# The University of Kansas

## Information and Telecommunication Technology Center

Technical Report

# Impact of Traffic Aggregation on Network Capacity and Quality of Service: Preliminary Results

Towela P.R. Nyirenda-Jerre
Victor S. Frost

ITTC-FY2000-TR-13200-11

February 2000

**Executive Summary**

This report contains preliminary results obtained in the study of traffic handling mechanisms for support of Quality of Service (QoS) in the Internet. Traffic handling mechanisms can be broadly classified as being aggregate, semi-aggregate or per-flow (zero-aggregation). In aggregate traffic handling there is no differentiation between traffic flows and resources are allocated to the entire set of flows as a whole. With semi-aggregate traffic handling, traffic is grouped into a small number of predefined classes based on some criteria such as the nature of delay guarantees required by the traffic. Resources are then allocated to each class of traffic. With per-flow handling, there is no grouping of traffic and each flow is allocated its own dedicated resources. Each of these traffic handling schemes can be used to meet service guarantees of different traffic types, the major difference being in the quantity of resources that must be provided in each case. For instance aggregate schemes in general require more bandwidth than per-flow schemes. The choice of which traffic handling strategy to use requires a methodology that can be used to capture the trade-off between the different schemes which is the purpose of this study.

One of the objectives of the study is to quantify the difference in capacity requirements between aggregate, semi-aggregate and per-flow traffic handling schemes. A second objective is to determine the sensitivity of the traffic handling schemes to changes in network load above the design specification. Four traffic handling schemes were used in the study, namely Weighted Fair Queueing (WFQ), First-in-First-Out (FIFO), Class-Based-Queueing (CBQ) and Strict Priority Queueing (PQ). WFQ is a per-flow scheme while FIFO is an aggregate traffic handling scheme. CBQ and PQ represent semi-aggregate schemes. Four classes of traffic were considered which have delay QoS requirements representative of Voice, Video, E-mail and WWW traffic. The voice and video represent real-time (RT) traffic while the e-mail and WWW represent non-real-time (NRT) traffic. The difference in capacity requirements was assessed by varying the load levels of all four traffic types and calculating the capacity required by the other three schemes to obtain performance equivalent to WFQ. For studying the sensitivity to load changes we calculated the delay performance of the four traffic classes when network capacity was fixed and load changes were due to voice and WWW traffic. We considered two types of networks: one which was designed with voice as the dominant traffic and the other which was designed with WWW as the dominant traffic. We summarize the key findings in the paragraphs that follow.

- Weighted Fair Queueing (WFQ)

  WFQ requires the least capacity among the four schemes and its sensitivity to changes in load depends on how bandwidth is re-allocated when the load changes. In the simplest bandwidth re-allocation strategy we assumed that each traffic class had static bandwidth allocations and any increase in traffic would result in the bandwidth of

that class being shared between the old and new connections of that class. Using this approach, the delay QoS of voice deteriorated when voice traffic increased and the delay QoS of WWW deteriorated when WWW traffic increased. A second approach to re-allocation of bandwidth assumes that the goal is to maintain the QoS of voice at all costs so that an increase in voice traffic is accommodated by reducing the allocation to e-mail and WWW and using the "stolen" bandwidth for the new voice traffic. Using this approach resulted in the e-mail and WWW QoS deteriorating when the voice traffic was increased.

- First-In-First-Out (FIFO)

  In general, FIFO required up to 2 orders of magnitude (up to 400 times) more bandwidth than WFQ. Regarding sensitivity to network load, increases in voice traffic up to 90% above the design point did not affect the delay QoS of any traffic class. This was true irrespective of whether the network was designed with voice as the dominant traffic or WWW as the dominant traffic. When the increase in load was due to WWW traffic, the delay QoS of voice was violated considerably.

- Class-Based-Queueing (CBQ)

  This traffic handling strategy required up to twice the bandwidth of WFQ. In this case increasing the voice traffic affected the delay QoS of voice only while increasing WWW traffic affected the e-mail traffic's QoS. This is because CBQ provides isolation between the RT class and the NRT class.

- Priority Queueing (PQ)

  The bandwidth requirements of PQ are similar to CBQ and are within twice the bandwidth of WFQ. In a network designed for voice, the QoS of all the classes is maintained for increases in voice traffic of up to 75% above the design point. Increasing WWW traffic only affects the e-mail QoS. In a network designed for WWW traffic, increasing the voice traffic does not affect the QoS of any class while increasing the WWW traffic affects the e-mail QoS.

One conclusion that can be drawn from this study is that on the basis of capacity requirements, there is no significant difference between semi-aggregate traffic handling and per-flow traffic handling. For voice traffic, CBQ exhibits the same sensitivity to changes in load as WFQ and the QoS of voice can be maintained by stealing bandwidth from the non-real time traffic.

a

# 1 Introduction and Motivation

When the Internet first came into being it was used primarily as a research tool and was designed to deliver uniform best-effort service to all users. The majority of traffic carried at this time was primarily data, which did not have very stringent requirements on delivery delay. During this decade the Internet has evolved into being more of a commercial entity than a research network and has experienced tremendous growth in both the volume of traffic carried as well as diversity in the type of traffic carried.

The major tool that was used to engineer the Internet was over-engineering (often referred to as "throwing bandwidth at the problem") which refers to providing more bandwidth than the aggregate demand so that every subscriber is given ample access to network resources. The engineering philosophy behind the Internet was based on the model of a homogenous community that had common interests rather than on a model of service providers and customers [14]. The best-effort Internet can be considered as consisting of just one user group in which everyone is allowed to use the network for any purpose and limits are imposed only when the capacity is not enough to satisfy demand. It is also assumed that all users behave agreeably during times of congestion by limiting their usage. The recent growth in network usage both at the commercial and public level coupled with the advances in high-speed applications however tends to stretch the limits of over-booking as more and more customers are demanding and using more bandwidth from the networks while at the same time having high expectations on the service that they receive.

The emergence of applications with diverse throughput, loss and delay requirements requires a network that is capable of supporting different levels of service as opposed to the single best-effort service that was the foundation of the Internet. Quality of Service (QoS) has become the buzzword and umbrella term that captures the essence of this shift in paradigm. IP Telephony is a good example of an application that is driving the push towards QoS on the Internet and is in fact being touted as today's killer application for the Internet [18]. Latency rather than bandwidth is the primary issue in providing voice services in the Internet, thus the traditional approaches of simply over-engineering may not work as well for this type of application. To provide a network that caters to these different levels of service requires changes to network control and traffic handling functions. Control mechanisms allow the user and network to agree on service definitions, identify users that are eligible for a particular type of service and let the network allocate resources appropriately to the different services. Traffic handling mechanisms are used to classify and map information packets to the intended service class as well as controlling the resources consumed by each class. Notable results of the effort to provide Quality of Service in the Internet are the definition of Integrated Services and Differentiated Services by the IETF [2, 3, 4, 12, 13].

The Integrated Services (Intserv) model uses resource reservation to provide delay and

throughput guarantees. The Intserv model is based on the idea that bandwidth must be explicitly managed in order to meet application requirements therefore resource reservation and admission control are a must [4, 5]. Advocates of the Intserv model claim that high fidelity interactive audio and video applications need higher quality and more predictable service than that provided by the best-effort Internet and that this can only be achieved through explicit resource reservation [6].

The Differentiated Services model takes a different approach from Intserv in that it does not promote the use of resource reservation. Proponents of Diffserv argue that a simple priority structure will be sufficient to provide Quality of Service in the Internet.

One of the arguments against resource reservation is that in the future bandwidth will be infinite, therefore there is no need to reserve it. Advances in fiber-optic communication may suggest that bandwidth will be so abundant, ubiquitous and cheap that it will not benefit network operators to undertake resource reservation. However, one cannot ignore the fact that increases in available bandwidth are always followed by corresponding development of applications that will consume and exhaust this bandwidth [4, 11]. Trends in the history of communications indicate that regardless of how much bandwidth is made available, applications are always created that quickly exhaust the supply.

Another argument against resource reservation models is that simple priority will be sufficient to meet the needs of real-time traffic. This may be true under some conditions but not always. For instance if the number of high priority real-time transmissions increases then they will all have degraded performance.

A third argument against resource reservation is that it is too expensive because reservation of resources is wasteful in that not all the reserved resources are used. This is true if all of the resource is exclusively reserved and thus it must be ensured that there is a limit on how much guaranteed traffic is allowed and provisions must be made for non-real time traffic to utilize bandwidth unused by real-time traffic [11].

Lastly, it has been suggested that delay bounds are not necessary and throughput bounds are enough. However, guaranteeing minimum throughput does not automatically result in better delay performance. Delay bounds must be explicitly guaranteed and enforced.

Opponents of reservation contend that the issue boils down to one of provisioning and that reservation-enabled networks can only provide satisfactory service if the blocking rate is low. It is believed that by adequate provisioning, a best-effort network can achieve the same performance as a reservation-based network [6, 10]. As an example consider IP telephony users who require the network to guarantee to carry 64kbps with a maximum end-to-end latency that is no larger than 100msec. If an IP network is provisioned to accommodate N users simultaneously with the end-to-end latency within 100msec, an increase in traffic beyond N would result in the service of all the current users being degraded and the resources wasted since no user would attain acceptable performance [15]. Thus significant over-provisioning is required. The higher the quality of guarantee, the more over-provisioning that must be done for the same level of user satisfaction and hence the lower the efficiency of

network utilization. Consequently the quality of guarantees must be traded-off against the efficiency of network resource usage. The case for over-provisioning is that declining prices in bandwidth will make the extra capacity required in a best-effort Internet more economical than the complexity of supporting reservations.

Neither a pure best-effort model such as the current Internet, nor a pure guaranteed service model such as the Integrated Services model can provide an efficient solution in a multiple service environment [14]. Having a large number of service classes increases the management overhead and impairs cost efficiency. An integrated network must balance the trade-off between performance and flexibility while ensuring that performance of traffic with real-time guarantees is not degraded. Providing QoS in the Internet requires providers to re-evaluate the mechanisms that are used for traffic engineering and management in their networks. Over-engineering is an attractive option because it is simple and it has been said that within a well-defined scope of deployment it can prove to be a viable solution [10]. Recent proposals are calling for more active traffic management in the Internet that will be used to make more efficient use of resources while allowing providers to offer varying levels of service suited to the different applications being supported. These traffic management mechanism range from simple admission policies to complex queuing and scheduling mechanisms within routers and switches.

We can envision several alternative paths for future networks to follow in their quest to provide QoS. These are:

1. Inefficient use of network bandwidth with no traffic management. This approach assumes that bandwidth is abundant and cheap and thus traffic management is not needed.

2. Moderately efficient use of network bandwidth with simple traffic management

3. Efficient use of network bandwidth with complex traffic management. With this approach the assumption is that the cost of bandwidth justifies the use of traffic management.

Knowledge of the network capacity required to achieve comparable user perceived performance will indicate the importance of traffic management as the network evolves. For example, if an aggregate network capacity of 10Gb/s is needed given no traffic management while only 100Mb/s is needed when the traffic is controlled, then the cost of traffic management can be justified. However, if the difference in required capacities is "small" then it may not be time to deploy complex traffic management functionality. There is a need for a clearer understanding of the issues surrounding the provision of QoS in IP-based networks as well as guidelines on how traffic management and network capacity can be used to provide QoS.

In this report we consider the issue of finding the cross-over point at which the three approaches of no traffic management, simple traffic management and complex traffic man-

agement become equivalent. Specifically we would like to determine the network capacity required to achieve equivalent levels of performance under a variety of traffic management schemes. Knowledge of this crossover point will help network engineers and decision-makers determine the suitability of IP QoS traffic management as well as the type of traffic management to use.

In section 2 we provide a discussion on the correspondence between traffic management schemes and traffic aggregation and consider some of the questions that need to be addressed in comparing traffic management strategies. Section 3 describes an analytic study that was undertaken to illustrate how the issues raised in Section 2 could be addressed using a single-node network for illustration. In Section 4 we describe how this work can be extended to carrier-size networks and we conclude with the significance of this work to Sprint in Section 5.

# 2 Traffic Aggregation, Quality of Service and Network Capacity

The Internet's need to support traffic with diverse requirements and with differing levels of service coupled with the transition of the Internet from a research network to a commercial one has resulted in the re-definition of the Internet's architecture. The major change is in the definition of new services and traffic handling mechanisms that can be used to provide differentiated and guaranteed quality of service in the Internet.

The challenge facing the deployment of integrated services is to satisfy the strict delay and loss guarantees required for real-time services while realizing the economics of statistical multiplexing which are essential for high-speed bursty data. One objective is to be able to support both voice, video and data traffic on one network in such a way that the performance of voice is equivalent to that on a Public Switched Telephone Network (PSTN) network.

Providing guaranteed QoS today can be achieved in one of three ways. The first technique is to over-provision the network which is the classical "throw bandwidth at the network" solution. This is based on the premise that bigger pipes mean less congestion and hence better performance. The second alternative is to reduce delay by introducing the notion of precedence and treating certain types of traffic with higher priority than others. Delay for higher priority traffic in this case will be better than best-effort but will depend on the traffic load in each priority level. The last technique is to use dedicated resources for each flow in the network, recently referred to as "throwing hardware at the network". This gives the most predictable performance [1, 14].

The above solutions can be related to the level of aggregation of flows used by traffic handling mechanisms within the network. We define three levels of aggregation as shown in Figure 1.

As can be seen from the figure, in a total aggregation environment, all flows are enqueued in the same buffer and share the buffer and link resources. This is the simplest and most prevalent form of traffic handling. The link must be configured with enough capacity to meet the most stringent QoS and the typical approach to maintaining QoS in this situation is to add more capacity to the link - "throwing more bandwidth".

In the partial aggregation environment, flows are divided into classes based on some criteria, the most obvious one being to group flows with similar QoS requirements. In this way, the QoS needs of a class of flows can be ensured in isolation from other flows. This type of aggregation corresponds to the precedence solution.

In an environment with zero aggregation, each flow is assigned its own set of resources and thus attains its QoS independent of other flows. This is the best means of ensuring QoS but it is also the most complex to administer. This environment corresponds to the dedicated resources solution. The common term for zero aggregation is per-flow queueing.

Scheduling mechanisms are used to achieve the levels of aggregation that we have outlined.
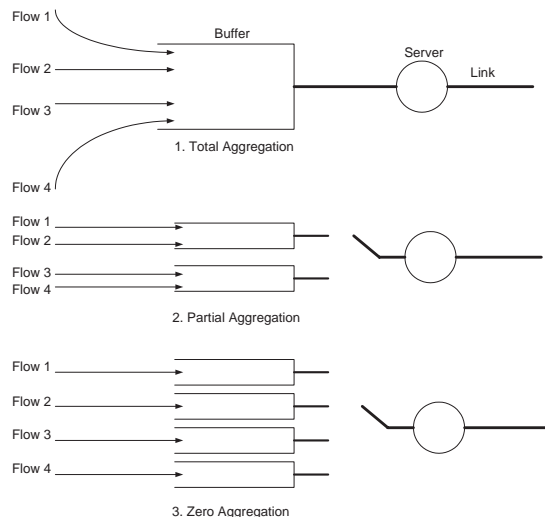
Figure 1: Levels of Aggregation

Total aggregation can be achieved with First-In-first-Out (FIFO) scheduling in which packets are served in the order of arrival to a queue. For partial aggregation Priority Queueing (PQ) and Class Based Queueing (CBQ) are typical approaches. Priority Queueing imposes a strict service order by assigning each queue to a fixed priority level and serving the queues accordingly. With Class-Based Queueing, flows are mapped to classes based on some predefined attribute and service weights are assigned to each class. Per-flow queueing can be implemented using (Weighted) Fair Queueing, (Weighted) Round Robin and their many variants.

Given the levels of aggregation and the associated scheduling mechanisms which we couple under the umbrella term of traffic handling, the question facing the network engineer is that of determining the equivalence of the different traffic handling mechanisms in terms of their ability to support traffic with varying QoS requirements. Of particular interest is the trade-off between the complexity of traffic handling mechanisms and the network capacity required to support QoS.

In addition to the traffic aggregation in traffic handling, the solution to providing QoS depends on the network capacity. It is widely accepted that the use of aggregate schemes may necessitate the provisioning of more network capacity than per flow schemes but it is not clear just how much more capacity is needed nor is it clear how the complexity of per-flow management measures up against the cost of additional capacity with aggregate traffic handling. To provide an adequate answer to this problem requires some quantification of the gain in performance obtained by using complex traffic handling with smaller network capacity versus using simple traffic handling with abundant network capacity. A pertinent issue also has to do with the sensitivity of the selected solution to changes in network load both in terms of the total load and in terms of the relative mix of different classes. Suppose that using aggregate traffic handling requires high capacity links but the resulting network

is insensitive to fluctuations in network traffic whereas using a complex scheme with limited capacity results in a network that is very sensitive to network variations, what would be the better option? It is issues such as these that need to be addressed.

Based on the foregoing, four objectives have been identified. The first objective is to examine the trade-off between complexity of traffic handling and required network capacity by comparing the bandwidth required for a given level of performance under traffic handling schemes that range from complex to simple. A second objective is to determine to what extent the analytical methods we intend to use are able to scale with network size and capacity and what modifications if any must be made to ensure that they do. In evaluating the performance under different traffic handling schemes we must ensure that the analysis is robust and scalable. Results obtained should be consistent in any network topology or configuration. If the analysis is not robust or scalable then it will provide results that are misleading. A third objective is to provide insight into how connection-less networks such as the Internet can be used to support traffic with diverse QoS requirements and to provide the analytic framework for deciding on a traffic handling and capacity provisioning strategy. A final objective is to study the sensitivity of the traffic handling algorithms to changes in network load and traffic mix.

We anticipate two main results from this study. The first result is a quantification of the trade-off between complexity of traffic management and network capacity. Such a quantification would take the form of a graph showing the trend in the capacity requirements of the different traffic handling requirements. The simplest representation is the capacity required by the three traffic handling models for the same network load and performance as shown in Figure 2.
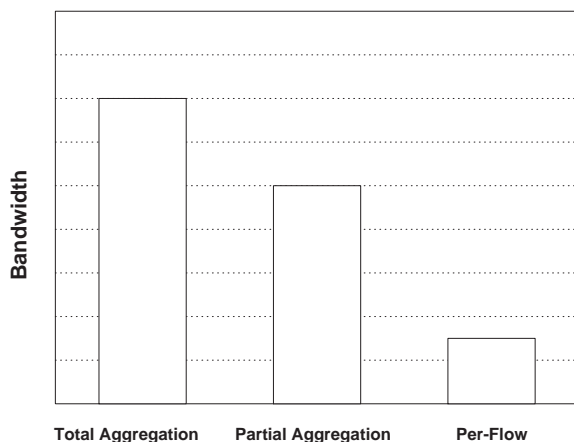


Figure 2: Simple Traffic handling and Network Capacity Trade-off

From Figure 2 we can obtain quantification of the extra bandwidth required by aggregate schemes when compared to a per-flow scheme. By taking measurements of the required capacity for equivalent performance over a variety of network loads we can obtain a graph

that shows how the difference in performance depends on the network load (level of utilization in the network). A hypothetical example of such a plot is shown in Figure 3.
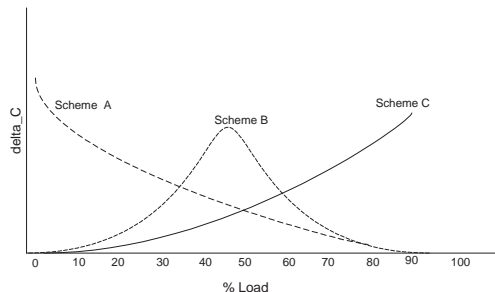


Figure 3: Traffic handling and Network Capacity Trade-off with varying Network Load

In this figure, we plot the difference in capacity (delta_C) of three traffic handling schemes A,B,C as a function of network load with reference to a per-flow scheme such as WFQ. From the plot we are able to immediately identify the points and regions where the different mechanisms provide equivalent performance and are also able to assess how this equivalence translates into a difference in network capacity requirements.

A second result that we anticipate is in the difference in sensitivity of the traffic handling parameters to network conditions and one way of illustrating this difference is as shown in Figure 4. In this figure, the design point represents the point at which the delay objectives are satisfied for a given network capacity and load and the figure illustrates how the delay perceived by a candidate traffic class may vary when the network load is varied above and below the design point for three traffic handling schemes. The sensitivity can thus be measured by the ratio of change in delay to change in network traffic and this can be used to determine which scheme is more preferable. It is apparent that we would like to pick the scheme with the least sensitivity especially at loads above the design point and in this case Scheme B would be the likely candidate.
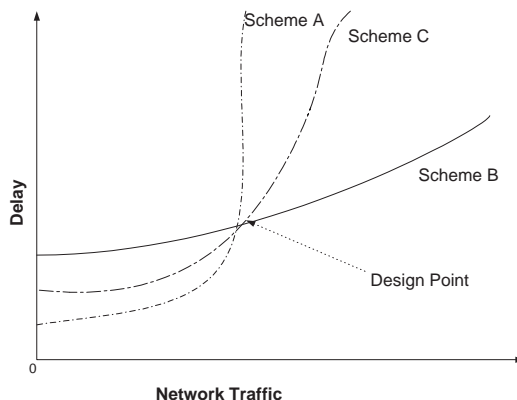


Figure 4: Comparison of Traffic Handling Sensitivity

By combining the observations from the capacity-traffic handling trade-off and the sensitivity analysis, we can provide a quantitative answer to the issue of selecting an appropriate traffic handling mechanism that meets the objectives of supporting traffic with diverse requirements in an efficient manner. In Section 3 we describe an analytical study within the context of a single-node network that was performed to demonstrate how the issues raised in this section could be addressed.

# 3 Analytic Study of Traffic Aggregation in a Single Network Node

In this section we describe the methodology and results that were obtained from analysis of traffic handling schemes in a single network element.

## 3.1 Methodology

### 3.1.1 User Characterization

We considered three aspects of user characterization. The first is the identification of the applications that are likely to prevail in a network offering differentiated and guaranteed quality of service. Having identified the applications the second aspect to characterization is the specification of the nature of quality of service guarantees that are required for each application. The third aspect of characterization is with respect to the way in which the application is described to the network, often referred to as traffic modeling.

In Table 1 we list the four applications that were used and their characteristics.

| Application | RT/NRT | Rate type | QoS |
|---|---|---|---|
| Telephony | RT | Stream | low delay |
| Interactive Video | RT | Stream | low delay, low loss |
| E-mail | NRT | Burst | delay tolerant |
| WWW | NRT | Burst | delay tolerant |

Table 1: Network Applications

Based on the literature we also identified various parameters for each class as shown in Table 2.

| Class No. | Class | Delay (s) | Average Rate $\rho$ (Mbps) | Burstiness $\sigma$ (Bytes) | Packet Size (Bytes) |
|---|---|---|---|---|---|
| 1 | Voice | 0.002 | 0.064 | 64 | 64 |
| 2 | Video | 0.005 | 3 | 2560 | 512 |
| 3 | E-mail | 0.5 | 0.128 | 320 | 64 |
| 4 | WWW | 1.0 | 2 | 5120 | 512 |

Table 2: Traffic Class Parameters

The voice and video applications belong to the Real-Time (RT) traffic class while the e-mail and WWW traffic belong to the Non-Real Time (NRT) traffic class. Our choice of these applications was based on the fact that they are representative of current network usage and they provide diversity in their attributes and QoS.

For characterization of the traffic sources we used the burstiness constraint model of Cruz [7] in which traffic is characterized by two parameters, a burstiness parameter $\sigma$ and an average rate parameter $\rho$. We assume that the network uses regulator elements or shapers to ensure that the traffic entering it conforms to these parameters. We chose to use this bounded model for the traffic processes so that the results obtained are general and applicable to a variety of situations and do not depend on specific traffic assumptions. The model is very appealing because both the IETF and ATM Forum have defined network elements which can convert an arbitrary traffic process into a process that is bounded in this way.

### 3.1.2   Traffic Handling Mechanisms

We classified traffic handling mechanisms as simple, intermediate and complex depending on whether they are used for total aggregation, partial aggregation or per-flow handling respectively. We identified four candidate traffic handling mechanisms as shown in Table 3:

| Classification | Mechanisms | Abbreviation |
|---|---|---|
| Simple | First-In-First-Out | FIFO |
| Intermediate | Strict Priority Queueing | PQ |
| | Class-Based Queueing | CBQ |
| Complex | Weighted Fair Queueing | WFQ |

Table 3: Traffic Handling Mechanisms

We chose these mechanisms because they are representative of current and future implementations in network routers and switches. Figures 5 to 8 illustrate how the applications are handled by the four schemes.



Figure 5: Aggregate Traffic Handling using FIFO



Figure 6: Per-Flow Traffic Handling using WFQ

For the CBQ scheme, each class is assigned a guaranteed rate g_rt and g_nrt respectively, while in WFQ each application is assigned its own guaranteed rate g_voice, g_video, g_email and g_www respectively
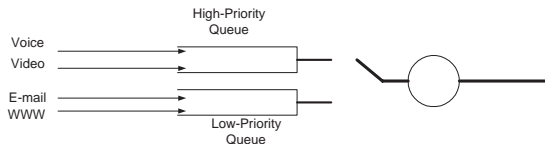
11

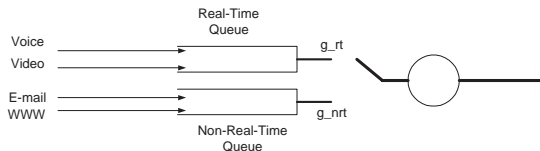Figure 7: Partial Aggregate Traffic Handling using PQ



Figure 8: Partial Aggregate Traffic Handling using CBQ

## 3.2   Analysis

In order to obtain results that are easily understood and verified we focused on the simplest model of a network with a single network router or switch. In addition to the application parameters in Table 2, other parameters that were used are:

- Link Capacity $C$ Mbps

- Reservation Factor $f$

- Reserved Bandwidth $C_{resv} = f * C\,Mbps$

- Expected utilization of class i $\lambda_i$ with $\sum_{i=1}^{4} \lambda_i = 1$

- Total Bandwidth allocated to class i $g_{total}(i) = \lambda_i * C_{resv}$ Mbps

Note that consistent with the notation in Table 2, we use the indices 1,2,3 and 4 to represent Voice, Video, E-mail and WWW traffic respectively. It should be noted that WFQ is used as the reference mechanism and that the reservation and utilization parameters are with respect to the link capacity used with WFQ. Note that the reservation factor represents the amount of traffic that is considered to be reserved and that will be shared among the four applications according to each class's expected utilization. We also assumed that whatever bandwidth is not used after reservations have been accounted for will be used by Best Effort (BE) traffic. The results presented in this report assume a link of OC-3 capacity and use a reservation factor of 0.96 to give a reserved bandwidth of 150Mbps. We use the reservation factor to capture the effect of a link whose bandwidth is partitioned into bandwidth reserved for guaranteed traffic and bandwidth that is available as best-effort.

We used three different values of utilization parameters for video to give $\lambda_2 = [0, 0.1, 0.2]$ expressed as a fraction of the link bandwidth. For each of these three values, the voice utilization $\lambda_1$ was varied in increments of 0.1 from 0.1 to $(1 - \lambda_2)$. We used 5 different weights to control how the remaining bandwidth after the voice and video were accounted for was shared between e-mail and WWW traffic. Denoting the weight vector as $w = [0.1, 0.3, 0.5, 0.7, 0.9]$, in each case the e-mail and WWW utilization was calculated as:

$$\lambda_3 = w * (1 - (\lambda_1 + \lambda_2)) * C \tag{1}$$

12

$$\lambda_4 = (1 - w) * (1 - (\lambda_1 + \lambda_2)) * C \tag{2}$$

where w is one of previously mentioned weights. Using these parameters allows us to examine the effects of varying the proportions of the four traffic classes.

Note that the utilization parameters capture the manner in which the reserved bandwidth is shared between the four classes. Also note that these proportions are with reference to the link capacity for WFQ which is used as a datum for the other schemes. For each reservation factor and expected utilization parameter we found the number of sources that could be supported for each traffic class using WFQ assuming an OC-3 link(155Mbps). This was done by first finding the guaranteed rate required for a single source from each class using the formula:

$$g_i = \frac{\sigma_i}{D_i} \tag{3}$$

where $D_i$ is the delay for class i. The number of connections for class i is then given by:

$$N_i = \left\lfloor \frac{\lambda_i * C_{resv}}{g_i} \right\rfloor = \left\lfloor \frac{g_{total}(i)}{g_i} \right\rfloor \tag{4}$$

We then determined how much capacity would be required to support the same traffic using the other three schemes. For CBQ, the required bandwidth $C_{CBQ}$ was found as :

$$C_{CBQ} = \frac{\sigma_{RT}}{D_{RT}} + \frac{\sigma_{NRT}}{D_{NRT}} \tag{5}$$

where

$$
\begin{aligned}
\sigma_{RT} &= N_1\sigma_1 + N_2\sigma_2 \\
D_{RT} &= min\{D_1, D_2\} \\
\sigma_{NRT} &= N_3\sigma_3 + N_4\sigma_4 \\
D_{NRT} &= min\{D_3, D_4\}
\end{aligned}
$$

For Priority Queueing, the required capacity $C_{PQ}$ is found as:

$$C_{PQ} = max\{C_1, C_2\} \tag{6}$$

$$C_1 = \frac{\sigma_{RT}}{D_{RT}} \tag{7}$$

$$C_2 = \frac{\sum_{i=1}^{4} N_i \sigma_i}{D_{NRT}} + \rho_{RT} \tag{8}$$

where

$$\sigma_{RT} = N_1 \sigma_1 + N_2 \sigma_2$$

$$D_{RT} = min\{D_1, D_2\}$$

$$D_{NRT} = min\{D_3, D_4\}$$

$$\rho_{RT} = N_1 \rho_1 + N_2 \rho_2$$

For FIFO, the capacity $C_{FIFO}$ is given by:

$$C_{FIFO} = \frac{\sum_{i=1}^{4} N_i \sigma_i}{\min_i D_i} \tag{9}$$

Note that in all three cases, the required capacity was found using the relation below which is derived from the work of Parekh and Gallager [16, 17] and from the network calculus rules of Cruz [7]:

$$C_{required} = \sum_i N_i g_i = \frac{\sum_i N_i \sigma_i}{\min_i D_i} \tag{10}$$

For the PQ scheduler, we assume that high priority traffic is not affected by lower priority traffic. For CBQ and PQ, the formula in Equation 10 is applied to the RT and NRT queues separately and for PQ we take the maximum over the two queues while for CBQ we take the sum since both queues must get a guaranteed rate.

## 3.3 Comparison of Bandwidth Requirements

In this section we present results on the difference in bandwidth requirements of the four schemes using the methodology of Section 3.2. We plot the capacity requirements of CBQ, PQ and FIFO relative to WFQ for each of the five weights used for e-mail and WWW traffic in Equations 1 and 2. For each weight we plot the ratio of the CBQ, PQ or FIFO

14

Figure 9: CBQ Capacity Requirement with 0% Video



Figure 10: CBQ Capacity Requirement with 20% Video

requirements to the WFQ bandwidth respectively for each setting of the voice utilization. Separate graphs are plotted for each setting of the video utilization.

In Figures 9 and 10 we show the capacity requirements of CBQ for the case of 0% and 20% video. From these figures we note that the capacity requirements of CBQ are of the same order of magnitude as WFQ and are within twice the capacity of WFQ. The general trend is that increasing the proportion of voice results in a decrease of the difference in capacity between CBQ and WFQ and at 100% voice, the capacity requirements are the same. Regarding the proportion of e-mail and WWW traffic we observe that the capacity requirements are greater when the weight is smaller which is when the proportion of WWW traffic is greater. This is because within the NRT queue, increasing the proportion of WWW traffic requires more capacity to ensure that the e-mail traffic is still able to meet its delay objective. Comparing the two figures also shows that the introduction of video traffic increases the capacity requirements by less than 10%.

With PQ, Figures 11 and 12 show that the capacity requirements of PQ are not monotonic with changes in the voice load. Consider the case when the weight is 0.9 in Figure 11. We observe from 10% to 60% voice load, the capacity requirements decrease as we approach 60%. Above 60%, the capacity requirements start to increase. With other values of weight and with 20% video load, we observe the same trend although the point of inflexion changes. The reason for this trend is that with the priority system we calculate two separate capacities, one for the high priority queue and the other for the low priority queue and we take the maximum of the two as the required capacity. The results thus indicate that below the critical voice load(point of inflexion), the bandwidth requirements are largely due to meeting the needs of the lower priority NRT traffic whereas above the critical value the capacity requirements are determined by the needs of the higher priority RT queue. This is further verified by the observation that above the critical voice load the capacity requirements do not depend on
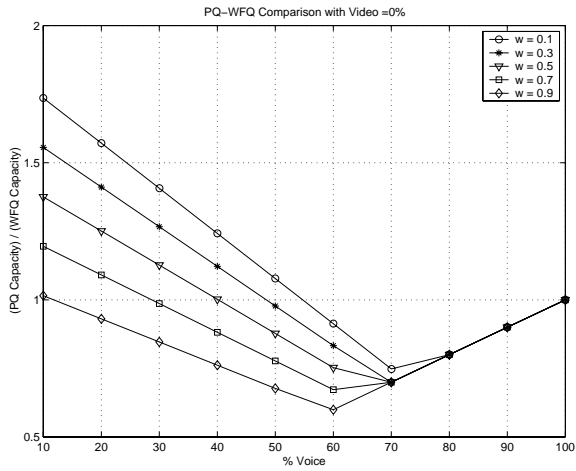
15
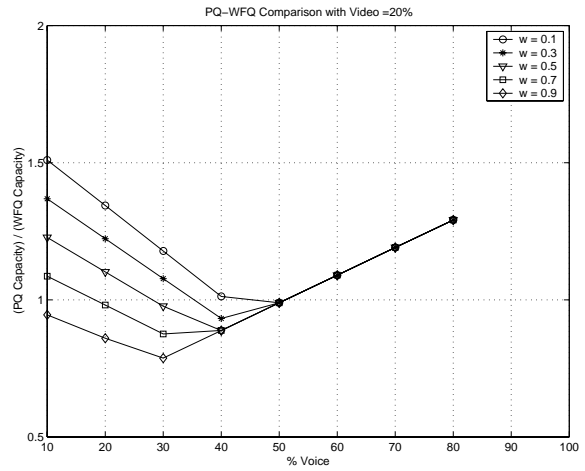
Figure 11: PQ Capacity Requirement with 0% Video



Figure 12: PQ Capacity Requirement with 20% Video

the relative proportions of e-mail and WWW traffic. In general the bandwidth requirements of PQ are of the same order of magnitude as WFQ and can in some cases be less than WFQ. The introduction of video traffic increases the capacity requirements by as much as 30%. In Figures 13 and 14 we compare the capacity requirements of CBQ and PQ.
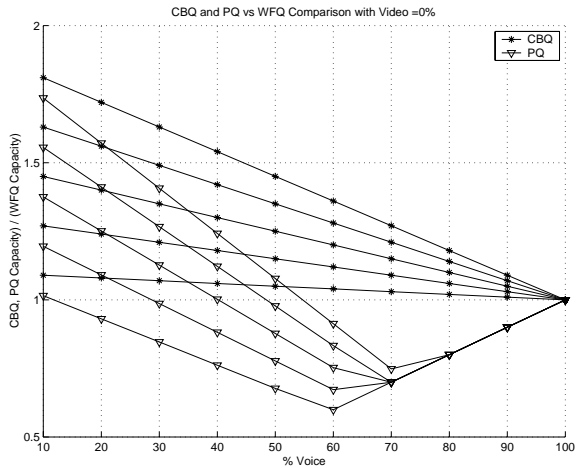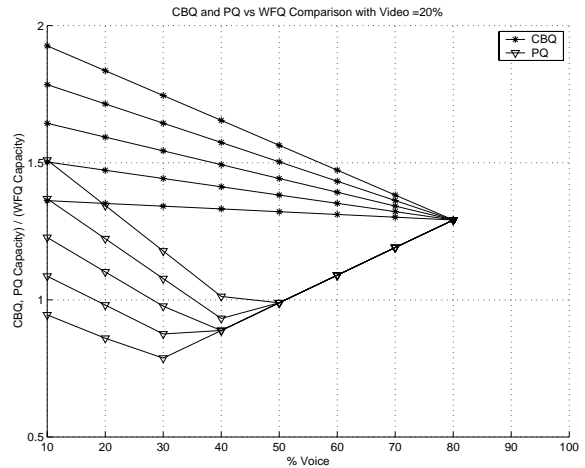


Figure 13: CBQ and PQ Capacity Requirement with 0% Video



Figure 14: CBQ and PQ Capacity Requirement with 20% Video

Note that we have plotted the results using the weights as before but have omitted the labeling for clarity. Below the critical voice load, the decrease in capacity requirements for PQ is much greater than with CBQ as evidenced by the slopes of the graphs. Above the critical voice load, the PQ capacity increases but is still less than the CBQ capacity. We also note that above the critical voice load PQ is not sensitive to the proportion of e-mail
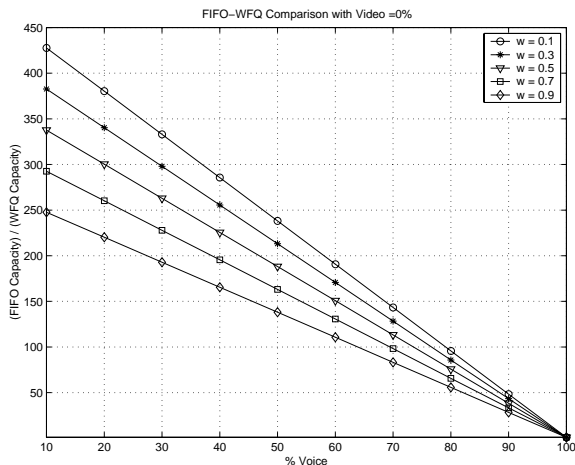
and WWW load.



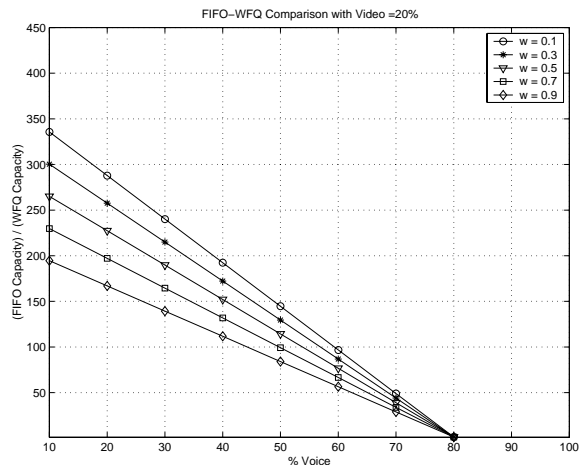Figure 15: FIFO Capacity Requirement with 0% Video



Figure 16: FIFO Capacity Requirement with 20% Video

In Figures 15 and 16 we show the capacity requirements of FIFO. The most striking observation is how the capacity requirements of FIFO are in all cases two orders of magnitude greater than WFQ. The general trend of the results is monotonic as was the case for CBQ and the capacity requirements decrease with increasing voice load. Similar to the other schemes, increasing the weight reduces the capacity requirements. In contrast to the other results, introducing video traffic reduces the capacity requirements by as much as 40%. The reason for this is that adding video traffic reduces the proportion of e-mail and WWW traffic and hence reduces the capacity required to ensure that the e-mail and WWW traffic get the same performance as voice traffic which is required in a FIFO environment. In Figures 17 and 18 we compare the absolute bandwidth requirements of WFQ, CBQ, PQ and FIFO when the video load is 20%.

From Figure 17 we notice that the FIFO capacity requiements far exceed those of the other three schemes and in Figure 18 we reduce the scale of the plot for better visualization of the CBQ, PQ and WFQ results. To obtain these graphs we took the maximum capacity required for each setting of voice load over the 5 different weights. Thus to some approximation, these results are somewhat independent of the relative proportions of e-mail and WWW traffic. In general, there is no significant difference between the four schemes when the voice traffic is between 80 to 100%. The bandwidth requirements of FIFO and CBQ decrease monotonically with increasing voice load while the PQ bandwidth is monotonically decreasing at voice loads below 50% and increasing at voice loads above 50%. The bandwidth requirements of CBQ and PQ do not exceed twice that of WFQ while FIFO bandwidth is of the order of 100's more than WFQ.
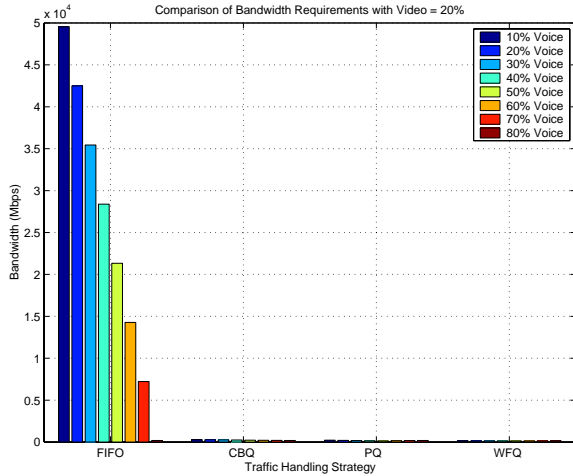
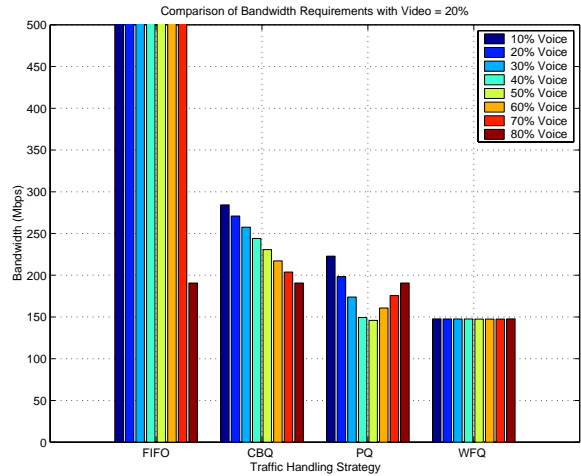Figure 17: Comparison of Capacity Requirement with 20% Video



Figure 18: Comparison of Capacity Requirement with 20% Video

## 3.4 Sensitivity to Design Point

The goal of this analysis was to explore the ability of the three schemes to provide acceptable delay QoS guarantees when the traffic submitted exceeded the traffic for which the network was designed. For a fixed allocation of bandwidth between the four classes, the capacity required by each of the four schemes was calculated using the procedures in Section 3.2. The number of sources, the link capacities and the delay performance are collectively referred to as the design point. Using these capacities, either the volume of voice or WWW traffic was varied and the delay for each traffic class was calculated by inverting the formulas in Equations 3 through 9. We considered two broad cases: one in which the majority of traffic at the design point was voice and the other in which the majority of traffic at the design point was WWW. For each of these cases we used three values of video load: 0, 10 and 20%.

For WFQ, we used two approaches for re-allocation of bandwidth when the load changed. In the first method which we call WFQ1, an increase in the traffic of a particular class resulted in the design bandwidth allocation for that class being re-distributed equally among the sources (old and new) of that class. As a result, the allocations for each class remained the same as at the design point. In the second approach called WFQ2, an increase in voice bandwidth was accommodated by "stealing" bandwidth from the e-mail and WWW classes to guarantee the voice traffic its delay QoS. The same amount of bandwidth was taken from the e-mail and WWW classes and allocated to the new voice sources. We present the results obtained in the sections that follow.

### 3.4.1 Network with Voice as the dominant traffic class

In this case the proportion of voice traffic at the design point was either 50% , 40% or 30% corresponding to video loads of 0, 10 and 20% while e-mail and WWW were 25%. We present results only for the case of 10% video load.
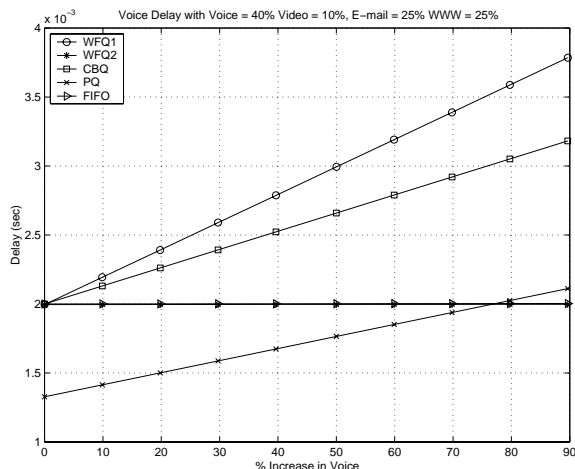


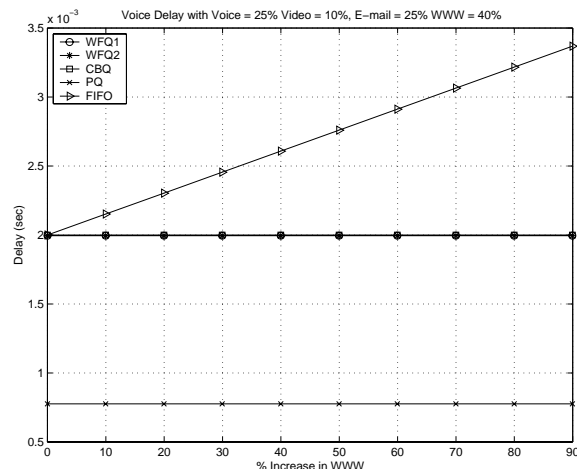Figure 19: Variation in Voice Delay with increase in Voice load

Figure 20: Variation in Voice Delay with increase in WWW load

In Figures 19 and 20 we show how the delay of voice traffic changes with increasing voice and WWW load respectively. We observe that FIFO and WFQ2 are able to maintain the delay guarantees for voice when voice traffic increases. With PQ, the delay guarantees are met for an increase of up to 75% above the design point. With WFQ1, the delay increases linearly with increasing voice load since the bandwidth available to each source decreases linearly. CBQ exhibits the same behavior as WFQ1 since the new sources have to share the same capacity that was allocated at the design point. When WWW traffic is increased, FIFO is not able to maintain the delay requirements for voice whereas all the other three meet the voice delay objectives. Increasing the voice and WWW loads affects the e-mail and WWW delays differently as shown in Figures 21 to 24.

Increasing the voice load increases the e-mail and WWW delays exponentially when WFQ2 is used. With the other schemes, increasing the voice traffic has no noticeable effect. Increasing the WWW load deteriorates the performance of e-mail when using PQ and CBQ whereas with WWW the performance deteriorates only with WFQ. The reason is that with PQ and CBQ the capacity is chosen so as to meet the requirements of e-mail which are more stringent than those of WWW. So increasing WWW traffic will impact the performance of e-mail more than it impacts WWW under PQ and CBQ.
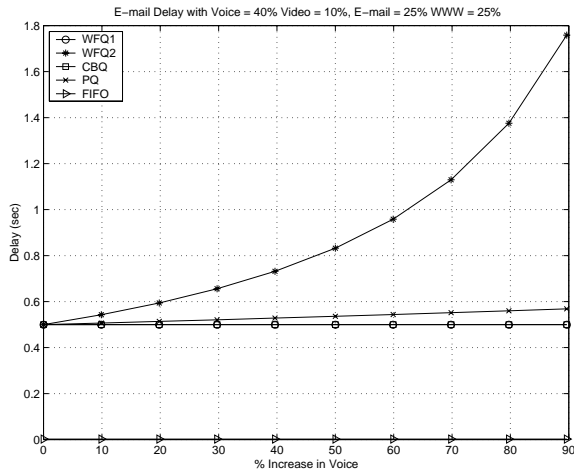
Figure 21: Variation in e-mail Delay with increase in Voice load
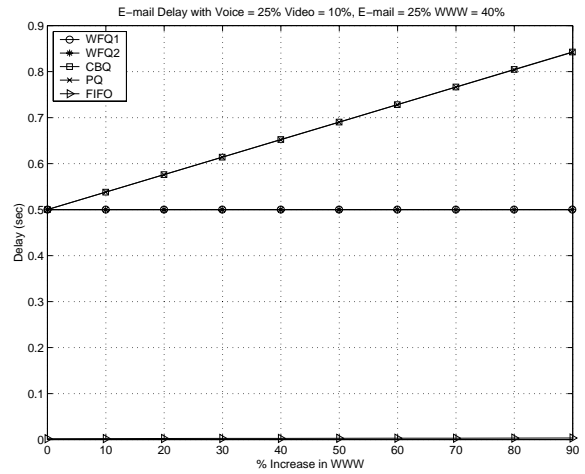


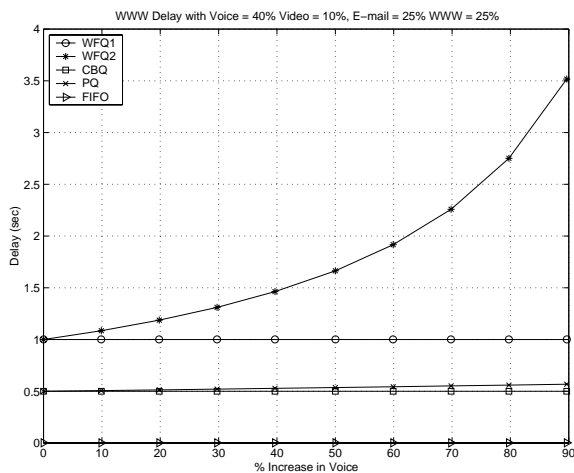Figure 22: Variation in e-mail Delay with increase in WWW load



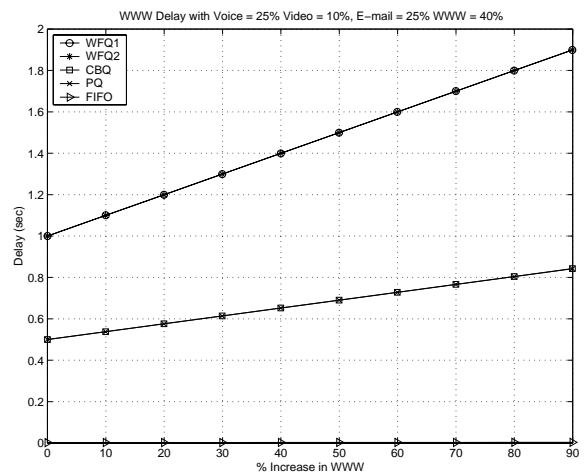Figure 23: Variation in WWW Delay with increase in Voice load



Figure 24: Variation in WWW Delay with increase in WWW load

## 3.5 Network with WWW as the dominant traffic class

In this case the proportion of WWW traffic at the design point was either 50% , 40% or 30% corresponding to video loads of 0, 10 and 20% while voice and e-mail were both 25%. We present results only for the case of 10% video load.
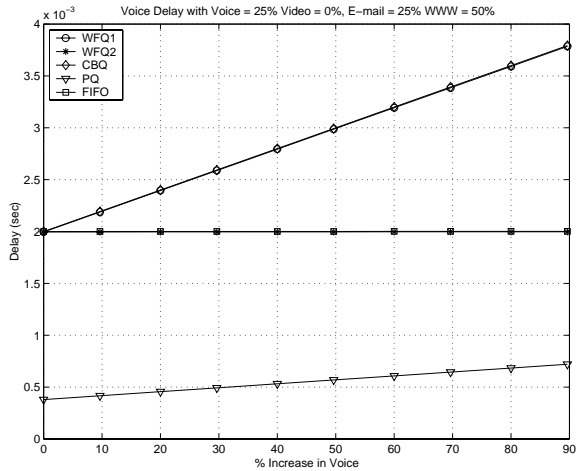


Figure 25: Variation in Voice Delay with increase in Voice load
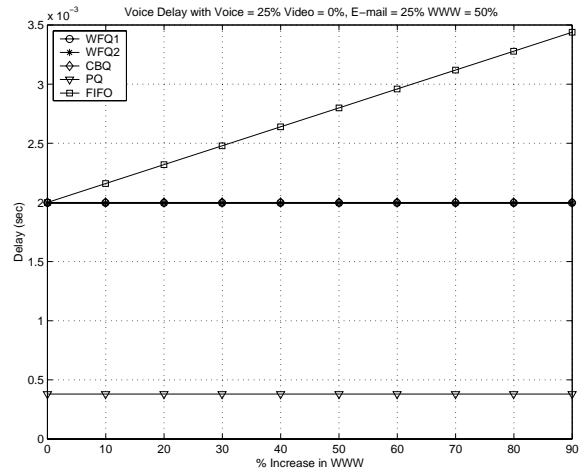


Figure 26: Variation in voice Delay with increase in WWW load

The trend of the results for voice traffic as shown in Figures 25 and 26 is the same as when voice was the dominant traffic class with the exception of the behavior with PQ. In this case we find that with PQ we can increase the voice traffic up to 90% above the design point and still meet the delay objective for voice.
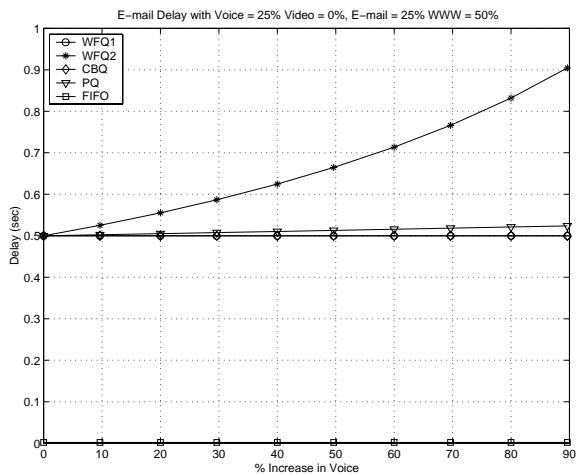


Figure 27: Variation in e-mail Delay with increase in Voice load
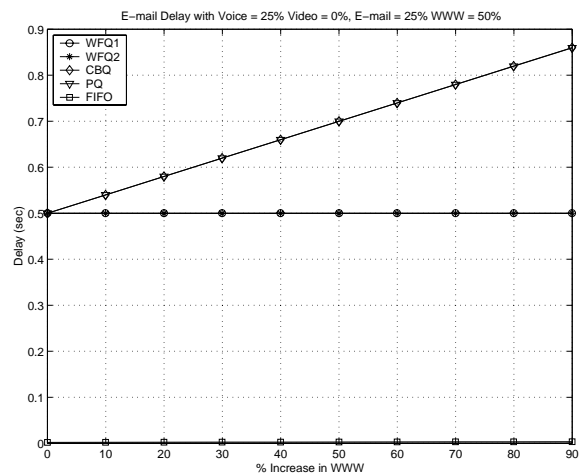


Figure 28: Variation in e-mail Delay with increase in WWW load

For e-mail and WWW we find from Figures 27 to 30 that increasing the voice traffic gives the same results as before except that with WFQ2 the increase in delay is more linear than exponential.
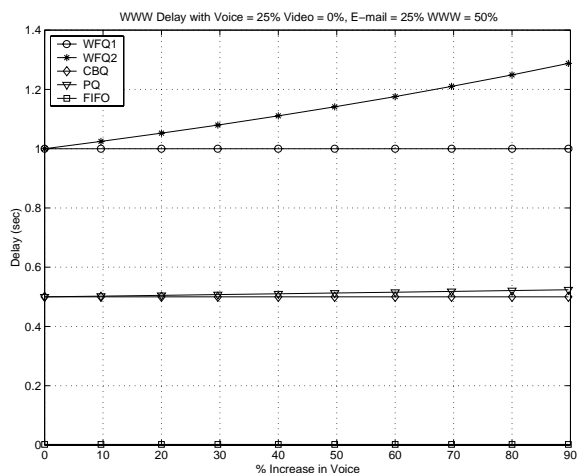


Figure 29: Variation in WWW Delay with increase in Voiceload
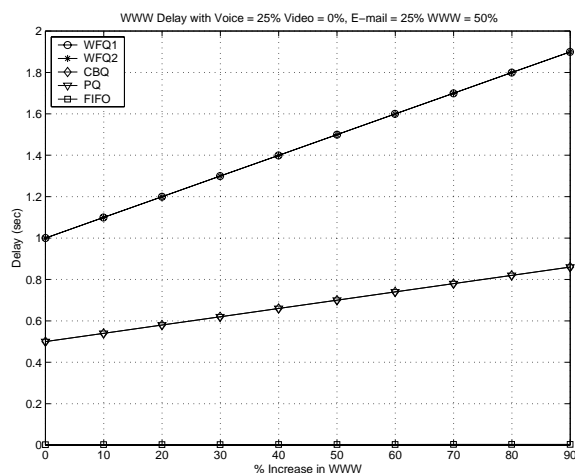


Figure 30: Variation in WWW Delay with increase in WWW load

The picture emerging from these results is that the traffic handling schemes are both sensitive to the type of traffic that dominates the network at the design point as well as to the type of traffic that increases the load on the network. For a network designed with voice as the dominant class, FIFO is the least sensitive to increases in voice traffic and the most sensitive to increases in WWW traffic when considering the delay objectives of voice. PQ is also sensitive to increases in the voice traffic but is able to meet the delay requirements up to a 75% increase in voice load. WFQ and CBQ are both sensitive to increases in the voice load and if the goal is to maintain the delay objectives of voice at all costs, the use of a scheme like WFQ2 can achieve this with a corresponding exponential increase in the delay of e-mail and WWW traffic. Both PQ and CBQ affect e-mail delay performance when WWW traffic is increased thus it is instructive to determine how much of variance in the delay objectives can be tolerated by e-mail traffic. When WWW is the dominant traffic class, using WFQ2 the voice load can be increased up to 90% to obtain the same delay performance for e-mail and WWW as a 50% increase when voice is the dominant traffic. If we assume that e-mail and WWW can tolerate delays up to twice their objective, this result suggests that using a re-allocation of bandwidth as in WFQ2 can allow for an increase in voice traffic of at most 90% in a network designed with WWW as the dominant traffic class. The value of these results are best demonstrated by taking into account the permissible variances in the delay objectives which means using statistical objectives as opposed to deterministic ones and this will be explored in extensions to this research.

Several lessons were learned from this analysis. The most important lesson is that it is possible to quantify the trade-off between network capacity and traffic management. We

also found that the sensitivity of the traffic handling schemes depend on the assumptions made in designing the network as well as the traffic class contributing to the growth in traffic. We note however that the results presented apply to a single node and we anticipate that the effort to extend the results to networks of arbitrary topology will be significant. A second lesson is that there might be a need to review the methodology. In this simple study, we assigned fixed utilization levels to each QoS class which limits the applicability of the results to specific configurations. A more useful approach would be to consider a wider variety of feasible mixes of traffic and base the comparison on this. In this way we would get a better idea of how the traffic mix affects the performance of the different traffic handling mechanisms.

# 4  Future Work

There are several ways in which we propose to apply and extend our analysis in order to fully address the traffic management complexity versus network capacity tradeoff.

 To begin with, we need to extend the analysis to a network whose size is representative of carrier networks and determine incrementally how the analysis scales with increasing network size. This will involve augmenting the analysis with simulation at each stage.

 A second extension is to consider the sensitivity of the four schemes to changes in network traffic. This would be done concurrently with the iterations on network size. Another way in which we will extend the results will be to consider the performance when the delay bounds are statistical and not deterministic. Lastly we will consider the use of stochastic bounds in the traffic models and compare how the performance differs from that of deterministic bounds. Based on the foregoing, we have identified the following tasks:

### 4.0.1  Review and formalization of methodology

The first task will be to review the methodology used in the proof-of-concept which is based on a single-node network. We will need to identify and make adjustments to the methodology to cater for network topologies of arbitrary size and numerous flows. We also intend to address the issue of how to capture the notion of the capacity of a network as opposed to the capacity of a single link. This concept will become an important basis of comparison between the traffic handling schemes in networks of arbitrary topology.

### 4.0.2  Design of network topology

In this task we will design a network topology that is representative of carrier networks. Our concern here will be to capture the size of the networks in terms of the number of nodes as well as in terms of the number of attached hosts and the profile of applications that are

supported. This will enable us to apply our methodology to realistic scenarios.

### 4.0.3 Extension of analysis to carrier network topology

This task will look at the application of Network Calculus to carrier-sized networks of arbitrary topology. In particular we will address the issues of how to determine and use service and arrival curves in order to obtain bounds on performance which will then lead us to quantifying the traffic-handling and network capacity trade-off as well as the sensitivity of the traffic-handling schemes.

### 4.0.4 Implementation of simulation model

We will implement a simulation model of the carrier network using Opnet and run simulations to validate the trends predicted by the analysis.

### 4.0.5 Research on use of stochastic bounds for traffic models

In this task we will examine existing models that use stochastic bounds in the description of the user traffic. We will select candidate models for analysis and may also explore modifications to existing models.

### 4.0.6 Extension of analysis to cover stochastic bounds on traffic models

This task is related to the previous one and will extend the analysis to use probabilistic descriptions of the delay objectives. The objective is to determine to what extent the use of stochastic bounds on the traffic models affects the difference in capacity requirements of the traffic handling mechanisms.

### 4.0.7 Research on use of statistical bounds for performance objective

Similar to Section 4.0.5, we will select candidate models for analysis which use statistical descriptions on performance.

### 4.0.8 Extension of analysis to cover statistical bounds on performance objectives

In this task we will extend the analysis to the case of statistical bounds on delay objectives. We will then combine this with the analysis from Section 4.0.6 to derive an analysis that includes statistical bounds on both the traffic models and the performance objectives.

### 4.0.9 Extension of analysis to variants of Class-Based Queueing and Weighted Fair Queueing

There are many variations on how Class-Based Queueing and Weighted Fair Queueing can be implemented. This task will consider a subset of these variations and determine whether the capacity requirements are linked to the way in which the mechanisms are implemented.

### 4.0.10 Review of research objectives and results

This last task will review the objectives and results of the research and identify any open issues for future work.

## 5 Significance to Sprint

This work is significant to Sprint in three main ways. The first is that it addresses an important question in network engineering and design: that of identifying the tradeoffs associated with the use of traffic handling mechanisms with respect to network capacity. Network engineers are faced with a multitude of options when it comes to deciding what traffic handling mechanisms to use and how much capacity to provision in the network. In most cases decisions are reached in an ad hoc manner either by trial-and-error or desired performance is obtained by over-provisioning. With the results that we expect to obtain, network engineers will be able to obtain an understanding of the tradeoffs and base their decisions on quantitative data. We also note that by incorporating sensitivity analysis we provide a tool for long-term planning since we are able to show how the traffic handling mechanisms will react to growth in network traffic.

The second benefit to Sprint which is related to the first is in the development of a methodology which can be used to compare traffic handling schemes in general. In most cases comparisons between traffic handling mechanisms are based on regions of schedulability which are simply representations of the amount of traffic of each supported QoS class that can be admitted into the network such that their QoS is satisfied for a given link bandwidth. The regions of schedulability are usually represented in graphical form and this limits their use to networks having 2 classes of service. With the methodology that we are proposing, one can get an idea of the relative difference in performance between the traffic handling mechanisms easily for any number of QoS classes.

Lastly, we expect this work to prove valuable in the design of Sprint's Edge-Core Network Architecture. In general, networks have a sturctured hierarchy comprising of the access layer, the distribution layer and the core layer. The access layer is the outermost part of the network to which customers are directly connected. The distribution layer handles aggregation of traffic from multiple access points and provides transit between the core and access. The core is the innermost part of the network responsible for high-speed transfer of

customer traffic. Given this hierarchy, network providers have several options in deciding how to implement traffic handling mechansisms at the different layers. One of the models that is emerging is one in which the level of aggregation increases towards the core of the network and the edges (access and distribution) use per-flow and/or per-class traffic handling. A critical issue that will need to be addressed in such networks is how to allocate the end-to-end delay between the different devices at each layer of the hierarchy in order to meet the delay requirements of network users. The way in which these delays are allocated will be directly related to the available bandwidth in the different layers as well as the traffic handling mechanisms that are used. The methodology and analysis that we are proposing will allow Sprint to evaluate different edge-core architectures and obtain a solution that will meet their customers' requirements.

# References

[1] Black D., *Building Switched Networks*, Addison-Wesley-Longman Inc.,1999.

[2] Blake S. et al., *A Framework for Differentiated Services*, Oct. 1998.

[3] Blake S. et al., *An Architecture for Differentiated Services*, Oct. 1998.

[4] Braden R., D. Clark, S. Shenker, *Integrated Services in the Internet Architecture: an Overview* IETF RFC 1633, June 1994.

[5] Braden R. et al, *RFC 2205: Resource Reservation Protocol (RSVP) - v1. Functional specification*, Sept. 1997.

[6] Breslau L., S. Shenker, *Best Effort versus Reservations: A Simple Comparative Analyis*, Proc. ACM SIGOCMM 1998, p.3-16.

[7] Cruz R.L., *A Calculus for Network Delay, Part I: Network Elements in Isolation*, IEEE Transactions on Information Theory, Vol 37, no. 1, Jan, 1991.

[8] Cruz R.L., *A Calculus for Network Delay, Part II : Network Analysis*, IEEE Transactions on Information Theory, Vol 37, no. 1, Jan, 1991.

[9] Cruz R.L., *Quality of Service Guarantees in Virtual Switched Networks*, IEEE JSAC, vol 13, no. 6, Aug. 1995.

[10] Ferguson P.,G. Huston, *Quality of Service, delivering QoS on the Internet and in Corporate Networks*, J. Wiley & Sons Inc, 1998.

[11] Ferrari D., *Client Requirements for Real-time Communication Services*, IEEE Communications Magazine Nov. 1990 p. 65-72.

[12] IETF RFC 2211, *Specification of the Controlled-Load Network Element Service*, ftp://ftp.isi.edu/in-notes/rfc2211.txt.

[13] IETF RFC 2212, *Specification of Guaranteed Quality of Service*, ftp://ftp.isi.edu/in-notes/rfc2212.txt.

[14] Kilkki K., *Differentiated Services for the Internet*, Macmillan Technical Publishing, 1999.

[15] Microsoft, *Quality of Service Technical White Paper*,

[16] Parekh A., R. Gallager, *A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: The Single Node Case*, IEEE/ACM Transactions on Networking, Vol. 1, no. 3, June 1993.

[17] Parekh A., R. Gallager, *A Generalized Processor Sharing Approach to Flow Control in Integrated Services Networks: The Multiple Node Case, IEEE/ACM Transactions on Networking*, Vol. 2, no. 2, April 1994.

[18] Stardust Forums, *Internet Bandwidth Management Whitepaper*, Proc. iBAND2, May 1999.