Measurement- Based Admissions Control for Controlled- Load Services", Technical Report, 1996. http://www.aciri.org/floyd/admit.html

[16.] R. Guerin, H. Ahmadi, and M. Naghshineh, "Equivalent Capacity and Its Application to Bandwidth Allocation in High-Speed Networks", IEEE Journal on Selected Areas in Communications, Vol. 9, No. 7, pp.968-981, Sep. 1991.

[17.] S. Jamin, P. B. Danzig, S. J. Shenker, and L. Zhang, "A Measurement- Based Admission Control Algorithm for Integrated Service Packet Networks", IEEE/ACM Transactions on Networking, Vol. 5, No. 1, pp.56-70, Feb. 1997.

[18.] S. Shenker, C. Partridge, and R. Guerin, "Specification of Guaranteed Quality of Service", IETF RFC2212, Sep. 1997.

[19.] J. Wroclawski, "The Use of RSVP with IETF Integrated Services", IETF RFC 2210, Sep. 1997.

[20.] "Cisco IOS NetFlow Technical Documents", http://www.cisco.com/warp/public/732/Tech/netflow/netflow_techdoc.shtml

[21.] The RSVP project, http://www.isi.edu/div7/rsvp/

[22.] The Zebra project, http://www.zebra.org/

[23.] G.apostolopoulos, R. Guerin, and S. Kamat, "Implementation and Performance Measurements of QoS Routing Extensions to OSPF," Proceedings of IEEE INFOCOM'99, New York, NY, March 1999.

[24.] Yun-Wen Chen, Ren-Hung Hwang, and Ying-Dar Lin. "Multipath QoS Routing with Bandwidth Guarantee,"

[25.] RFC 2676, "QoS Routing Mechanisms and OSPF Extensions," Jan 1998.

[26.] RFC 2205, R.Braden Ed, L.Zhang, S.Berson, S.Herzog, and S.Jamin, "Resource ReSerVation Protocol (RSVP)," Sep. 1997.

[27.] RFC 2328, J.Moy, "OSPF Version 2," April 1998.

[28.] Internet Draft, R.Braden, D.Hoffman, "RAPI -- An RSVP Application Programming Interface Version 5," August 1998.

[29.] NetScreen Appliances, http://www.netscreen.com/international/products/appliances.html#ns5

[30.] Intel NetStructure VPN Gateway Family,

http://www.intel.com/network/idc/products/vpn_gateway.htm.

[31.]Intel Layer 3 Switching, "High speed LAN routing in an affordable switching solution,"
http://www.intel.com/network/tech_brief/layer_3_switching.htm.

[32.]eQoS Solutions for Service Providers using Riverstone Networks' Switch Routers,
http://www.riverstonenet.com/technology/eqos.shtml.

[33.]Technical report on Hardware-Based Layer5 load balancer,
http://www.nwfusion.com/research/2000/0501feat2.html.

[34.]Intel Electronic Design Kit,
http://developer.intel.com/design/edk/product/ixp1200_edk.htm.

[35.]T.V. Lakshman, and D. Stiliadis, "High-Speed Policy-based Packet Forwarding Using
Efficient Multi-dimensional Range Matching," ACM SIGCOMM'98.

[36.]M. Shreedhar, and G. Varghese, "Efficient Fair Queuing Using Deficit Round-Robin,"
IEEE/ACM Transactions on Networking, June 1996, vol. 4, no. 3, pp. 375-385.

[37.]Data Sheets of RDRAM, http://www.rambus.com/developer/support_rdram.html

[38.]IQ2000 Network Processor, VITESSE Corp,
http://www.vitesse.com/products/categories.cfm?family_id=5&category_id=16

[39.]T. Spalink, S. Karlin, L. Peterson, "Evaluating Network Processors in IP Forwarding,"
*Technical Report TR-626-00*, Computer Science, Princeton University, Nov. 1999.