（II）

89

2

# LIST OF TABLE

**Table**                                                                                                            **page**

# LIST OF FIGURE

# 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 consisting of five sub-projects constantly improves the results in the first year and further investigates the following key technologies:

A. High-capacity (Gigabit) routers: Sub-project I investigated two key technologies of developing high-capacity routers: switch architecture and longest prefix matching with hardware (or hardware routing). We employed space-division architecture, such as the crossbar, to build a real large-capacity router. In addition, we designed queue management and fast scheduling algorithms to (partially) remove head-of-line blocking in an attempt to improve the router's throughput. This study mainly focuses on designing a switch for variable-length packets to reduce idle time on output port in achieving better performance in terms of throughput and packet delay. Through simulation, the result showed that switches designed for fixed-length packets are inadequate for network environment of variable-length packets on the fly. Besides, to efficiently classify arriving packets, this study proposed a three-phase packet classification algorithm based on destination/source IP addresses, destination/source port numbers and protocol ID fields and designed efficient hardware routing schemes to speed up routing decision.

B. Admission control/QoS scheduling: QoS scheduling for broadband Internet is aimed to provide bounded delay and fairness while retaining a minimum of computational complexity. Prevailing weight-based scheduling disciplines advocate the use of multiple

queues and engage in timestamp computation. These disciplines achieve either superior QoS performance at the expense of higher complexity or degraded performance in return for computational simplicity. In the sub-project IV of this year, we have designed a weight-based Versatile QoS Scheduler (VQS) and its feasible VLSI hardware implementation architecture. VQS is capable of being implemented in various network elements in broadband Internet facilitating proper trade-off balance between performance and complexity. Taking advantage of simpler single-queue management and lack of timestamp computation, VQS governs packet insertion in a shared data structure comprising a sequence of fixed-size *windows* based on weights. Within a given widow, the maximum number of packets from a session is proportional to the session weight and the Window Size (*WS*). Simulation results demonstrate that, applying a smaller *WS* for high-power network elements, VQS performs as superior as WF$^2$Q with respect to throughput fairness, mean delay, and worst-case delay fairness. Moreover, compatible to WF$^2$Q, VQS outperforms WFQ with respect to 99-percentile delay bound and jitter in the presence of traffic burstiness.

C. Traffic measurement and statistics collection/admission control: Sub-project III 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 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." Simulation results show that in order to maintain a maximum packet loss rate, it is important that the system avoids operating at heavy loads, i.e. high link utilization. In a link sharing system, dynamic bandwidth allocation based on input loads can effectively avoid congestion. Furthermore, when combined with active queue management, it can further accommodate transient traffic bursts for non-self-similar traffic.

D. QoS routing: For reducing network information to achieve scalability in large ATM

networks, ATM Private Network-to-Network Interface (PNNI) adopts hierarchical routing. Consequently, although routing complexity is significantly reduced, numerous issues in PNNI routing require further study to achieve more efficient, accurate, scalable, and QoS-aware routing. In this year, we proposed several methods to achieve efficient, scalable, and QoS-aware ATM PNNI routing. First, an efficient aggregation scheme, referred to as Asymmetric Simple, is proposed. The aggregated routing information includes available bandwidth, delay and cost. Second, two approaches for defining link costs are investigated, namely, the Markov Decision Process (MDP) approach and the Competitive On-Line (COL) routing approach, and these are compared with the Widest Path (WP) approach. Third, a dynamic update policy, referred to as the dynamic cost-based update (DCU) policy, is proposed to improve the accuracy of the aggregated information and the performance of hierarchical routing, while decreasing the frequency of re-aggregation and information distribution. Finally, we proposed CIS (Crankback Information Stack) and CT (Cost Threshold) approaches to reduce crankback overhead. Simulation results demonstrate that the proposed Asymmetric Simple aggregation scheme yields very good network utilization while significantly reducing the amount of advertised information. Between these two links cost functions, the MDP approach provides a systematic method of defining call admission function and yields better network utilization than the COL approach. The proposed DCU policy also yields an enhanced network utilization while significantly reducing the frequency of re-aggregation. Meanwhile, the proposed CIS and CT approaches reduce crankback overhead significantly. Especially, the combination of CIS and CT approach achieves further improvement.

E. RSVP (ReSource reserVation Protocol) to PHB (Per-Hop Behavior) mapping: Subproject V investigates how granularity of routing decision significantly affects the scalability and blocking performance of QoS routing based on QoS routing extensions to OSPF. Three mechanisms, overflowed-cache, two-phase routing, and per-class routing mark, are also proposed to achieve computational and storage scalability as well as low blocking probability

in wire-speed packet-switching networks. Simulation results of various routing and forwarding granularities, including per-destination, per-pair, per-flow, per-pair/ overflowed-cache, per-pair/two-phase, per-pair/class, indicate that the proposed mechanisms can significantly lower blocking probability, increase fairness, as well as lower storage and computational overhead. Also, two or three classes are sufficient for per-class routing which is suitable for DiffServ core networks. Comparing flow driven mechanisms like per-flow and per-pair/overflowed-cache with topology driven mechanisms like per-destination, per-pair, per-pair/two-phase, and per-pair/class reveals that the former usually perform better in blocking probability, fairness, and state accuracy, while the latter result in less overhead.

A.          Gigabit

port trunking

crossbar

hashing   search tree   linear search

filter

B.          QoS  (Quality-of-

Service)                              Prevailing

weight-based          multiple  queues   engage   timestamp

QoS

weight-based Versatile QoS Scheduler (VQS)

VLSI                              VQS

VQS          single-queue      timestamp

weight

windows          window          session

session   weight   window                    window

high-power       VQS   throughput fairness  mean delay   worst-case delay

10

fairness   WF²Q   WF²Q   traffic burstiness

VQS   WFQ   99%   delay bound   jitter

C.

IETF

D.   ATM Private Network-to-Network Interface

(PNNI)   ATM

PNNI

PNNI

Asymmetric Simple

(MDP)   (COL)

(WP)   (DCU)

CIS   CT   crankback overhead

Asymmetric Simple

MDP   call admission   calld   COL   WP

MDP   DCU

CIS   CT

crankback overhead

--

flow driven                          per-flow    per-pair/overflowed-cache

                    data driven                          per-destination    per-pair    per-pair/two-

phase    per-pair/class                overhead

                                                    Gigabit

Circuit

switching                          Packet switching

Quality  of  Service

Fast  Ethernet

100Mbps                     Ethernet Switch

:

port trunking

crossbar

hashing  search tree  linear search

filter

QoS  (Quality-of-Service)

weight-based Versatile QoS Scheduler (VQS)                VLSI

VQS

(                                        )

QoS

metrics  routing algorithm              RSVP                                    re-

routing

PNNI                                                                Asymmetric

Simple                                   (MDP)

(COL)              (WP)

(DCU)

CIS     CT                crankback overhead

   (                                                    )

Internet                                        Quality
of Service                    QoS control



                              crossbar

crossbar                              Switch scheduling algorithm

                        PIM  Parallel Iteration Matching    RRM  Round-robin

matching      SLIP



                 hashing   search tree   linear search                    filter



Scheduling disciplines proposed in the literature have been either single-queue or

multiple-queue-based. Single-queue- based disciplines advocate the maintenance of a single

shared queue for each output link. Different-session packets destined to the same output link

are inserted in the shared queue in accordance with, for instance, the deadlines or priorities of

packets. Packets are then transmitted in a FIFO manner. Consequently, scheduling complexity completely resides in the enqueueing process. Examples of this class include Earliest Deadline First (EDF), Threshold Based Priority (TBP), and Precedence with Partial Push-Out (PPP). The EDF discipline was shown to successfully support tight delay bound. However, it undergoes two major limitations- a priori deadline assumption and high implementation complexity due to packet sorting. Although TBP and PPP were justified effective for switches supporting two priorities, they fail to provide bounded delay and throughput fairness in the presence of malicious sessions.

Multiple-queue-based disciplines, on the other hand, adopt multiple queues maintained at each output link, one for each session. Packets arriving from different sessions are simply placed at the end of their corresponding queues. Scheduling complexity in this class resides in the dequeueing process instead. Prevailing disciplines in this class, which are weight-based, include Weighted Fair Queueing (WFQ), Worst-case Fair Weighted Fair Queueing (WF$^2$Q), Self-Clocked Fair Queueing (SCFQ), and Frame-based Fair Queueing (FFQ).

In this project, we aim to design a weight-based, highly versatile QoS scheduler, referred to as VQS, capable of being implemented in diverse network elements facilitating proper trade-off balance between performance and complexity. Taking advantage of simpler single-queue management and lack of timestamp computation, VQS governs the insertion of packets belonging to the same output link in a shared data structure comprising a sequence of fixed-size *windows*. Within a given widow, the maximum number of packets from a session is proportional to the session weight and the Window Size (*WS*). Packets being placed at the same window are transmitted on a FIFO basis, limiting short-term unfairness to within a window interval. Packets being arranged outside of the window trigger new windows to be activated, enforcing weight-proportional service to be exerted.

IETF

(Integrated  Service)                              (Differentiated

Service)                                        (absolute)           (delay

bound)                          (queueing)

(Guaranteed  Service)

*vat   nv   vic*

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

(measurement-based)

(admission control)                          (parameter- based)          equivalent

capacity        effective  bandwidth              (a  priori)

(measurement  window)

(predictive service)                    equivalent

capacity                                          (average arrival

rate)        Hoeffding bound       equivalent capacity

The ATM PNNI standard adopts a source-based hierarchical routing for supporting scalability and security in a large network. The main advantage of the hierarchical routing is reducing large communication overhead while achieving efficient routing. The scalability and performance of hierarchical networks depend on various design schemes, such as the aggregation scheme for aggregating routing information, the cost functions for defining link and path costs, and update policies for advertising the aggregated information. Nevertheless, how to design these schemes remains an open issue. In this project, we propose several solutions, including the efficient Asymmetric Simple aggregation scheme, the QoS-capable COL and MDP link cost approaches, the dynamic cost-based update policy, and CIS and CT for reducing crankback overhead to achieve an efficient hierarchical QoS routing in large ATM networks. These methods are briefly described below.

## A. Hierarical Routing

In this project, we study the PNNI standard, source-based hierarchical routing in ATM networks. The source-based hierarchical routing problem can be decomposed into two issues: how to aggregate routing information and how to perform hierarchical routing. For routing information aggregation, Iwata *et al.* proposed two aggregation schemes, star and simple node, with three aggregation versions, aggressive, conservative and simple no-aggregation. These schemes transform a non-linear programming problem into a linear problem for the corresponding QoS parameters. Meanwhile, Lee proposed a spanning tree aggregation

scheme, and Awerbuch *et al.* compared the performances of several aggregation schemes, including a star with radius equal to half the cost of the network diameter (DIA), a star with radius equal to half the average cost between border nodes (AVE), Minimum Spanning Tree (MST), Random Spanning Tree (RST), and t-spanner. Our earlier work proposed a novel aggregation scheme, called Asymmetric Simple and compared it with two existing aggregation schemes (Simple Node and Full-Mesh) using various performance metrics, such as representation size and representation accuracy for routing information, and network revenues.

For hierarchical routing, Guo *et al.* applied probabilistically routing to hierarchical networks. Meanwhile, Mieghem presented the unicast hierarchical routing based on PNNI standard. Furthermore, , Montgomery *et al.* and Xie *et al.* applied the theory of reduced load approximation to analyze the blocking probability of PNNI hierarchical networks. Finally, Hao *et al.* investigated the call rejection probability for routing with crankback.

Our numerical results first demonstrate that an effective aggregation scheme reduces overhead for call set up while yielding high traffic throughput. The proposed Asymmetric Simple aggregation scheme can yield competitive performance compared to the Full Mesh aggregation scheme.

## B. Cost Functions for Hierarchical QoS Routing

In this project, we define the objective of hierarchical QoS routing as to maximize network revenue under the constraint that each established connection is guaranteed with certain QoS requirements. This optimization problem is generally re-formulated as, for each new arriving connection, to find the path with minimum cost while satisfying certain QoS requirements. Two issues can be identified that differ from QoS routing in flat networks. First, owing to the inaccuracy of aggregated information, a chosen hierarchical path may not satisfy the end-to-end QoS requirement. In this case, the crankbank scheme can be employed to

reroute the connection to an alternate path. The crankback scheme is further discussed in chapter 2. The second issue is how to define the link cost function and aggregate link cost functions in a peer group. A good link cost function should comprise two properties; maximizing network revenue by minimizing path cost, and providing a systematic call admission function.

In this project, we study two approaches for defining link cost functions in a hierarchical network, namely the Markov Decision Process (MDP) approach and the Competitive On-Line (COL) approach. The theory of Markov decision process is a pledge method in a lot of network-related issues. Various network control schemes have been developed based on the Markov decision process. For example, many MDP-based routing algorithms, which compute the link cost based on the MDP theory, have been proposed and demonstrated to perform very well. However, Gawlick, *et al.*, proposed an on-line optimal routing algorithm, referred to as the Competitive On-Line (COL) algorithm. This approach defines the link cost function as an exponential function of the residual bandwidth. They have shown good routing performance based on this cost function. We proposed MDP-based cost function for hierarchical QoS routing.

In this project, we compare the performance of the MDP-based and the COL-based cost functions with the Widest Path (WP) approach, which routes an incoming connection to the path with maximum residual bandwidth. The residual bandwidth of a path is defined as the minimum of the residual bandwidth of all links on the path. The simulation results show that the MDP and the COL approaches outperform the WP approach. The MDP approach yields the best network utilization. A further advantage of the MDP approach is that it provides a systematic call admission function.

C.  Update Poicies

On the other hand, PNNI adopts time-based update policy which is inadequate to cope

with dynamic network traffic. Furthermore, the accuracy of aggregated information is depended on the update interval, with a reduced update interval meaning more accurate aggregated information. However, in this situation the overhead of re-aggregation and information distribution increases. Awerbuch et al. proposed the logarithmic update approach, which is based on the residual bandwidth of a link, to reduce the computational overhead of re-aggregation.

An event-based update policy typically suffers from oscillation, which can be avoided by hystersis. The technique of hysteresis has been applied to various areas of high-speed networks. For example, Jong applied hysteresis to ATM rate control to enhance system stability, and Orda et al. proposed an adaptive virtual path allocation policy using hysteresis to prevent excessive processing of requests due to oscillations around thresholds. Meanwhile, Shun-Ping Chung and Jin-Chang Lee propose a dynamic reservation with hysteresis as CAC for cellular multiservice networks.

We proposed an event-driven update policy based on the link cost. Furthermore, to avoid oscillation, the hystersis technique was applied. The proposed policies are called the Dynamic Cost-based Update (DCU) policy and the DCU with hysteresis (DCUH) policy. The performance of the proposed policies is compared to other time-driven and event-driven update policies, including the PNNI time-based update approach (PNNIU), full update approach (FU), logarithm of residual bandwidth update approach (LRBU), and dynamic cost-based update policy without hysteresis. Simulation results show that the DCUH policy performs best among these update policies.

### D. Crankback Approaches

Alternate path routing approaches can be used to achieve lower connection blocking probabilities and higher network throughput in ATM networks. Conventional alternate path routing techniques tend to fall into two classes: progressive control. Hwang, R.H. *etc.* noted

that progressive control, as compared to the originating control, provides a higher blocking probability but smaller connection setup time The primary advantage of progressive control is the ability to provide fast connection setup, while the disadvantages of progressive control are the use of sub-optimal alternate path. Chung proposed an ICD (information about the crankback destination nodes) approach in a connection setup messages, hence they could predict the crankback destination node. Meanwhile, Spieqel *et al.* presented an approach to combine the features of progressive and originating control approaches, allowing for fast connection setup as well as near minimum cost paths. Felstaine *et al.* proposed to allocate a "quota" to the PGs along the message path and then to "sub-allocate" quota to the son PGs of these PGs.

The purpose of this section is to study mechanisms for reducing the crankback overhead in PNNI hierarchical networks. Two heuristics are proposed herein to reduce crankback overhead. First, if we can predict where to crank back, referred to as the destination node, the call setup message need not to crank back to the source node, thus the call set up overhead at the common nodes along the path can be reduced. Second, the aggregated path cost represents the expected cost for setting up the call on the path of O-D pair. Therefore, a call setup message on a path with high path cost is likely to be blocked. Consequently, avoiding setting up a call on a path with high path cost will reduce the crankback overhead. Based on these two heuristics, we proposed two approaches to reduce crankback overhead. The first approach adds additional information, referred to as CIS, to predict the crankback destination node. In this approach, we keep CIS in the call setup message while crossing the PG. When call setup message encounters a block link, the intermediate node using CIS to find the crankback destination node. Hence, the call setup message needs not to crank back to the ingress node. The second approach uses aggregated path cost, referred to as CT, to determine whether call setup on alternate paths should be tried. In PNNI hierarchical network, paths with smaller hierarchical path cost are tried first. If the aggregated path cost information is accurate and a

call setup message failed on the path with minimum cost, then call set up on alternate paths with higher path cost will likely to be failed. Simulation results show that both of these two approaches reduce the crankback overhead significantly. Furthermore, when we combine CIS with CT, referred to as CIS_CT, the crankback overhead can be reduced more conspicuously.

To achieve both low blocking probability and high scalability, this study proposes three QoS routing extensions to OSPF, *overflowed-cache*, *two-phase routing*, and *per-class routing mark*. The *overflowed-cache* mechanism divides the packet-forwarding cache into a per-pair cache (*P-cache*) and an overflowed per-flow cache (*O-cache*). The flows that the P-cache indicated paths cannot satisfy with the required QoS are routed individually and their forwarding decisions are *overflowed* into the O-cache. The *two-phase routing* reserves a block of bandwidth that exceeds the bandwidth requirement of a flow when the first flow between an S-D pair is established. The *per-class routing* aggregates QoS flows into a number of classes via a marking technique. Therefore, flows with the same *mark* for the specific S-D pair are routed on the same path. This routing mechanism is suitable for DiffServ networks, where packets are marked at edge routers and fast-forwarded in the core network.

# Gigabit

## A.                                    Crossbar

3.1

Compact PCI



**Backplane**

**Crossbar
Switch
Fabric**

**Bus**

**Interface cards**

**Memory**

**Network
Interface**

**Port
Controller**

**Bus
Interface**

3.1

## A.1

3.2

I/O Port Controller        3.3

memory                                    P        P_n in & P_n out

3

(RSP)

MC                                                                                          CTRL

```
        ┌─────────┐
        │   RAM   │
        └────┬────┘
             ↕
        ┌─────────┐  P_0 in
        │         │  P_0 out
        │  Port   │◄── RSP
        │Controller│◄── MC         P_n in : input port of n th
        │         │◄── CTRL        port controller
        └─────────┘
             ⋮                     P_n out : output port of n th
        ┌─────────┐                port controller
        │   RAM   │
        └────┬────┘  ┌──────────┐
             ↕       │ Crossbar │  RSP : response channel
        ┌─────────┐  │  Switch  │
        │         │  │ Element  │  MC : multicast channel
        │  Port   │  └──────────┘
        │Controller│  P_N in       CTRL : control channel
        │         │  P_N out
        └─────────┘◄── RSP
                  ◄── MC
                  ◄── CTRL
        ┌─────┐    ┌──────┐
        │ CPU │◄──►│ DRAM │
        └─────┘    └──────┘
```

3.2

## A.2

Fast Ethernet

3.3          Data  Path

PHY+MAC                                                    Table Lookup Engine

Queue  Manager

VOQ

26

VOQ　　　　　　　　　　　　　　　　　　　　　　　Crossbar Interface

　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　　P_n

out　　　　　　　　　　　　　　　　　　MC

　　　　　　　　　　　CPU Interface　　　　　　　　　　　　　routing table

　　(classification table)

　　　　　　　　　　　　　　RSP  Monitor　　　　RSP

　　　　　　　　　　　　　VOQ　　　　　　　　　　　P_n in



Table Memory

PHY — MAC ⋮ PHY — MAC — Data Path — Table Lookup Engine / Queue Manager — Crossbar Interface / RSP Monitor — P_n in / RSP / P_n out / MC — CPU Interface — CTRL

Packet Memory

3.3

## A.3

　　　　　　　　　　(Command Decoder)　　　　　(Arbitration Logic)

　　　Response encoder　　　　　　　Status Monitor　　　　　　Request Queue

　　　　　　Crossbar matrix　　　　　　(Matrix controller)　　　　3.4

VOQ

BUS1

BUS2

RSP

BUS3

(T+I)

RSP

## A.4

SLIP

3.4    a        b                                s_size=16   MaxPktSize=1500

MinPktSize=64   PSN=5000        CellSize= 64

## B.

## B.1

■            Hashing

bits              filter                          104    bit      $m$

bit            filters        space        $2^m$

$N/2^m$    filter        $N$    filter                                                              filter

                        Source IP address        Destination IP address                bits            bit

                                                a
                                                b
        3.4  SLIP

                a                    b

        UDP    TCP    protocol            7    bits        Hashing                        filter

space

        ■                Search Tree

filter    IP prefix        search tree                Source IP address> Destination Ip address>Protocol

$m$ bits              $2^m$-ary  Search  Tree

search tree                                              filter

search tree

■              Conflict Check and Sorting

filters         cost                        search

IP address        longest prefix matching     mask        filter              priority

mask          filter              match          conflict  problem

conflict                          conflict check            1)    mask

Filter $i$                  £ mask        Filter $j$                  Filter $i$  2)    mask

Filter $i$          > mask        Filter $j$                          costs

conflict check        filters' cost                        filter

bubble sort

## B.2

1        Hashing              bits              search tree        root

2        Source IP   Destination IP              prefix                      filter

3      Linear Search                      filter                filter      best
matching filter

## B.3

filter      $n$  hashing key    $b$ bits   $2^m$-array search tree    depth    $h(n)$

filter      $\dfrac{n}{2^{mh\left(n/2^b\right)}}$      search      filter

$$O\!\left(h\!\left(\frac{n}{2^b}\right)\right)+O\!\left(\frac{n}{2^{mh\left(n/2^b\right)}}\right)$$

$2^m$-ary search tree                IP prefix                depth=$K$
30

prefix address mask $m \cdot K <$ *Address Mask* $< m \cdot (K+1)$ depth=$K+1$

$[(K+1)\cdot m\text{-}Address\ Mask]\cdot 2$ filters depth=$P$ $P>K+1$

$$\left(2^{m}\right)^{P\cdot(K+1)}[(K+1)\cdot m\text{-}Address\ Mask]\cdot 2 = 2^{m(P-K-1)+1}[(K+1)\cdot m\text{-}Address\ Mask]$$

$n$ filters $i$ filters depth $K_i$ address mask

$aMask_i$ search

$$O\left(h(n/2^{b})\right) + O\left(\frac{n - i + \sum_{a=1}^{i} 2^{m\left(h(n)-K_i-1\right)+1}\left[(K_i+1)\cdot m - aMask_i\right]}{2^{mh(n/2^{b})}}\right)$$

**A.**

(traffic load)

oc3 link

(sampling interval) (aggregate)

( ) 3.5 (pdf,

probability density function)

3.6

(burst)

pdf

(bell shape)

**B.**

3.7

(        *T*)    *S*                                    *T/S*                              (0<

<1)                                              *AvgLoad*          *StdLoad*

   (1)                              *EC*                          *EC*        *EC*

                              (1)    *z*                      1-

                              (*T/S*)×

*EC*                *EC*                              (*T/S*)×

   *EC*                                                          *MaxLoad*

      *EC*                  (*MaxLoad-EC*)×*S*

   (drop)                                              (2)

$$EC = AvgLoad + z \cdot StdevLoad \qquad (1)$$

$$BufferSize = \frac{T}{S} \cdot u \cdot [(MaxLoad - EC) \cdot S]$$

$$= T \cdot u \cdot (MaxLoad - EC) \qquad (2)$$

   (self-similar)

         0



### 19990825_0200_0.100ms

3.5. 1999/8/25        2:00                  100ms

## 19990825_0200_0



3.6. 1999/8/25　　2:00　　　　　　　(1sec, 100ms, 20ms)



3.7

## A. Hierarchical Routing

In this project, we study the PNNI standard, source-based hierarchical routing in ATM networks. The source-based hierarchical routing problem can be decomposed into two issues: how to aggregate routing information and how to perform hierarchical routing. For routing

information aggregation, Iwata et al. proposed two aggregation schemes, star and simple node, with three aggregation versions, aggressive, conservative and simple no-aggregation. These schemes transform a non-linear programming problem into a linear problem for the corresponding QoS parameters. Meanwhile, Lee proposed a spanning tree aggregation scheme, and Awerbuch et al. compared the performances of several aggregation schemes, including a star with radius equal to half the cost of the network diameter (DIA), a star with radius equal to half the average cost between border nodes (AVE), Minimum Spanning Tree (MST), Random Spanning Tree (RST), and t-spanner. Our earlier work proposed a novel aggregation scheme, called Asymmetric Simple and compared it with two existing aggregation schemes (Simple Node and Full-Mesh) using various performance metrics, such as representation size and representation accuracy for routing information, and network revenues.

For hierarchical routing, Guo et al. applied probabilistically routing to hierarchical networks. Meanwhile, Mieghem presented the unicast hierarchical routing based on PNNI standard. Furthermore, Montgomery et al. and Xie et al. applied the theory of reduced load approximation to analyze the blocking probability of PNNI hierarchical networks. Finally, Hao et al. investigated the call rejection probability for routing with crankback.

Our numerical results first demonstrate that an effective aggregation scheme reduces overhead for call set up while yielding high traffic throughput. The proposed Asymmetric Simple aggregation scheme can yield competitive performance compared to the Full Mesh aggregation scheme.

## B. Cost Functions for Hierarchical QoS Routing

In this project, we define the objective of hierarchical QoS routing as to maximize network revenue under the constraint that each established connection is guaranteed with certain QoS requirements. This optimization problem is generally re-formulated as, for each new arriving connection, to find the path with minimum cost while satisfying certain QoS requirements. Two issues can be identified that differ from QoS routing in flat networks. First,

owing to the inaccuracy of aggregated information, a chosen hierarchical path may not satisfy the end-to-end QoS requirement. In this case, the crankbank scheme can be employed to reroute the connection to an alternate path. The crankback scheme is further discussed in section 2. The second issue is how to define the link cost function and aggregate link cost functions in a peer group. A good link cost function should comprise two properties; maximizing network revenue by minimizing path cost, and providing a systematic call admission function.

In this project, we study two approaches for defining link cost functions in a hierarchical network, namely the Markov Decision Process (MDP) approach and the Competitive On-Line (COL) approach. The theory of Markov decision process is a pledge method in a lot of network-related issues. Various network control schemes have been developed based on the Markov decision process. For example, many MDP-based routing algorithms, which compute the link cost based on the MDP theory, have been proposed and demonstrated to perform very well. In this project, we adopt the idea for its simplicity and efficiency. However, Gawlick, et al., proposed an on-line optimal routing algorithm, referred to as the Competitive On-Line (COL) algorithm. This approach defines the link cost function as an exponential function of the residual bandwidth. They have shown good routing performance based on this cost function. We proposed MDP-based cost function for hierarchical QoS routing.

In this project, we compare the performance of the MDP-based and the COL-based cost functions with the Widest Path (WP) approach, which routes an incoming connection to the path with maximum residual bandwidth. The residual bandwidth of a path is defined as the minimum of the residual bandwidth of all links on the path. The simulation results show that the MDP and the COL approaches outperform the WP approach. The MDP approach yields the best network utilization. A further advantage of the MDP approach is that it provides a systematic call admission function

## C. Update Poicies

On the other hand, PNNI adopts time-based update policy which is inadequate to cope with dynamic network traffic. Furthermore, the accuracy of aggregated information is depended on the update interval, with a reduced update interval meaning more accurate aggregated information. However, in this situation the overhead of re-aggregation and information distribution increases. Awerbuch et al. proposed the logarithmic update approach, which is based on the residual bandwidth of a link, to reduce the computational overhead of re-aggregation.

An event-based update policy typically suffers from oscillation, which can be avoided by hystersis. The technique of hysteresis has been applied to various areas of high-speed networks. For example, Jong applied hysteresis to ATM rate control to enhance system stability, and Orda et al. proposed an adaptive virtual path allocation policy using hysteresis to prevent excessive processing of requests due to oscillations around thresholds.

We proposed an event-driven update policy based on the link cost. Furthermore, to avoid oscillation, the hystersis technique was applied. The proposed policies are called the Dynamic Cost-based Update (DCU) policy and the DCU with hysteresis (DCUH) policy. The performance of the proposed policies is compared to other time-driven and event-driven update policies, including the PNNI time-based update approach (PNNIU), full update approach (FU), logarithm of residual bandwidth update approach (LRBU), and dynamic cost-based update policy without hysteresis. Simulation results show that the DCUH policy performs best among these update policies.

## D. Crankback Approaches

We further study mechanisms for reducing the crankback overhead in PNNI hierarchical networks. Two heuristics are proposed to reduce crankback overhead. First, if we can predict where to crank back, referred to as the destination node, the call setup message need not to crank back to the source node, thus the call set up overhead at the common nodes along the path can be reduced. Second, the aggregated path cost represents the expected cost for setting

up the call on the path of O-D pair. Therefore, a call setup message on a path with high path cost is likely to be blocked. Consequently, avoiding setting up a call on a path with high path cost will reduce the crankback overhead. Based on these two heuristics, we proposed two approaches to reduce crankback overhead. The first approach adds an addition information, referred to as CIS, to predict the crankback destination node. In this approach, we keep CIS in the call setup message while crossing the PG. When call setup message encounters a block link, the intermediate node using CIS to find the crankback destination node. Hence, the call setup message needs not to crank back to the ingress node. The second approach uses aggregated path cost to determine whether call setup on alternate paths should be tried. In PNNI hierarchical network, paths with smaller hierarchical path cost are tried first. If the aggregated path cost information is accurate and a call setup message failed on the path with minimum cost, then call set up on alternate paths with higher path cost will likely to be failed.

## E. Simulation Results

This section evaluates the performance of hierarchical routing with different aggregation schemes and cost functions via simulations. Figure 3.8 shows the network topology of the simulations. The capacity of each link is 155 Mbps, except for link(A.3.2-B.1.3) and link(B.1.3-A.3.2), which are 622 Mbps. Each link has a delay of 1ms.

In the simulations, the network supports two classes of traffic. The class 1 traffic has a bandwidth requirement of $b_1=1$ Mbps, and an end-to-end delay bound of $d_1=15$ ms. Class 1 traffic is assumed to arrive at any O-D pair with the same arrival rate $\lambda_1^s = \lambda$ , while the class 2 traffic is assumed to have a bandwidth requirement of $b_2=5$ Mbps, an end-to-end delay bound of $d_2=15$ ms, and an arrival rate of $\lambda_2^s = \lambda$ / $b_2$ between each O-D pair. The average holding times of these two classes of traffic are normalized to unity.

The 95% confidence intervals of the simulation results shown in the following figures are obtained from 25 independent runs. For each run, the simulated time is 1100 units of mean

Figure 3.8 Network topology

call holding time. The initial 100 time units are estimated as the transient period and, thus, performance samples are discarded.

Figure 3.9 shows the fractional reward loss of calls under different arrival rates, when two aggregation schemes are used for the WP and the MDP and the COL link cost functions. The figure shows that the fractional reward loss increases as the arrival rate increases. Meanwhile, the fractional reward loss of the Asymmetric Simple scheme of the COL link cost is slightly higher than for the Full Mesh scheme. The fractional reward loss of Asymmetric Simple scheme of MDP link cost is almost the same as Full Mesh scheme. Meanwhile, the MDP approach yields a better performance than the COL in both the Asymmetric Simple and Full Mesh schemes. This superiority is because the MDP approach considers the arrival rate, and so can provide more accurate cost information and call admission function. However, the MDP and COL cost-based approaches significantly outperform the WP approach. The main reason is that the bandwidth-based approach only considers the maximum residual bandwidth

38

Figure 3.9 Fractional reward loss.

of bottleneck link among all candidate paths. However, the cost-based approaches first transfer the residual bandwidth of link into a reasonable cost and then route the minimum path cost.

Figure 3.10 shows the average number of crankbacks per connection request under different arrival rates, when two aggregation schemes are used. The figure shows that the average number of crankbacks increases as the arrival rate increases. The average number of crankbacks of the Asymmetric Simple of the COL link cost function is higher than for the Full Mesh scheme. Meanwhile, the average number of crankbacks of the Asymmetric Simple of MDP link cost function is almost the same as that of the Full Mesh scheme. The COL approach of all schemes requires less crankbacks than the MDP. The reasons for this phenomenon are that the cost function is less accurate and the call admission of COL approach is less conscientious than that of the MDP approach. Consequently, a call may be routed to non-optimal path that requires more network resources and thus increases call blocking probability. The call admission function rejects a call if the aggregated path cost exceeds the admission threshold, which also results in less paths (crankbacks) being tried and

Figure 3.10 Crankback per requested calls.

higher call blocking probability. However, the MDP approach provides more accurate aggregated cost and precise call admission policy because the link cost is computed based on Markov decision theory. This approach results in less calls being blocked due to call admission at the hierarchical routing procedure, but more blocking due to the call admission at the physical network level. Hence, the MDP approach results in more crankbacks. Additionally, the number of crankbacks of the WP approach is more than that of the COL and MDP approaches. Clearly, the reason is that the WP approach is not a good approach for hierarchical routing.

Figure 3.11 shows the fractional reward loss yielded by the five update policies under different arrival rates. As figure 3.11 illustrates, the FU policy yields the lowest fractional reward loss, because it provides the most accurate routing updates. However, the LRBU, DCU, and DCUH policies all yield very competitive fractional reward loss as compared to the FU policy. Figure 3.11 shows that the DCUH policy yields slightly lower fractional reward loss than the LRBU and DCU policies. Meanwhile, figure 3.11 also illustrates that PNNIU policies yield worse performance than the other four policies. The performance of PNNIU

Figure 3.11 Fractional reward loss of five update policies.

policy can be improved if the update interval reduces. For example, the PNNIU_1 policy outperforms the PNNIU_4 policy.

Figure 3.12 shows the average number of crankbacks per connection request under different arrival rates. Intuitively, the more accurate the aggregated routing information, the less the average number of crankbacks. Hence, the average number of crankbacks per connection request is an important performance metric for update policies. Figure 3.12 illustrates that FU, LRBU, DCU, and DCUH have almost the same number of crankbacks per connection while the PNNIU policies have a much higher number of crankbacks.

Figures 3.13 and 3.14 show the number of information re-aggregations and distributions per unit of time under different arrival rates. Intuitively, the FU policy should have the highest re-aggregation frequency, and the two figures confirm this. On the other hand, adjusting the update interval gives the PNNU policy the least update frequency. Comparing the LRBU, DCU, and DCUH policies reveals that the DCU and DCUH policies yield less overhead re-aggregation than the LRBU policy. However, the DCU policy with hysteresis does not significantly reduce the update frequency. One possible reason for this phenomenon is that the

41

Figure 3.12 Average number of crankback per connection request.



Figure 3.13 Average number of re-aggregation per unit of time.

DCU policy does not suffer much from oscillation. The effect of hysteresis requires further investigation.

Figure 3.14 Average number of re-aggregation per unit of time yielded by DCU, DCUH, and LRBU policies.

The above figures reveal that the Asymmetric Simple and Full Mesh aggregation scheme of hierarchical routing perform very competitively under different simulation parameters. However, the complexity of advertised information of Asymmetric Simple is $O(n)$ while that of Full Mesh is $O(n^2)$. The MDP link cost function outperforms the COL link cost function due to more accurate cost information and more systematic call admission function. We also observe that FU, LRBU, DCU, and DCUH perform competitively in fractional reward loss and average number of crankbacks per connection request. However, the DCU and DCUH policies require the least re-aggregation frequency. Furthermore, the time-based PNNIU policies yield the highest fractional reward loss and number of crankbacks, but require significantly less update frequency. Meanwhile, the LRBU performs very competitively performance compared to the cost-based update policies, and thus can be considered a simple but efficient policy. Finally, DCU appears not to suffer from the oscillation problem, while the DCUH policy does not significantly reduce re-aggregation frequency.

Finally, the performance of PNNI routing with the proposed crankback schemes, CIS and CT, is evaluated via simulations and compared with two baseline crankback schemes, NAIVE and ABL. When crankback procedure is initiated, the NAIVE scheme tries to route the call setup message on alternate paths one by one. ABL(Avoids Block Link) is similar to NAIVE but will avoid routing on alternate paths with blocked link. Both schemes crank back the blocked setup message to the ingress node within a PG or the first PG of the same level. The blocked link information is required in the ABL scheme. Figure 3.15 shows the network topology for our simulations. The detail capacity of each link is shown in Table 3.11 and each link has 1 ms delay.

Table 3.1:   Link capacity assignments

| Link id | Cap. | Link id | Cap. | Link id | Cap. |
|---------|------|---------|------|---------|------|
| e(A11,A12) | 250 | e(A42,A47) | 400 | e(B21,B22) | 155 |
| e(A11,A13) | 250 | e(A43,A44) | 350 | e(B21,B24) | 155 |
| e(A12,A21) | 350 | e(A43,A46) | 350 | e(B22,B23) | 155 |
| e(A13,A51) | 350 | e(A44,A45) | 350 | e(B22,B24) | 155 |
| e(A21,A22) | 155 | e(A46,A47) | 400 | e(B23,B24) | 155 |
| e(A21,A24) | 155 | e(A47,B41) | 622 | e(B23,B34) | 350 |
| e(A22,A23) | 155 | e(A51,A52) | 155 | e(B31,B32) | 350 |
| e(A22,A24) | 155 | e(A51,A56) | 155 | e(B32,B33) | 155 |
| e(A23,A34) | 350 | e(A52,A53) | 155 | e(B32,B36) | 155 |
| e(A23,A24) | 155 | e(A52,A54) | 155 | e(B33,B34) | 350 |
| e(A31,A32) | 350 | e(A53,A54) | 155 | e(B33,B35) | 155 |
| e(A31,A53) | 350 | e(A54,A56) | 155 | e(B35,B36) | 155 |
| e(A32,A33) | 155 | e(A55,A56) | 350 | e(B36,B37) | 350 |
| e(A32,A36) | 155 | e(B11,B12) | 155 | e(B37,B45) | 350 |
| e(A33,A34) | 350 | e(B11,B17) | 155 | e(B41,B42) | 400 |
| e(A33,A35) | 155 | e(B12,B13) | 155 | e(B41,B46) | 400 |
| e(A35,A36) | 155 | e(B12,B15) | 155 | e(B42,B43) | 350 |
| e(A36,A37) | 350 | e(B13,B14) | 350 | e(B42,B47) | 350 |
| e(A37,A45) | 350 | e(B14,B21) | 350 | e(B43,B44) | 350 |
| e(A41,A42) | 400 | e(B15,B17) | 155 | e(B44,B46) | 400 |
| e(A41,A55) | 350 | e(B16,B17) | 350 | e(B44,B47) | 350 |
| e(A42,A43) | 350 | e(B16,B31) | 350 | e(B45,B46) | 400 |

Figure 3.15 Network topology.

In this simulations, the network supports two classes of traffic. The class 1 traffic has a bandwidth requirement of $b_1=1$ Mbps, and an end-to-end delay of $d_1=30$ ms, with arrival rate $\lambda_1^w = \lambda$. The class 2 traffic has a bandwidth requirement of $b_2=5$ Mbps, and an end-to-end delay of $d_2=30$ ms, with arrival rate $\lambda_2^w = \lambda / b_2$. The average holding time of these two classes of traffic is normalized to unity. The PNNI protocol adopts time-based PTSE update policy, referred to as the PNNIU policy. We assume that the PTSE update interval is 0.2 times of mean call holding time. Each simulation result is observed over 10 independent runs. The length of each run is 1100 units of mean call holding time. For each run, the initial 10% of the

samples were discarded. Crankback overhead is measured by the average number of nodes that a call setup message traversed during the crankback procedure.

Note that, NAIVE, ABL, and CIS, use the same call admission control policy. The difference of these schemes is where to crank back, referred as crankback destination node. Therefore, these three schemes have the same fractional reward loss. The CT scheme predicts that an alternate path will be blocked if its aggregated path cost is larger than a cost threshold. In our simulations, the threshold is set such that crankback overhead can be reduced as much as possible without increasing the fractional reward loss. According to our experiments, the best cost threshold parameter is set to $1.25 * D_{k,\min}$. Figure 3.16 shows that these schemes have the same fractional reward loss. The vertical lines about each point in the figure indicate 95 percent confidence interval. Figure 3.17 shows that both of CIS and CT schemes yield much lower crankback overhead than the NAIVE and ABL schemes. Furthermore, when we combine CIS with CT, referred to as the CIS_CT scheme, the crankback overhead can be reduced further. The reduction on crankback overhead becomes more significant as traffic load goes higher.



Figure 3.16 Fractional reward loss of different crankback schemes under PNNIU update policy.

Figure 3.17 Average hops count of crankback of different crankback schemes under PNNIU update policy.

Chang et al. show that PTSE update interval will affect the accuracy of routing information. Therefore, in following simulations, we adopt LRBU (logarithm of residual bandwidth update), as the PTSE update policy. In LRBU, the bandwidth of link $l$, is divided into $\lfloor \log_2 Cap(l) \rfloor + 1$ states. For example, a link with capacity $Cap(l)=8$ has four states with residual bandwidth of $[0,1), [1,2), [2,4), [4,8]$, respectively. The routing information will be re-aggregated and distributed when the link state changes. Figure 3.18 shows that CT yields the same fractional reward loss as compared to the NAIVE, ABL, and CIS schemes under the LRBU update policy. The CT scheme yields much lower crankback overhead than those schemes, as shown in figure 3.19. The crankback overhead yielded by the CT scheme is only half of that of the CIS scheme.

From the above figures, we observe that CIS and CT schemes reduce the crankback overhead significantly. Furthermore, combining CIS and CT schemes can further reduce the crankback overhead.
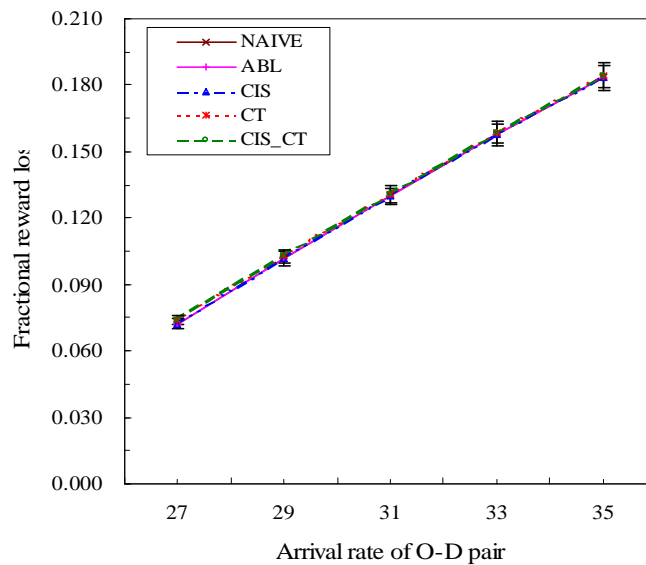
Figure 3.18 Fractional reward loss of different crankback schemes under LRBU update policy.



Figure 3.19 Average hops count of crankback of different crankback schemes under LRBU update policy.

## A. QOSPF with Overflowed Cache (PER-PAIR/OC)

According to Fig. 3.20, the overflowed-cache (OC) mechanism divides the forwarding cache into a per-destination cache (*D-cache*), a per-pair cache (*P-cache*) and an overflowed per-flow cache (*O-cache*). D-cache entries are looked up for the best effort packets, and the O-cache entry is created when a new flow arrives and cannot find sufficient bandwidth on the path of P-cache.



Figure 3.20 Overflowed Cache mechanism (Per-pair/OC)

## B. QOSPF With Overflowed Cache With Two-Phase Routing (PER-PAIR/OCTP)

Intuitively, QoS flow locality exists between node pairs as in circuit-switched networks. This section extends OC to OCTP, the Overflowed Cache with Two-Phase routing, by replacing the function **FindRouteLeastCost** with the function **FindRoute_TP**. TP can work independently of OC, known as Per-pair/TP in this study. OCTP uses three types of caches,

exactly the same as the OC scheme. However, a two-phase routing concept is used for finding the QoS path $s$ with bandwidth requirement $b$. In phase I, referred to as *soft-reservation*, OCTP tries to find a path $s_1$ with more bandwidth than $b$, i.e. where $width(s_1) \ddagger b+b_{more}$. Consequently, the subsequent incoming flows of the same S-D pair will be more likely to successfully reserve bandwidth on the path. Less misleading will increase the likelihood of success. If a soft QoS path $s_1$ cannot be found, OCTP will attempt to find a path $s_2$ with bandwidth $b$, i.e., $width(s_2) \ddagger b$, referred to as *hard-reservation* because it takes the actual required bandwidth into account. This study refers to the database that reflects the link state of soft-reservation as soft-RBDB and that of actual reservation as hard-RBDB.

## C. QOSPF Using Per-Class Routing Mark (PER-PAIR/PC)

When a new flow request with QoS requirement, Per-pair/PC first checks the forwarding cache (C-cache) if the number of sub-entries of the desired S-D pair, say $|P(s, d)|$, is zero. If yes, the C-cache will be missed and Per-pair/PC attempts to find the least costly feasible path, termed $s$. If $s$ is found, Per-pair/PC assigns a new *mark* to $s$, inserts it into the cache, and will forward packets of the flow through $s$. The link cost function could be defined according to the need of network administrators. In this paper, we simply make the cost function the inverse of path width. If $P(s, d)$ is full, Per-pair/PC simply finds the next hop $p$ of the least costly feasible path among the existing $|P(s, d)|$ paths, where $p \in P(s, d)$. If $p$ is found, the algorithm marks the flow and forwards it to $p$, otherwise it blocks the flow. If $P(s, d)$ is neither empty nor full, Per-pair/PC can either forward it to the $p$ led by the cache, or route it through a newly computed path $s$, whichever costs less. Consequently, flows between an S-D pair may be routed on a maximum of $m$ different paths where $m$ is the maximum number of routing classes.

## D. Performance Evaluation

### D.1 Network Model and Traffic Model

Simulations are run on a 40-node random graph based on the Waxman's model. In our

TABLE 3.2: The Traffic Model Of The Simulation

| *Class* | *Application* | *ratio* | *Bandwidth requirement* |
|---|---|---|---|
| $GS_1$ | Video, e.g. H.260 | 20% | 128Kbps |
| $GS_2$ | Voice, e.g. I- | 80% | 16Kbps |

TABLE 3.3: Cache Granularities Of The Simulation

| *Granularity* | *Scheme* | *Feature* | *Path computation* |
|---|---|---|---|
| Per-destination | OSPF | Lookup next-hop by destination | Topology driven |
| Per-flow | QOSPF/G* | Route each individual flow | Flow driven |
| Per-pair | QOSPF/Z | Same route between a src-dst | Topology driven |
| | Per-pair/TP | Two-phase routing | Topology driven |
| Per-pair with overflowed cache | Per-pair/OC | Dual caches | Flow driven |
| | Per-pair/OCTP | Dual caches, two-phase routing | Flow driven |
| Per-pair/class | Per-pair/PC | Diff. route for diff. class between a src-dst | Topology driven |

* QoS routing table indexed by (dst, hop_count).

* Only on-demand path-computation is used in our simulations

simulations assume that the token rate from *TSpec* of an RSVP PATH message is used as the bandwidth requirement of the flow. Furthermore, as Table 1 shows, this study assumes that there are two types of QoS traffic: $GS_1$ and $GS_2$. $GS_1$ models video sources where the bit rate is set to 128Kbps, for example videoconference, while $GS_2$ models voice source where the rate is set to 16Kbps.

## D.2 Granularity and Performance Metrics

In our simulations, five different granularities of forwarding caches used in various QoS routing schemes are studied, as shown in Table 3.2. Seven performance metrics are interesting here: (1) Request blocking probability, $P_{req}$, (2) Cache misleading probability, $P_{mish}$, (3) Fractional reward loss, $L_{rwd}$, (4) Forwarding cache size, or $N_{cache}$, is the total storage overhead for a caching scheme. (5) Number of path computations, or $N_{comp}$, is the total number of path computations in the simulated network. (6) Path length distribution is the histogram of the path length of the admitted flows. (7) Fairness of traffic requirement.

TABLE 3.4: Summary Of The Simulation Results

| Granularity | Mechanism | computation overhead | storage overhead | Blocking | misleading | fairness | path length |
|---|---|---|---|---|---|---|---|
| Per-dest | topology driven | low | low | high | high | poor | prefer short |
| Per-pair | topology driven | low | low | high | high | poor | prefer short |
| Per-flow | flow driven | very high* | very high | low** | no | fair | no discrimination |
| Per-pair /OC | flow driven | medium | medium | low** | no | fair | no discrimination |
| Per-pair /OCTP | flow driven | medium | medium | low** | no | fair | no discrimination |
| Per-pair /TP | topology driven | Low | low | medium | medium | medium | medium |
| Per-pair /PC | topology driven | medium | low | low | low | medium | medium |

\* scalable if pre-computation is used.

\*\* except heavy loading.

## D.3 Results

We has investigated the QoS routing extensions to the OSPF (QOSPF) and has proposed three mechanisms to achieve scalability with low blocking probability, overflowed cache (OC), two-phase routing (TP), and per-class routing mark (PC). OC divides the forwarding cache into a P-cache and an O-cache, and thus prevents the cache misleading effect. OC can be extended to OCTP with two-phase routing. Phase I soft-reserves more bandwidth for subsequent flows of the same S-D pair, while phase II hard-reserves actual bandwidth requirement if a flow is blocked in phase I. TP also can work independently of OC. PC aggregates the flows into several paths using routing marks, thus allowing packets to be fast forwarded in DiffServ core networks.

Extensive simulations using various routing and forwarding mechanisms found that per-destination routing has the worst blocking probability. This is because a coarser granularity is used, which reduces the accuracy of the network state. TP results in more flows running through their shortest paths than purely Per-pair. OC strengthens the path-finding ability as Per-flow scheme. OCTP combines the above two mechanisms and performs better than the alternatives. Note that, under heavy loading, the blocking probabilities performed by the flow driven mechanisms, including Per-flow, Per-pair/OC, and Per-pair/OCTP are as high as Per-dest and Per-pair. This is because many flows are traveling through longer paths which consume more resources per flow. Imposing a hop count limit $H$, where $H$ can be either the

value of network diameter or can be set explicitly, might solve the above unexpected behavior of flow driven mechanisms, which requires future studies. Additionally, Per-pair/PC has moderate blocking probability, fractional reward loss, with small forwarding cache. Per-pair/PC is suitable for the DiffServ networks because only 2 or 3 routing classes, i.e. marks, are needed.

Table 3.4 summarizes the simulation results. *Flow driven* mechanisms perform better in blocking probability, fairness, and state accuracy, while *topology driven* mechanisms result in less overheads. QoS routing and forwarding in a wire-speed core network, may use coarser granularity to achieve cheaper computation and storage cost, and forward packets faster.

Results presented herein can hopefully be applied to evaluate the overheads and performance of real network topologies. Moreover, we plan to extend the scalability issues studied in this paper to multicast QoS-based routing in IntServ and DiffServ networks.

## A. The VQS System

VQS is assumed non-cut-through and non-preemptive. In other words, a packet is not served until its last bit has arrived, and once a packet is being served, no interruption is permitted until the whole packet is completely served. It is also work conserving in the sense that the server remains busy as long as there are packets in the queue. Packets are served under a normalized service rate of one cell/slot. Given a backlogged session, $i$, assigned with weight $w_i$, VQS allocates the session a minimum service rate of $w_i / W$ (cells/slot), where $W = \sum_{i=1}^{N} w_i$ and N is the total number of sessions in the system. For ease of description, we assume the packet size is fixed (=$L$ cells). The VQS algorithm, as will be shown, requires little modification for supporting variable-size packets.

For generalization, we consider two different VQS systems- a standalone VQS and an embedded VQS. While the former directly accepts input traffic from each session, the latter exerts a leaky-bucket regulator between each session's input traffic and VQS. First, the input

traffic from each session is generally modeled by a discrete-time Switched Bernoulli Process (SBP). The process alternates between the High and the Low states. Second, the leaky-bucket regulator for session $i$ is defined by ( $\omega_i$, $\tau_i$), where $\omega_i$ (cells/slot) is the token generation rate and $\tau_i$ (cells) is the maximum token bucket size. Thus, under the embedded system, traffic from session $i$ exhibits a mean arrival rate of $\omega_i / L$ (packets/slot) and burstiness which increases with $\tau_i$.

## B. Implementation Architecture

The architecture (see Figure 3.21) includes a VLSI chip, called the Sequencer, as a key component. The Sequencer is essentially a sorting-memory chip. By cascading multiple Sequencer chips in series or in parallel, a large linked list type packet pool could be implemented. As depicted in Figure 3.22, when a packet arrives, the packet is stored in the packet pool at the address provided by the Idle-Address FIFO, which contains the addresses of unused space in the packet pool. Before the packet is written into the packet pool, its session identifier is extracted and used as an index into the Session (S)-cache. The S-cache maintains a separate entry for each session, including the normalized weight and credit. Notice that since we assume $WS$=1 in this architecture, the sum of normalized weights of all sessions is equal to 1. The Session Controller is responsible for determining the index of the window in which this arriving packet can be placed, based on the normalized weight of the session to which the packet belongs.

## C. Results and Merit Review of the Project

The performance of VQS is evaluated via simulation. Simulation results demonstrate that, applying a smaller $WS$ for network elements with sufficient computation power, VQS performs as superior as the optimal scheme, WF$^2$Q, with respect to mean delay, throughput fairness, and worst-case delay fairness (see Figure 3.22). Moreover, compatible to WF$^2$Q, VQS outperforms WFQ with respect to 99-percentile delay bound and jitter in the presence of traffic burstiness. For network elements with limited power, VQS provides the best possible

QoS using a larger window size. The design and experimental results have been presented and demonstrated in various conferences and meetings, including IEEE ICC'00. Moreover, we have designed several networking control systems making use of the mechanism, which has been submitted to IEEE ICC 2001.
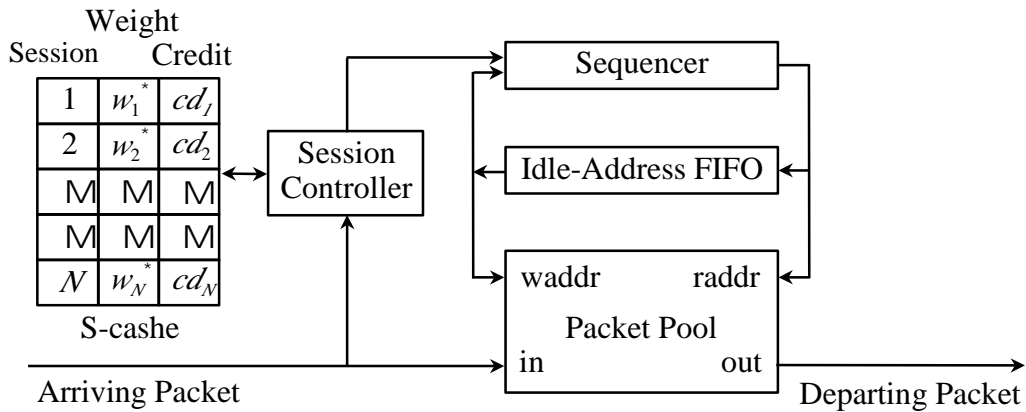


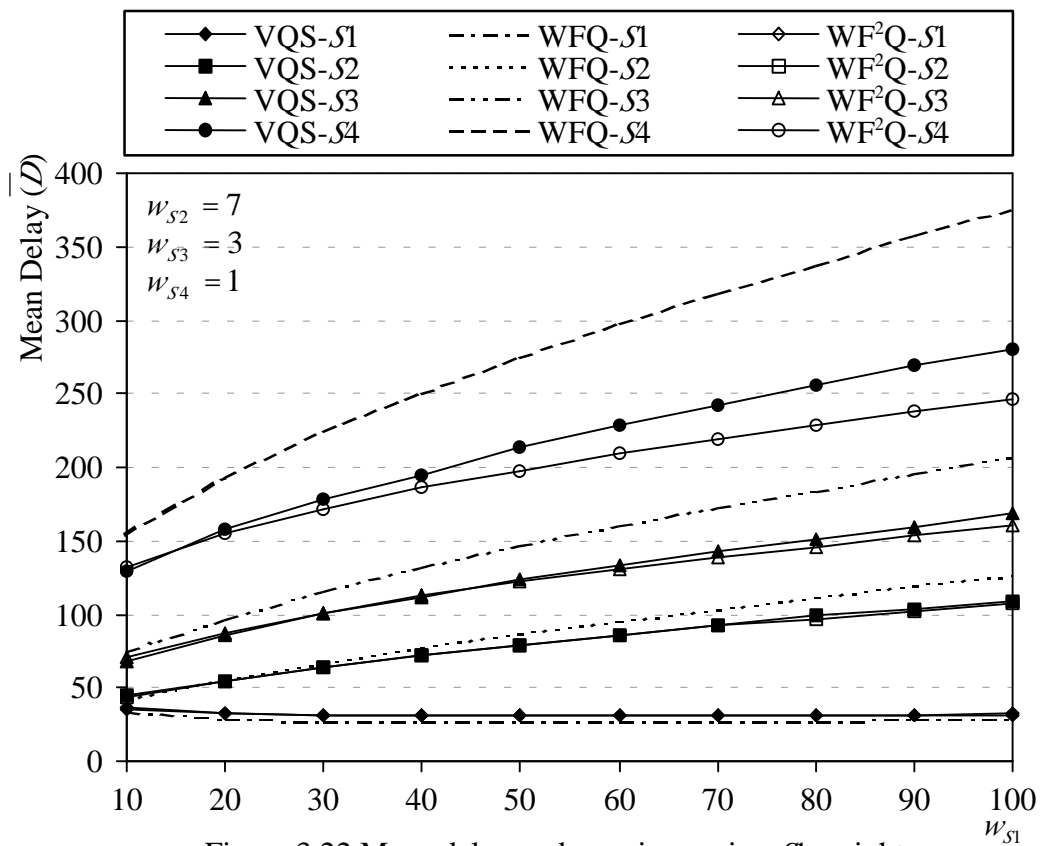Figure 3.21 VQS implementation architecture.



Figure 3.22 Mean delay under an increasing $S$1 weight.

port trunking

crossbar

hashing   search tree   linear search                    filter

QoS  (Quality-of-Service)

Prevailing weight-based

multiple queues    engage    timestamp

QoS

weight-based Versatile QoS Scheduler (VQS)            VLSI

VQS

VQS              single-queue       timestamp            weight

windows             window          session                       session

weight    window                                    window    high-power

VQS        throughput fairness   mean delay      worst-case delay fairness

WF$^2$Q                    WF$^2$Q            traffic burstiness          VQS

WFQ        99%    delay bound    jitter

(                                                )

QoS

metrics    routing algorithm                RSVP                                re-

routing

PNNI                                                            Asymmetric

Simple                                (MDP)

(COL)            (WP)

(DCU)

CIS    CT            crankback overhead            Asymmetric Simple

MDP

call admission          call            COL   WP            MDP

DCU

CIS      CT

crankback overhead

(                                                )

flow driven            per-flow    per-pair/overflowed-cache

57

data driven                 per-destination   per-pair   per-pair/two-phase

per-pair/class              overhead

[1]  ATM Forum Technical Committee, "PNNI Specification Ver 1.0," March, 1996.

[2]  A. Iwata, H. Suzuki, R. Izmailov, and B. Sengupta, "QOS Aggregation   Algorithms in Hierarchical ATM Networks," *IEEE ICC'98*, 1998.

[3]  B. Awerbuch, Y. Du, B. Khan, and Y. Shavit, "Routing Through Networks with Hierarchical Topology Aggregation," *IEEE Symposium on Computer and Communication*, pp. 406-412, 1998.

[4]  B. Awerbuch, Y. Du, and Y. Shavit, "The Effect of Network Hierarchy Structure on Performance of ATM PNNI Hierarchical Routing," *IEEE Computer Communications and Networks*, pp. 406-412, 1998.

[5]   Ben-Jye Chang, Hsien-Kang Chung, and Ren-Hung Hwang, "Hierarchical QoS Routing in ATM Networks," *The 14th International Conference on Information Networking*, Jan. 2000.

[6]  M. Montgomery and G. D. Veciana, "Hierarchical Source Routing Through Clouds," *IEEE INFOCOM'98*, 1998.

[7]  F. Hao, E. W. Zegura, and S. Bhatt, "Performance of the PNNI Protocol in Large Networks," *IEEE ATM Workshop*, pp. 315-323, 1998.

[8]  Ren-Hung Hwang and Youn-Chen Sun, "Adaptive Multicast Routing in Broadband Networks," *SPIE Conference on Performance and Control of Network Systems*, Nov. 1998.

[9]  R. A. Howard, Dynamic programming and Markov processes. John Wiley & Sons, Inc., 1960.

[10]  Ren-Hung Hwang, James F. Kurose, and Don Towsley, "State Dependent Routing for Multirate Loss Network," *Globecom'92*, pp. 565-570, 1992.

[11]  Ren-Hung Hwang, May 1993,"Routing in High-Speed Networks", Ph. D. Thesis, also available from Computer Science Dept. Technical Report 93-43, Univ. of Mass.

[12]  Borodin Allan and Ran El-Yaniv, "Online computation and competitive analysis," Cambridge University Press, 1998.

[13]  Ren-Hung Hwang and Youn-Chen Sun,  "Effect of Link Cost Function and Call Admission Control on Multicast Routing in Broadband Networks," submit for publication.

[14]  Ben-Jye Chang and Ren-Hung Hwang, "Hierarchical QoS Routing in ATM Networks Based on MDP Cost Function," to appear in *IEEE ICON 2000.*

[15] B. Awerbuch, Y. Du, B. Khan, Y. Shavit, "Routing Through Networks with Hierarchical Topology Aggregation," *IEEE Symposium on Computer and Communication*, pp. 406-412, 1998.

[16] Shun-Ping Chung and Jin-Chang Lee, "Dynamic Reservation with Hysteresis in Cellular Multiservice Networks," *The 14th International Conference on Information Networking*, Jan. 2000.

[17] Ben-Jye Chang and Ren-Hung Hwang, "Dynamic Update of Aggregated Information for Hierarchical QoS Routing in ATM Networks," submitted to *IEEE GLOBECOM 2000.*

[18] Jae-Yeul Chung, "A Predictive Alternate Path Routing Approach Supporting the Best QOS in ATM Networks", *ICCT'98*, pp. 22-24, October 1998.

[19] Felstaine, E., Cohen, R., and Hadar, O., "Crankback Prediction in Hierarchical ATM Networks", *INFOCOM '99*, pp. 671–679, 1999.

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

[21] E. Crawley, R. Nair, B. Rajagopalan, H. Sandick, "A Framework for QoS-based Routing in the Internet," *RFC 2386*, August 1998.

[22] R. Guerin, A. Orda and D. Williams, "QoS Routing Mechanisms and OSPF extensions," *Internet Draft*, draft-guerin-qos-routing-ospf-03.txt, January 1998.

[23] G. Apostolopoulos, D. Williams, S. Kamat, R. Guerin, A. Orda and T. Przygienda, "QoS Routing Mechanisms and OSPF extensions," *RFC 2676*, August 1999.

[24] G. Apostolopoulos, R. Guerin, S. Kamat, and S. Tripathi, "Quality of Service Based Routing: A Performance Perspective," *SIGCOMM'98*, September 1998.

[25] G. Apostolopoulos, S. Kamat, and R. Guerin, "Implementation and Performance Measurements of QoS Routing Extensions to OSPF," *INFOCOM'99*, July 1999.

[26] Z. Zhang, C. Sanchez, B. Salkewicz, and E. Crawley, "QoS Extensions to OSPF," *Internet Draft*, draft-zhang-qos-ospf-01.txt, September 1997.

[27] D. Black, S. Blake, M. Carlson, E. Davies, Z. Wang, and W. Weiss, "An Architecture for Differentiated Services," *RFC 2475*, December 1998.

[28] R. Braden, L. Zhang, S. Berson, S. Herzog, S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Spec.," *RFC 2205*, September 1997.

[29] P. Chemouil, J. Filipiak, and P. Gauthier, "Analysis and Control of Traffic Routing in Circuit-Switched Networks," *Computer Networks and ISDN Systems*, 11, pp. 203-217, 1986.

[30] B. M. Waxman, "Routing of Multipoint Connections," *IEEE JSAC*, Vol. 6, No. 9, pp. 1617-1622, December 1988.

[31] McKeown, N. "The iSLIP Scheduling Algorithm for Input-queued Switches," in *IEEE/ACM Transactions on Networkin*g, Vol. 7 , April 1999.

[32] Karol, M.; Hluchyj, M.; and Morgan, S. "Input versus Output Queueing on a Space Division Switch," in *IEEE Transactions Of Communications*, vol. 35, pp. 1347-1356, 1987.

[33] V. Srinivasan, S. Suri, and G. Varghese. "Packet Classification Using Tuple Space Search," in *Proceedings of ACM SIGCOMM'99*, Sept. 1999.

[34] V. Srinivasan, G. Varghese, S. Suri, and M. Waldvogel. "Fast and Scalable Layer Four Switching," in *Proceedings of ACM SIGCOMM' 99*, Sept. 1999.

[35] P. Gupta and N. Mckeown. "Packet Classification on Multiple Fields," in *Proceedings of ACM SIGCOMM'99*, Sept. 1999.

[36] Thomas Y.C. Woo. "A Modular Approach to Packet Classification Algorithms and Results," in *Proceedings of INFOCOM 2000*.

[37] Wei-Che Chen. "Design and Implemen- tation of Traffic Classification for Layer3 Switch," July, 1999.

[38] M.Buddhikot, S. Suri, and M. Waldvogel. "Space Decomposition Techniques for Fast Layer-4 Switching," in *Proceedings of IFIP Workshop on Protocols for High Speed Networks*, Salem, Massachusetts, August 1999.

[39] S. T. Chuang, A. Goel, N. McKeown, and B. Prabhakar, "Matching Output Queueing witch a Combined Input Output Queued Switch," *Technical Report*, CSL-TR-98-758, April 1998.

[40] W. Doeringer, G. Karjoth, and M.Nassehi, "Routing on Longest Matching Prefixes," in IEEE/ACM Transactions on Networking, vol. 4, Feb. 1996.

[41] D. C. Stephens and H. Zhang, "Exact Emulation of an Output Queueing Switch by a Combined Input Output Queueing Switch," in *Proceedings of IWQOS'98*.

[42] P. Krishna, N. Patel, A. Charney, and R. Simcoe, "On the Speedup Required for Work-Conserving Crossbar Switches," in *Proceedings of IWQOS'98*.

[43] N. McKeown, "Scheduling Algorithms for Input-queued Cell Switches," *Ph.D. dissertation*, Univ. California, Berkeley, May 1995.

[44] A. Charny, P. Krishna, N. Patel, and R.simcoe, "Algorithms for Providing Bandwidth and Delay Guarantees in Input-Buffered Crossbar with Speedup," in *Proceedings of IWQOS'98*.

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

[46] 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.

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

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

[49] R. Guerin, H. Ahmadi, 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.

[50] S. Jamin, P. B. Danzig, S. J. Shenker, 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.

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

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

[53] J. Liebeherr, D. Wrege, and D. Ferrari, "Exact Admission Control for Networks with a Bounded Delay Service," *IEEE/ACM Trans. on Networking*, vol. 4, no. 6, Dec. 1996, pp. 885-901.

[54] D. Lee, and B. Sengupta, "Queueing Analysis of a Threshold Based Priority Scheme For ATM Networks," *IEEE/ACM Trans. on Networking*, vol. 1, no. 6, Dec. 1993, pp. 709-717.

[55] J. Hah, and M. Yuang, "A Delay and Loss Versatile Scheduling Discipline in ATM Switches," *Proc. IEEE INFOCOM*, 1998, pp. 939-946.

[56] A. Demers, S. Keshav, and S. Shenker, "Analysis and Simulation of a Fair Queueing Algorithm," *Proc. SIGCOMM*, 1989.

[57] J. Bennett, and H. Zhang, "WF$^2$Q: Worst-case Fair Weighted Fair Queueing," *IEEE INFOCOM*, 1996, pp. 120-128.

[58] S. Golestani, "A Self-Clocked Fair Queueing Scheme for Broadband Applications," *IEEE INFOCOM*, 1994, pp. 636-646.

[59] D. Stiliadis, and A. Varma, "Efficient Fair Queueing Algorithms for Packet-Switched Networks," *IEEE JSAC*, vol. 6, no. 2, April 1998, pp. 175-185.

[60]  R. Cruz, "A Calculus for Network Delay, Part I: Network Elements in Isolation," *IEEE Trans. on Information Theory* , vol. 37, no. 1, Jan. 1991, pp. 114-131.

[61]  H. Saito, *Teletraffic Technologies in ATM Networks*, Artech House, 1994.

[62]  H. Chao, and N. Uzun, "An ATM Queue Manager Handling Multiple Delay and Loss Priorities," *IEEE/ACM Trans. on Networking*, vol. 13, no. 6, Dec. 1995, pp. 652-659.